



اقتصادسنجی
جلد اول
تک معادلات با فروض کلاسیک
جزء اول

دکتر مسعود درخشان



اقتصادسنجی

جلد اول

تک معادلات با فروض کلاسیک

جزء اول

دکتر مسعود درخشان

تهران

۱۳۷۴



سازمان مطالعه و تدوین کتب علوم انسانی دانشگاهها (سمت)

درخشان، مسعود

اقتصادسنجی / مسعود درخشان. — تهران: سازمان مطالعه و تدوین کتب علوم انسانی دانشگاهها (سمت)، ۱۳۷۴.

ج: جدول، نمودار. — (سازمان مطالعه و ...؛ ۱۴۲ و ۱۴۳: اقتصاد؛ ۹ و ۱۰)

عنوان پشت جلد به انگلیسی: *Econometrics; vol. 1, part 1: Single*

Equations with Classical Assumptions.

جلد اول شامل ۲ جزء: تک معادلات با فرض کلاسیک.

واژه‌نامه.

کتابنامه.

۱. اقتصادسنجی. ۲. اقتصاد ریاضی. الف. فروست. ب. عنوان.

الف ۱۳۹ د ۴ HB

الف ۳۳۷ د ۳۲۷ ۱۵۱۹۵ / ۰۳۳۰



اقتصادسنجی: جلد اول؛ تک معادلات با فرض کلاسیک (جزء اول)

دکتر مسعود درخشان

سازمان مطالعه و تدوین کتب علوم انسانی دانشگاهها (سمت)

چاپ اول: پاییز ۱۳۷۴

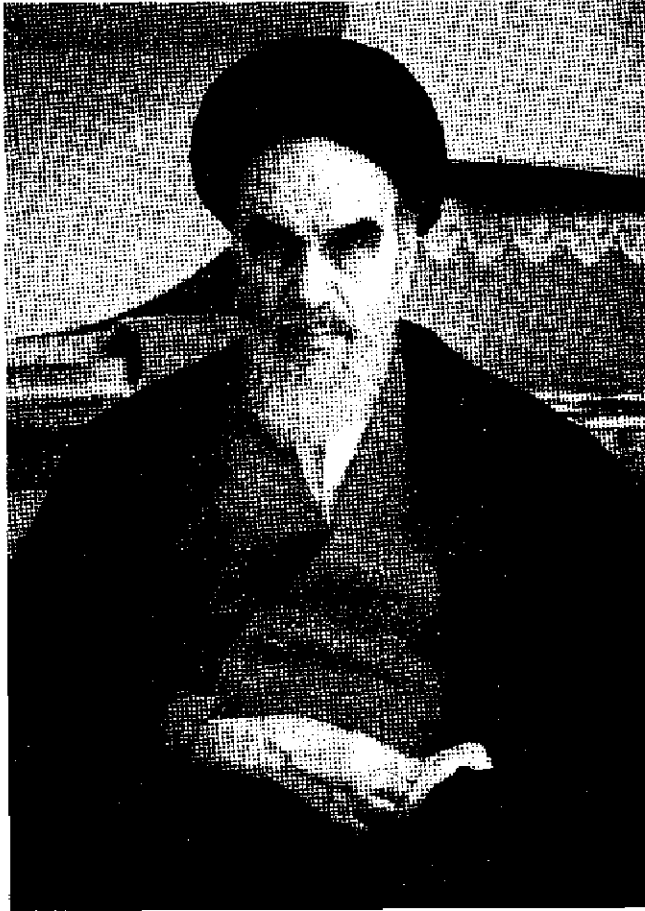
تعداد: ۵۰۰۰

حروفچینی، صفحه‌آرایی و لیتوگرافی: سمت

چاپ: مهر (قم)

کلیه حقوق اعم از چاپ و تکثیر، نسخه برداری، ترجمه، اقتباس، کلید راهنما و غیره برای «سمت» محفوظ است (نقل مطالب با ذکر مأخذ بلامانع است).

ابن‌الحسن



بدون شک جهانخواران به همان میزان که از شهادت طلبی و سایر ارزشهای ایثارگرانه ملت ما واهمه دارند، از گرایش و رواج اقتصاد اسلام به طرف حمایت از پابرهنگان در هراسند و مسلم هر قدر کشور ما به طرف فقرزدایی و دفاع از محرومان حرکت کند، امید جهانخواران از ما منقطع و گرایش ملتهای جهان به اسلام زیادتر می‌شود.

سخن «سنت»

یکی از اهداف مهم انقلاب فرهنگی، ایجاد دگرگونی اساسی در دروس علوم انسانی دانشگاهها بوده است و این امر، مستلزم بازنگری منابع درسی موجود و تدوین منابع مبنایی و علمی معتبر و مستند با در نظر گرفتن دیدگاه اسلامی در مبنایی و مسائل این علوم است. ستاد انقلاب فرهنگی در این زمینه گامهایی برداشته بود، اما اهمیت موضوع اقتضا می کرد که سازمانی مخصوص این کار تأسیس شود و شورای عالی انقلاب فرهنگی در تاریخ ۶۳/۱۲/۷ تأسیس «سازمان مطالعه و تدوین کتب علوم انسانی دانشگاهها» را که به اختصار «سنت» نامیده می شود، تصویب کرد.

بنابراین، هدف سازمان این است که با استمداد از عنایت خداوند و همت و همکاری دانشمندان و استادان متعهد و دلسوز، به مطالعات و تحقیقات لازم بپردازد و در هر کدام از رشته های علوم انسانی به تألیف و ترجمه منابع درسی اصلی، فرعی و جنبی اقدام کند. دشواری چنین کاری بر دانشمندان و صاحب نظران پوشیده نیست و به همین جهت مرحله کمال مطلوب آن، باید بتدریج و پس از انتقادهای و یادآوریهایی پایایی ارباب نظر به دست آید و انتظار دارد که این بزرگواران از این همکاری دریغ نوزند.

کتاب حاضر که دربرگیرنده مفاهیم اقتصادسنجی است برای دانشجویان رشته های اقتصاد در مقطع کارشناسی و کارشناسی ارشد به عنوان منبع اصلی درس «اقتصادسنجی» به ارزش ۷ واحد تدوین شده است و برای رشته های مدیریت، آمار و مهندسی صنایع نیز قابل استفاده است. امید است علاوه بر جامعه دانشگاهی سایر علاقه مندان نیز بتوانند از آن استفاده کنند.

از استادان و صاحب نظران ارجمند تقاضا می شود با همکاری، راهنمایی و پیشنهادهای اصلاحی خود، این سازمان را در جهت اصلاح کتاب حاضر و تدوین دیگر آثار مورد نیاز جامعه دانشگاهی جمهوری اسلامی ایران یاری دهند.

فهرست مطالب

صفحه	عنوان
	پیشگفتار
۱	فصل اول: مفاهیم و تخمین مدل رگرسیون خطی ساده
۱	۱-۱ مقدمه
۱۶	۱-۲ خصوصیات آماری مدل رگرسیون خطی ساده
۲۸	۱-۳ تخمین مدل رگرسیون خطی: روش حداقل مربعات معمولی (OLS)
۳۹	۱-۴ ضریب تعیین
۵۰	۱-۵ خطای معیار تخمین
۵۶	مسائل فصل اول
۶۵	حل مسائل فصل اول
	فصل دوم: آزمونهای آماری خصوصیات مطلوب تخمین زنده‌ها در مدل رگرسیون خطی ساده
۸۵	۲-۱ مقدمه
۸۵	۲-۲ خصوصیات آماری $\hat{\alpha}$ و $\hat{\beta}$
۸۶	۲-۳ آزمون فرضیه برای هر یک از پارامترها
۱۱۰	۲-۴ آنالیز واریانس و آزمون معنی دار بودن مدل رگرسیون
۱۳۹	۲-۵ خصوصیات مطلوب تخمین زنده‌ها
۱۴۸	۲-۶ خصوصیات مطلوب تخمین زنده‌های حداقل مربعات معمولی: قضیه گاس-مارکف
۱۸۲	مسائل فصل دوم
۱۸۸	حل مسائل فصل دوم
۲۰۸	فصل سوم: پیش‌بینی و مباحث تکمیلی در مدل رگرسیون خطی ساده
۲۰۸	۳-۱ مقدمه
۲۰۹	۳-۲ پیش‌بینی در مدل‌های رگرسیون ساده
۲۲۷	۳-۳ تغییر مقیاس در متغیرها
۲۳۷	۳-۴ رگرسیون معکوس
۲۴۴	۳-۵ مشاهدات دورافتاده
۲۵۰	۳-۶ رابطه‌های غیرخطی و تبدیل متغیرها
۲۶۹	مسائل فصل سوم
۲۷۳	حل مسائل فصل سوم
۲۸۹	فصل چهارم: مدل رگرسیون خطی با دو متغیر توضیحی
۲۸۹	۴-۱ مقدمه

۲۹۱	۴-۲ تخمین مدل
۲۹۸	۴-۳ ضریب تعیین
۳۰۵	۴-۴ آزمون فرضیه
۳۱۳	۴-۵ پیش‌بینی
۳۲۰	۴-۶ ضرایب همبستگی و ضرایب تعیین
۳۳۹	مسائل فصل چهارم
۳۴۴	حل مسائل فصل چهارم
۳۵۹	فصل پنجم: مفاهیم و تخمین مدل رگرسیون خطی چند متغیره
۳۵۹	۵-۱ مقدمه
۳۶۱	۵-۲ بیان ماتریسی رگرسیون چند متغیره
۳۷۰	۵-۳ تخمین مدل رگرسیون خطی چند متغیره: روش حداقل مربعات معمولی
۳۸۷	۵-۴ ضریب تعیین
۳۹۷	۵-۵ ضریب تعیین تعدیل شده یا \bar{R}^2
۴۰۴	۵-۶ نکاتی پیرامون R^2 و \bar{R}^2
۴۱۳	۵-۷ بیان ماتریسی R^2 بر حسب مشاهدات اصلی
۴۲۴	مسائل فصل پنجم
۴۲۸	حل مسائل فصل پنجم
۴۴۱	فصل ششم: آزمون پارامترهای مدل رگرسیون خطی چند متغیره
۴۴۱	۶-۱ مقدمه
۴۴۳	۶-۲ خصوصیات آماری $\hat{\beta}$
۴۵۹	۶-۳ آزمون هر یک از پارامترها
۴۷۱	۶-۴ آنالیز واریانس و آزمون معنی‌دار بودن مدل رگرسیون
۴۸۷	۶-۵ آزمون یک ترکیب خطی از پارامترها
۴۹۲	۶-۶ آزمون همزمان پارامترها و تعیین نواحی اطمینان
۴۹۸	۶-۷ آزمون تساوی چند پارامتر یا یکدیگر
۵۰۴	۶-۸ آزمون همزمان چند ترکیب خطی از پارامترها: آماره عمومی آزمون و کاربرد آن
۵۱۸	مسائل فصل ششم
۵۲۳	حل مسائل فصل ششم

تقدیم به:

شهدای دانشکده‌های اقتصاد

و آنان که

«محرومیت‌زدایی» را «عقیده و راه و رسم زندگی خود می‌دانند».

پیشگفتار

امیدوارم مطالب این کتاب برای دانشجویان عزیز و سایر پژوهشگران اقتصادی سودمند بوده، بتواند زمینه مطالعات دقیقتر و جامعتری را در مباحث اقتصادی فراهم آورد. کتاب حاضر اولین جلد از یک مجموعه است که برای سطوح کارشناسی و بالاتر در رشته اقتصادسنجی تدوین شده است. در این پیشگفتار به چند نکته اشاره می‌کنیم.

۱. ضرورت

اولین سؤال این است که آیا ضرورتی برای انتشار این کتاب وجود دارد؟ این سؤال بویژه با توجه به کتابهای متعددی که در سالهای اخیر در رشته اقتصادسنجی به زبان فارسی به چاپ رسیده است^۱ اهمیت فراوان دارد. اکثر کتابهای موجود، ترجمه‌های بسیار خوب از منابع ارزشمندی است که اغلب از متون درسی دانشگاههای معتبر جهان به‌شمار می‌آیند. تدریس یکی از این کتابها که متناسب با دوره کارشناسی است، و کتابی دیگر که هماهنگی بیشتری با دوره کارشناسی ارشد دارد، می‌تواند برای دانشجویان عزیز و اساتید بزرگوار تضمینی باشد که معیار آموزشی، حداقل در سطح دانشگاههای معتبر جهان حفظ شده است. این تضمین، دارای اهمیت فراوان بوده و آثار مثبت زیادی دارد. آنچه در حال حاضر مورد نیاز نظام آموزش عالی دانشگاههای کشور است اعتماد به نفس و تکیه بر تواناییهای علمی داخل برای استفاده از امکانات خارجی است؛ زیرا ظرفیت علمی دانشجویان عزیز و استادان محترم قطعاً از موارد مشابه در

۱. برای ملاحظه فهرستی از این کتابها، به پیوست «ه» مراجعه کنید.

دانشگاههای خارج کمتر نیست.

روش فوق می‌تواند نقطه ضعفی نیز داشته باشد، که دقیقاً ناشی از ماهیت آن است. آن دسته از تألیفات جدید اقتصادسنجی به زبان خارجی، که از کتابهای درسی دانشگاههای معتبر جهان در دوره‌های کارشناسی و کارشناسی ارشد به شمار می‌آیند، اولاً از سطح ریاضیات دانشجویان و ثانیاً از کیفیت تقاضا در بازار اشتغال برای فارغ‌التحصیلان متأثر است. متأسفانه سطح ریاضیات دانشجویان اقتصاد در کشورهای خارج کاهش یافته و بنابراین اکثر نویسندگان خارجی برای اینکه بتوانند کتابهای «موفق‌تری» عرضه کنند از عمق و گستره ریاضیات در تألیفات خود کاسته‌اند. همچنین نیاز مبرم اقتصاد کشورهای صنعتی به فارغ‌التحصیلان آشنا به استفاده از برنامه‌های کامپیوتری در حل مسائل اقتصادسنجی، موجب شده است که تحلیلهای نظری و دقتهای مبنایی در مطالعات اقتصادسنجی به نحو قابل ملاحظه‌ای محدودتر شود. در مواردی نیز که از ریاضیات در سطوح خوب و عالی برای ارائه مباحث نظری اقتصادسنجی استفاده شده است مطالب بقدری تخصصی ارائه گردیده که فقط قابل استفاده برای متخصصان این رشته است و دانشجویان دوره‌های کارشناسی و کارشناسی ارشد اقتصاد، از آن چندان بهره‌ای نمی‌برند. با وجود این، امتیاز بعضی از کتابهای جدید، یعنی آشنا کردن دانشجویان با مثالهای کاربردی و برنامه‌های کامپیوتری نباید فراموش شود؛ هرچند که به قیمت کاهش دقتهای نظری صورت گرفته است.

کتاب حاضر که با توجه به سطح خوب و ظرفیت عالی ریاضیات دانشجویان عزیز کشور و علاقه آنها به دقت در مباحث نظری تدوین شده است می‌تواند به عنوان مکملی برای اکثر کتابهای موجود محسوب شود. معمولاً نمی‌توان یک کتاب خارجی در سطوح کارشناسی یا کارشناسی ارشد یافت که از هر جهت نسبت به سایر کتابهای موجود برتری داشته باشد؛ اما مباحث یا مثالهایی در هر کتاب وجود دارد که جامع‌تر و عمیق‌تر از موارد مشابه در سایر کتابها است. برای تدوین هر بحث در این کتاب، اکثر منابع اصلی و معتبر ملاحظه شده و سعی بر این بوده است که از نکات مثبت موجود در هر یک از آنها استفاده شود. بنابراین شاید بتوان گفت که اگر ترجمه یکی از منابع خوب اقتصادسنجی مینا

قرار گیرد کتاب حاضر می‌تواند مکملی برای آن باشد. بر فرض درستی این نکته، ضرورت انتشار این کتاب اجمالاً به ثبوت می‌رسد.

۲. هدف

در تمام فصلهای این کتاب سعی بر این بوده است که مباحث اقتصادسنجی، در سطوح پایه و متوسطه، در کمال سادگی بیان شود به گونه‌ای که بتواند به قول دانشجویان عزیز یک «خودآموز اقتصادسنجی» باشد. مباحث پیشرفته اقتصادسنجی کتاب مستقلی است که ان‌شاءالله در آینده نزدیک منتشر خواهد شد.

بیان مباحث به زبانی بسیار ساده، هدفی است که آگاهانه در این کتاب تعقیب شده است. هرچند که ساده‌نویسی می‌تواند در ابعادی با اهداف پژوهشی و نوآوریهای علمی منافات داشته باشد اما به نظر می‌رسد که در مقطع فعلی، نمی‌توان آموزش و پژوهش برای نوآوری در اقتصادسنجی را به عنوان یک هدف اصلی در صدر برنامه‌های اقتصاد دانشگاهها قرار داد. به دانشجویان توصیه می‌شود که با توجه به شرایط و مقتضیات اقتصادی کشور، مطالعات اقتصادسنجی را تنها به عنوان «وسیله‌ای» بسیار خوب و قوی برای درک و تحلیل سیاستها و مدل‌های اقتصادی در کشورهای در حال توسعه، و نیز برای شناخت مکانیسم تحولات و سیاست‌گذاریهای اقتصادی کشورهای صنعتی، با دقت کامل فراگیرند.

اگر هدف، مطلق کردن مطالعات اقتصادسنجی و رسیدن به مرحله نوآوری در این قلمرو می‌بود، آنگاه به طور قطع نحوه نگارش اینجانب در این کتاب تغییر اساسی می‌کرد و از ساده‌نویسی و تبیین همه نکات و جوانب قضایا و حل مسائل، به شدت اجتناب می‌شد. بنابراین تأکید بر این است که دانشجویان در حداقل زمان و با فراگیری حداکثر مطالب، هر چه سریعتر با ابعاد مختلف روشهای اقتصادسنجی آشنا شوند تا به مرحله اصلی فراگیری و تحقیقات اقتصادی خود برسند که همانا تجزیه و تحلیل استراتژیهای رشد و توسعه و بررسی ساختار سیاستهای مختلف اقتصادی در کشورهای در حال توسعه و آثار بلندمدت آن در اقتصاد جهانی است.

۳. الگوی تدوین

کیفیت طبقه‌بندی و ارائه مباحث تابعی از ضرورت و هدف تدوین این کتاب بوده است. در هر فصل، مطالب در دو سطح کارشناسی و کارشناسی ارشد ارائه می‌شود. آن قسمت از مباحثی که به دوره‌های بالاتر از کارشناسی مربوط می‌شود با علامت (e) مشخص شده است. همین امر زمینه مناسبی را ایجاد می‌کند که دانشجویان دوره کارشناسی که علاقه یا فرصت بیشتری برای مطالعه دارند بتوانند قلمرو مطالعات خود را در هر مبحث افزایش دهند.

مطالب این کتاب به گونه‌ای تنظیم شده است که سطح فعلی ریاضیات و آمار عمومی که در حال حاضر در رشته‌های اقتصاد و بازرگانی در دانشگاهها تدریس می‌شود به عنوان پیش‌نیاز کفایت می‌کند. با وجود این، پیوست «الف» نکات ریاضی که در انتهای کتاب قرار دارد باید دقیقاً مطالعه شود. برای آن دسته از دانشجویان دوره کارشناسی اقتصاد که به هر دلیل ترجیح می‌دهند اقتصادسنجی را بدون استفاده از جبر ماتریسی فراگیرند مشکل خاصی وجود نخواهد داشت، زیرا تنظیم مطالب به نحوی است که می‌توان مباحث مطرح شده به زبان ماتریسی را از سایر عناوین جدا کرد به گونه‌ای که بدون آشنایی با جبر ماتریسی بتوان مفاهیم و روشهای اصلی اقتصادسنجی را فراگرفت. با وجود این، تحلیل مباحث اقتصادسنجی به کمک جبر ماتریسی، برای دانشجویان دوره‌های کارشناسی ارشد ضروری است.

با توجه به اینکه کتاب حاضر در واقع مکملی برای سایر کتابهاست، نه تنها سعی شده است که مطالب در چهارچوب مثالهای ساده مطرح شود، بلکه جمعاً ۸۳ مثال جامع در خلال ده فصل، نقش عمده‌ای در تبیین موضوعات داشته باشد. همچنین جمعاً ۱۱۳ مسأله در پایان فصلها آمده که همگی با روشنی و سادگی تمام حل شده است. از نظر کثرت مثال و مسأله و نبودن ابهام در راه‌حلیها و تبیین موضوعات تا سطوح پیشرفته در هر مبحث، شاید کتاب حاضر را بتوان به عنوان اولین کوشش در این زمینه معرفی کرد. با مطالعه این کتاب، دانشجویان عزیز می‌توانند از طریق تسلط بیشتر در درک و حل مسائل، تمرینات موجود در کتابهای دیگر را که استادان بزرگوار توصیه می‌نمایند

بهرتر حل و بحث کنند.

چند نکته تاریخی را در بخش پیوستها مطرح کرده‌ایم. مباحث تاریخی همواره شیرین و بسیار آموزنده است. در خلال بررسیهای تاریخی در زمینه روش حداقل مربعات، موفق شدم که در کتابخانه قدیمی دانشگاه آکسفورد، اولین کتابنامه در روش حداقل مربعات را به دست آورم که در سال ۱۸۷۷ یعنی همزمان با دوران سلطنت ناصرالدین شاه قاجار منتشر شده است. در این کتابنامه، جمعاً ۴۰۸ مقاله و کتاب گردآوری شده که درباره این روش از سال ۱۷۲۲ (مصادف با شکست سلطان حسین صفوی از محمود افغانی) تا سال ۱۸۷۶ (سه سال بعد از اولین سفر ناصرالدین شاه به فرنگستان) به چاپ رسیده است. همچنین موفق شدم اولین کتاب معروف درسی در روش حداقل مربعات را که در سال ۱۸۷۷ به چاپ رسیده و در سال ۱۸۸۵ با تجدید نظر کلی مجدداً منتشر شده است ملاحظه نموده و صفحاتی از آن را به همراه دو جدول آماری، در کتاب حاضر منعکس کنم تا بدین ترتیب دانشجویان عزیز با سطح مباحث در آن زمان بیشتر آشنا شوند.

در پیوست نکات تاریخی، سعی کرده‌ام که قدر و منزلت «گاست»^۱ آماردان معروف قرن بیستم و صاحب توزیع t را برای دانشجویان رشته اقتصاد، بیشتر روشن کنم. متأسفانه تنها نکته‌ای که در بعضی از تألیفات آماری درباره این آماردان بزرگ مطرح می‌شود اشتغال نامبرده در یک کارخانه آبجوسازی است. همین امر موجب شده است که این شخصیت بزرگ علمی، که توانست برای اولین بار توزیع و آزمون t را در سال ۱۹۰۸، یعنی همزمان با انقلاب مشروطیت و حکومت محمدعلی شاه قاجار، کشف کند مورد کم‌مهری و بی‌توجهی در بعضی از محافل علمی ایران قرار گیرد.

تدوین این کتاب، مبتنی بر متون اصلی اقتصادسنجی دانشگاه‌های معروف دنیا است. از تمام کتابها و مقالاتی که در قسمت منابع و مأخذ ذکر گردیده مستقیماً بهره‌برداری شده است. در مواردی نیز اثبات بعضی از قضایا به روشهای ساده‌تر، ارائه

1. William Sealy Gosset

نظریات یا اشکالاتی در تجزیه و تحلیل مباحث، و یا طراحی و حل بعضی از مثالها و مسائل آخر هر فصل از اینجانب است. با وجود این، منابع زیر بیشترین سهم را در تدوین این کتاب داشته است: جانستون^۱، کمنتا^۲، مادالا^۳، استیوارت و والیس^۴، کندی^۵، بیکن^۶، هاروی^۷، و گرین^۸.

یافتن واژه‌های فارسی برای اصطلاحات علمی یکی از سخت‌ترین مراحل تدوین کتاب و مقاله علمی به زبان فارسی است. متأسفانه کوشش جدی و مستمری از طرف جامعه متخصصان زبان فارسی در حل این مشکل ملاحظه نمی‌شود. رشد سریع علوم و تکنولوژی و تداوم ظهور اصطلاحات علمی جدید، و همچنین علاقه فراوان دانشجویان و استادان ارجمند به ترجمه آخرین دستاوردهای علمی جهان، زبان فارسی را با کمبود شدید واژه‌ها و اصطلاحات علمی مواجه کرده است. این امر دلالت بر فوریت اتخاذ تصمیماتی دارد که برای حفظ زبان فارسی و رشد و توسعه آن ضروری است.

در مقطع کنونی، که رابطه نزدیکی بین جامعه اقتصاددانان کشور و متخصصان ادبیات فارسی وجود ندارد، شاید روش بهینه، همفکری و تبادل نظر استادان و پژوهشگران اقتصاد با یکدیگر و کوشش در استفاده از ترجمه‌های یکسان برای اصطلاحات اقتصادی است. به همین دلیل سعی شده است که در این کتاب، تا حد امکان، از واژه‌هایی استفاده شود که استادان محترم در تألیفات و ترجمه‌های خود برگزیده‌اند. بنابراین در کتاب حاضر بسیاری از اصطلاحاتی را که در خلال چندین سال تدریس اقتصادسنجی از آنها استفاده کرده‌ام کنار گذاشته و به جای آنها، اصطلاحات موجود در کتابهای منتشر شده را قرار داده‌ام. با وجود این، نظر به ناهماهنگیهای رایج در ترجمه‌ها، طبعاً فقط به فصل مشترکات اکتفا شده است. در مواردی نیز انتقال مفاهیم در قالب جمله‌ها، به ساختن و پرداختن کلمات یا اصطلاحات جدید ترجیح داده شده تا بدین وسیله زمینه برای طراحی واژه‌های زیبا و صحیح فارسی توسط اهل فن از بین نرود.

1. J. Johnston (1984)

2. Jan Kementa (1986)

3. G. S. Maddala (1992)

4. M. B. Stewart and K. F. Wallis (1981)

5. P. Kennedy (1992)

6. R. Bacon (1988)

7. A. C. Harvey (1990)

8. W. H. Greene (1993)

۴. تاریخچه تدوین

ضرورت تدوین یک کتاب اقتصادسنجی به زبان فارسی، در سالهای ۱۳۴۶ - ۱۳۵۰، که اینجانب دانشجوی دانشکده اقتصاد دانشگاه تهران بودم، کاملاً احساس می‌شد، زیرا هیچ کتاب درسی در این زمینه وجود نداشت. تدوین مقدمات کتاب حاضر عملاً در سالهای ۱۳۵۱ و ۱۳۵۲ و همزمان با تحصیلات دوره فوق لیسانس نگارنده در مدرسه اقتصاد لندن (LSE) آغاز شد. توجه خاص آن مدرسه به آموزش و پژوهش در اقتصادسنجی، انگیزه بسیار مناسبی در تدوین این کتاب بود. در دوره تحصیلات دکتری و تهیه رساله در آکسفورد به نظر نگارنده رسید که تألیف کتابی در زمینه کاربرد روشهای سیستم و کنترل بهینه در اقتصادسنجی سریعتر به نتیجه خواهد رسید. بدین ترتیب اتمام کتاب اقتصادسنجی به تعویق افتاد.

حضور در ایران در سال ۱۳۵۵ به مدت یک ترم تحصیلی برای تکمیل یک مدل اقتصادسنجی کلان، فرصت مناسبی را فراهم آورد تا جزوه‌های تدوین شده در اقتصادسنجی و در کاربرد نظریه سیستم و کنترل در مدل‌های اقتصادی را در دانشکده اقتصاد دانشگاه تهران تدریس کنم. علاقه و استقبال بسیار زیاد دانشجویان موجب شد که بعد از بازگشت به آکسفورد، کوشش برای تکمیل و انتشار هر چه سریعتر کتاب اقتصادسنجی آغاز شود. خوشبختانه وقوع انقلاب اسلامی در سال ۱۳۵۷ باعث شد که از «قیل و قال» اقتصادسنجی و نظریه‌های کنترل رها شده و به «شور و حال» اقتصاد اسلامی، تاریخ تحولات اندیشه‌ها و نظامهای اقتصادی، که سالها مورد علاقه‌ام بود، بازگردم.

بعد از انقلاب فرهنگی و بازگشایی دانشگاهها، افتخار همکاری در تدریس اقتصادسنجی در دوره‌های کارشناسی، کارشناسی ارشد، و دکتری دانشگاههای تهران، تربیت مدرس، و امام صادق (ع) بر عهده اینجانب قرار گرفت، و در نتیجه جزوه‌های درسی سابق به مراتب کاملتر شد. اما در خلال این دوره، همکاران دانشمند و متخصصان بزرگوار اقتصادسنجی در دانشگاههای تهران، شهید بهشتی، اصفهان، مازندران، شیراز، و بعضی دیگر از مراکز مطالعات اقتصادی، کتابهای بسیار خوب و معتبری را در این زمینه

ترجمه یا تألیف کردند، در نتیجه دیگر هیچ دلیلی بر ضرورت انتشار آن جزوه‌ها، به عنوان کتابی مستقل وجود نداشت. به هر حال با توجه به شرحی که در بند اول آمد و با اصرار دانشجویان عزیزم، سرانجام انتشار این کتاب، بعد از تجدیدنظر اساسی و بهره‌برداری از جدیدترین متون، به تحقق پیوست، امیدوارم که مفید باشد.

۵. تشکر و قدردانی

از دانشجویان عزیزم که در خلال سالها تدریس اقتصادسنجی همواره در تدوین این کتاب مشوقم بوده‌اند صمیمانه تشکر می‌کنم. خاطرات شیرین کلاسهای درس، صفای قلب، اراده راسخ، و آمادگی کامل این جوانان عزیز برای حرکت در جهت رشد و توسعه اقتصادی کشور، انگیزه اصلی در تداوم خدمات ناچیز دانشگاهی نگارنده بوده است. از درگاه خداوند متعال موفقیت هرچه بیشتر برای این جوانان عزیز مسئلت می‌نمایم؛ جواناتی که با داشتن حداقل امکانات مادی می‌کوشند تا با اعتلای سطح علمی خود، دانشگاههای کشور را به مرحله رقابت با بهترین دانشگاههای معتبر جهان برسانند.

لازم است از حمایت‌های جناب آقای دکتر احمدی، ریاست محترم سازمان مطالعه و تدوین کتب علوم انسانی دانشگاهها (سمت)، در انتشار کتاب حاضر، صمیمانه سپاسگزاری کنم. همه می‌دانیم که سازماندهی و مدیریت نهادهای علمی و پژوهشی در سالهای اخیر، با توجه به کمبود منابع انسانی و مالی، تا چه حد طاقت‌فرسا است. تأسیس و توسعه نهاد علمی «سمت» مدیون فداکاریهای ایشان است؛ و اینجانب به سهم خود از این همه تلاش و کوشش صمیمانه تشکر می‌کنم. از جناب آقای عبدالرضا حسنی که با دقت، حوصله و تسلط فراوان متن کتاب را ویراستاری نمودند و از سرکار خانم دهقان‌نیری، دبیر محترم گروه اقتصاد، صمیمانه سپاسگزاری می‌کنم. از زحمات اداره تولید «سمت» که مسئولیت حروف چینی، نمونه‌خوانی، لیتوگرافی و چاپ کتاب را بر عهده داشته است قدردانی می‌شود.

مسعود درخشان

مرداد ۱۳۷۳

مفاهیم و تخمین مدل رگرسیون خطی ساده

۱-۱ مقدمه

اقتصادسنجی^۱ معمولاً با بحث مدل‌های رگرسیون خطی ساده آغاز می‌شود. در اینجا تعریف واژه‌های اقتصادسنجی، مدل رگرسیون، و مدل رگرسیون خطی ساده مفید است. بدیهی است این تعاریف در حد اجمال بوده و صرفاً برای ایجاد زمینه مناسب برای تبیین مفاهیم کلی اقتصادسنجی به کار می‌رود.

۱. اقتصادسنجی

در کتابهای اقتصادسنجی معمولاً فصل اول به تعریفی مبسوط از اقتصادسنجی و ماهیت آن اختصاص یافته است. با توجه به تخصصی بودن این تعاریف، به نظر می‌رسد در این کتاب بهتر است تعریف و ماهیت اقتصادسنجی بعد از مطالعه دقیق روشهای اقتصادسنجی (فصل آخر جلد دوم) ارائه شود.^۲ در این قسمت به تعریفی اجمالی از اقتصادسنجی اکتفا می‌کنیم.

اقتصادسنجی در لغت به معنای علم اندازه‌گیری در اقتصاد است. قلمرو این تعریف بسیار وسیع است؛ برای مثال، موضوع حسابداری ملی نیز به طور عمده اندازه‌گیری کمیّات اقتصادی مانند تولید ناخالص ملی، شاخص قیمتها، یا سرمایه‌گذاری ملی است. تعریف دقیقتر این است که «اقتصادسنجی علم تحلیلهای آماری از مدل‌های

1. Econometrics

۲. برای مطالعه شکل‌گیری و پیدایش اقتصادسنجی، و بررسی توسعه‌های اولیه و پیشرفتهای جدید، به مقاله محمدهاشم پسران در *The New Palgrave. A Dictionary of Economics*, 1987 مراجعه کنید.

اقتصادی است.) منظور از مدل‌های اقتصادی، صورت منظم و ریاضی توابع یا روابط^۱ اقتصادی است که در ذیل تعریف می‌شود.

۲. مدل

هرگاه متغیرهای اقتصادی مانند تولید، مصرف، سرمایه‌گذاری و قیمت در روابط معینی تعریف شوند زمینه مناسبی برای شکل‌گیری یک مدل اقتصادی فراهم می‌شود. در حالی که نظریه‌های اقتصادی سعی دارند تا این روابط و تحولات آنها را تفسیر کنند، اقتصادسنجی می‌خواهد این روابط را به زبان آماری بیان کند. یک مدل اقتصادی در واقع بیان دیگری از رابطه یا روابط اقتصادی است که معمولاً به زبان ریاضی بیان می‌شود. یک متغیر اقتصادی مانند مصرف را در نظر می‌گیریم؛ آنچه در خارج وجود دارد نوسانات یا تغییرات مصرف است. مشاهده تغییر، انگیزه تفسیر تغییر را به وجود می‌آورد که خود زمینه‌ساز مطالعات و تحلیلهای اقتصادی است. اولین قدم در تفسیر تغییر، ساختن مدل تغییر است؛ یعنی مدلی که در آن، تغییر می‌تواند تفسیر شود. چگونگی ساختن مدلها به روش مدل‌سازی و مبانی نظری مدل بر می‌گردد که از موضوع این کتاب خارج است. یکی از ساده‌ترین و در عین حال رایجترین روشهای مدل‌سازی، طراحی یک یا چند رابطه ریاضی است که بتواند چند متغیر را با یکدیگر مرتبط سازد. برای مثال، در تفسیر تغییرات مصرف، یک متغیر دیگر مانند درآمد را در نظر گرفته، فرض می‌کنیم که مصرف تابعی از درآمد است. اگر مصرف و درآمد در زمان را به ترتیب با C_t و Y_t نشان دهیم خواهیم داشت

$$C_t = f(Y_t) \quad (۱-۱)$$

به عبارت فوق که در واقع یک رابطه ریاضی بین مصرف و درآمد است، معادله یا مدل مصرف گفته می‌شود. برای تفسیر تغییرات مصرف می‌توان علاوه بر Y_t متغیرهای دیگری را نیز در نظر گرفت. برای مثال، معادله

۱. در اینجا تفاوت مختصر بین تابع و رابطه را در نظر نگرفته‌ایم.

$$C_t = f(Y_t, C_{t-1}), \quad (1.2)$$

مدل دیگری برای مصرف است که در آن C_{t-1} مقدار مصرف در سال گذشته است. در گذشته واژه مدل فقط برای سیستم معادله‌ها به معنای چند معادله مرتبط با یکدیگر به کار می‌رفت. برای مثال، اگر فرض کنیم مصرف تابعی از درآمد است در حالی که درآمد خود به صورت مجموع مصرف و سرمایه‌گذاری تعریف شود، آنگاه خواهیم داشت

$$C_t = f(Y_t), \quad (1.3)$$

$$Y_t = C_t + I_t,$$

که در آن I_t سرمایه‌گذاری است. به معادله‌های ۱.۳، اصطلاحاً یک «مدل» می‌گوییم. با وجود این، واژه مدل را می‌توان در حال حاضر علاوه بر سیستم معادله‌ها، برای تک معادله نیز به کار برد.

تک معادله‌های ۱-۱ و ۱-۲ و سیستم معادله‌های ۱-۳ بسیار عمومی است؛ پس نمی‌تواند کاربرد چندانی داشته باشد؛ بنابراین برای اولین قدم، باید شکل ریاضی این معادله‌ها معین شود. فرض کنیم که این معادله‌ها همگی خطی هستند، آنگاه تک معادله‌های ۱-۱ و ۱-۲ و سیستم معادله‌های ۱-۳ را می‌توان به ترتیب چنین نوشت،

$$C_t = \alpha + \beta Y_t, \quad (1.4)$$

$$C_t = \alpha + \beta Y_t + \gamma C_{t-1}, \quad (1.5)$$

$$C_t = \alpha + \beta Y_t, \quad (1.6)$$

$$Y_t = C_t + I_t,$$

که در آن α و β را اصطلاحاً «پارامترهای» مدل می‌گویند؛ بنابراین این مدلها هم

برحسب پارامترها و هم برحسب متغیرها، خطی است.

فرض خطی بودن این مدلها در سهولت ریاضیات و ابعاد محاسباتی آن بسیار مفید است؛ سؤالی که مطرح می‌شود این است که چگونه می‌توان فرض خطی بودن مدل مصرف را توجیه کرد؟ در پاسخ باید گفت که یکی از اهداف اصلی مدل مصرف این است که بتواند تغییرات مصرف را که در عینیت مشاهده شده است به نحو رضایت‌بخشی تفسیر کند. اگر با فرض خطی بودن مدل مصرف بتوان به این هدف رسید، این فرض بر هر فرض دیگری ترجیح دارد؛ زیرا از هر فرض دیگری ساده‌تر است. البته معیار و نحوه محاسبهٔ درجهٔ توانایی مدل در تفسیر عینیت، مسأله‌ای است که در قسمت ۱-۴ بررسی خواهد شد. در پایان مفید است به این نکته نیز توجه شود که به مدل‌های ۱-۴ و ۱-۵ و ۱-۶ در اصطلاح «مدل‌های ریاضی» مصرف می‌گویند؛ زیرا در این مدلها فقط به ابعاد ریاضی رابطهٔ مصرف با درآمد توجه شده است.

۳. رگرسیون

رگرسیون^۱ در لغت به معنای «بازگشت به مراحل قبلی در یک مسیر تحول و توسعه» است. نظریه رگرسیون تقریباً در قلمرو همین معنی در روانشناسی و جامعه‌شناسی به کار می‌رود. در آمار، اصطلاح رگرسیون را اولین بار گالتن^۲ در مطالعات جمعیت‌شناختی به کار برد. گالتن در نمونه‌های مختلف آماری سعی کرد رابطه‌ای بین قد فرزندان با قد والدین برقرار کند. با اینکه معمولاً والدین قد بلند فرزندان قد بلندی دارند و برعکس، گالتن به این نتیجه رسید که در جامعه گرایشی وجود دارد که قد فرزندان به سمت متوسط قد جامعه میل می‌کند.^۳ این میل کردن را گالتن «رگرسیون» نامید.

در متون اقتصادسنجی، رگرسیون هیچگاه در این معنی به کار نمی‌رود. در اقتصادسنجی بیشتر اصطلاحات «تحلیل رگرسیون» و «مدل رگرسیون» رایج است.

1. Regression

2. Sir Francis Galton (1822 - 1911).

۳. گالتن این حقیقت را به صورت regression towards mediocrity بیان می‌کند. برای توضیح بیشتر به مقاله گالتن در سال ۱۸۸۶ مراجعه شود.

تحلیل رگرسیون در واقع بدنه اصلی مطالعات اقتصادسنجی را تشکیل می‌دهد و به طور کلی درباره مدل‌های رگرسیون و نحوه تخمین آنها بحث می‌کند؛ بنابراین باید به تعریف مدل‌های رگرسیون بپردازیم تا بدین ترتیب مفهوم رگرسیون در اقتصادسنجی روشن شود.

۴. مدل رگرسیون خطی ساده

فرض می‌کنیم در یک جامعه مفروض تغییرات مصرف را در ۵ سال مشاهده کرده و داده‌های $C_t = 2, 3, 5, 6, 9$ را برحسب میلیارد تومان به دست آورده‌ایم. می‌خواهیم چگونگی این تغییرات را تفسیر کنیم. برای این منظور باید متغیر یا متغیرهایی را در نظر بگیریم که بتواند این تغییرات را توضیح دهد. فرض کنید مدل ۴-۱ مدل مطلوب مصرف است؛ یعنی این مدل می‌تواند تغییرات مصرف را به صورت رضایت‌بخشی توضیح دهد، مگر آنکه خلاف آن ثابت شود. این مدل خطی مصرف را، که در واقع یک مدل ریاضی است، یک بار دیگر می‌نویسیم:

$$C_t = \alpha + \beta Y_t$$

بدیهی است باید آمار و مشاهدات Y_t را نیز داشته باشیم؛ زیرا بر اساس مدل مفروض، تغییرات مصرف بر اساس تغییرات درآمد بیان می‌شود. فرض کنید تغییرات Y_t را نیز برای همان ۵ سال به ترتیب $Y_t = 3, 5, 9, 10, 13$ مشاهده کرده‌ایم. مفید است در اینجا به این نکته اشاره کنیم که داده‌ها و مشاهدات متغیرهای موجود در یک مدل معمولاً در سه نوع مختلف می‌تواند وجود داشته باشد: داده‌های سری زمانی^۱، داده‌های مقطع زمانی^۲ و داده‌های تلفیقی یا «پانل»^۳.

داده‌های سری زمانی، مقادیر یک متغیر را در نقاط متوالی در زمان، اندازه‌گیری می‌کند. این توالی می‌تواند سالانه، فصلی، ماهانه، هفتگی یا حتی به صورت پیوسته باشد. مثالی که تا به حال موضوع بحث ما بوده است، یک سری زمانی است که

مشاهدات را به صورت سالانه ارائه می‌کند. معمولاً از اندیس t برای سری زمانی استفاده می‌کنند. برای مثال، اگر متغیر مورد نظر، X ناپیوسته باشد سری زمانی آن را با X_t و در حالتی که X پیوسته باشد با $X(t)$ نشان می‌دهند؛ مثلاً داده مصرف در زمان، یعنی C_t ، یک متغیر ناپیوسته است، چون می‌توان فقط در فواصل معینی از زمان مقادیر آن را ملاحظه کرد. متغیرهای پیوسته در اقتصاد، کاربرد چندانی ندارد.

داده‌های مقطع زمانی، مقادیر یک متغیر را در زمان معین و روی واحدهای متعدد اندازه‌گیری می‌کند. این واحدها می‌توانند افراد، خانوارها، واحدهای تولیدی، صنایع، نواحی مختلف و حتی کشورهای مختلف باشد. می‌توان داده‌های درآمد و مصرف خانوارهای مختلف را در سال معینی جمع‌آوری کرد. معمولاً از اندیس i برای داده‌های مقطعی استفاده می‌کنند؛ بنابراین C_i و Y_i نشان‌دهنده مصرف و درآمد خانوارهای مختلف در زمان معین است.

داده‌های تلفیقی در واقع بیان‌کننده داده‌های مقطعی در طی زمان است. بنابراین حجم مشاهدات در داده‌های تلفیقی نسبتاً زیاد است. در سالهای اخیر، کاربرد داده‌های تلفیقی در اقتصادسنجی افزایش بسیاری یافته است. معمولاً داده‌های تلفیقی و داده‌های مقطعی در «اقتصادسنجی خرد»^۱ به کار می‌رود که موضوع آن بررسی روشهای اقتصادسنجی در اقتصاد خرد است، در حالی که داده‌های سری زمانی به طور کلی موضوع کار «اقتصادسنجی کلان»^۲ است که روشهای اقتصادسنجی را در سطح اقتصاد کلان مطالعه می‌کند.

حال که با انواع مختلف آمار و مشاهدات در اقتصادسنجی آشنا شدیم یادآوری می‌کنیم که به مدل ۱-۴ یک مدل ریاضی خطی ساده گفته می‌شود. واژه‌های ریاضی، خطی و ساده را به ترتیب تعریف می‌کنیم. می‌دانیم مدل فوق یک مدل ریاضی است؛ زیرا فقط رابطه ریاضی بین مصرف و درآمد را منعکس کرده است. همچنین یک مدل خطی است، زیرا پارامترها و متغیرهای این مدل در یک رابطه خطی تعریف شده است؛

و سرانجام یک مدل ساده است؛ زیرا در سمت راست فقط یک متغیر وجود دارد. اگر بیش از یک متغیر در سمت راست باشد، یک مدل خطی «چند متغیره» خواهیم داشت. برای سهولت ارائه مباحث، مشاهدات مصرف و درآمد را در جدول ۱-۱ تنظیم می‌کنیم.

جدول ۱-۱

۱	مصرف C_t	درآمد Y_t
۱	۲	۳
۲	۳	۵
۳	۵	۹
۴	۶	۱۰
۵	۹	۱۳

سؤال این است که آیا تغییرات

درآمد، در چهارچوب شکل ریاضی مدل ۱-۴ می‌تواند تغییرات مصرف را به نحو رضایت‌بخشی توضیح دهد؟ پاسخ منفی است:

اولاً، علاوه بر درآمد، به منزلهٔ عامل اصلی تغییرات مصرف، قطعاً

عوامل دیگری نیز وجود دارد که می‌تواند تغییرات مصرف را توضیح دهد. بسیاری از این عوامل، مانند مصرف سال قبل، را می‌توان شناخت که معمولاً آنها را وارد مدل می‌کنیم؛ و در نتیجه مدل رگرسیون ساده به رگرسیون چند متغیره تبدیل می‌شود. اما عوامل دیگری مانند چگونگی انتظارات مصرف‌کننده نسبت به تغییر در پارامترهای مختلف اقتصادی و درجهٔ عدم اطمینان نسبت به تحولات آیندهٔ اقتصادی وجود دارد که بیان کمی آنها معمولاً بسیار مشکل است.

ثانیاً، هیچ دلیلی وجود ندارد که رابطهٔ بین مصرف و متغیر یا متغیرهای سمت راست، خطی باشد.

ثالثاً، مصرف، خصوصیتی از رفتار انسانی است و می‌دانیم در رفتار انسان همواره عناصر تصادفی غیر قابل پیش‌بینی وجود دارد که اساساً نمی‌توان آنها را با مدل‌های معین ریاضی از نوع مدل ۱-۴ بیان کرد.

البته می‌توان دلایل دیگری مبنی بر عدم قابلیت مدل ۱-۴ در بیان دقیق تغییرات مصرف بیان کرد، مانند خطای اندازه‌گیری^۱ متغیر مصرف. این دلایل برای قبول این

حقیقت که مدل ۱-۴ دقیق نیست و خطا دارد کافی است. به این خطا اصطلاحاً «جمله اختلال»^۱ می‌گویند؛ زیرا تعادل ریاضی مدل مفروض را مختل می‌کند. اگر تأثیر تمام عوامل فوق را بر تغییرات مصرف با U_1 نشان دهیم، آنگاه مدل ۱-۴ را می‌توان چنین نوشت،

$$C_t = \alpha + \beta Y_t + U_t \quad (1.7)$$

که در آن U_t جمله اختلال بوده و منعکس‌کننده خطای ما در بیان مدل مصرف است. بدیهی است U_t یک متغیر تصادفی است؛ زیرا عناصری که مقادیر آن را تعیین می‌کند، تصادفی است. به معادله ۱-۷ مدل رگرسیون مصرف می‌گویند؛ بنابراین تفاوت کلی مدل‌های ریاضی و مدل‌های رگرسیون در جمله اختلال است. هر گاه به مدل‌های ریاضی یک جمله اختلال - که یقیناً تصادفی است - اضافه کنیم، مدل ریاضی به مدل رگرسیون تبدیل خواهد شد. با توجه به تعاریف قبلی می‌توان گفت که معادله ۱-۷ یک «مدل رگرسیون خطی ساده»^۲ است. همچنین به معادله ۱-۷ یک مدل ساختاری^۳ مصرف نیز می‌گویند؛ زیرا ساختار مصرف را مبتنی بر یک تئوری مصرف بیان می‌کند. به همین دلیل به پارامترهای α و β پارامترهای ساختاری^۴ گفته می‌شود.

۱-۲ خصوصیات آماری مدل رگرسیون خطی ساده مدل رگرسیون خطی زیر را در نظر می‌گیریم،

$$Y_t = \alpha + \beta X_t + U_t \quad (1.8)$$

در این مدل می‌خواهیم تغییرات Y_t را به کمک تغییرات X_t توضیح دهیم. به همین دلیل به X_t متغیر توضیحی^۵ و به Y_t متغیر توضیح داده شده^۶ می‌گویند. اصطلاحات دیگری نیز برای X_t و Y_t به کار می‌رود؛ برای مثال، به X_t متغیر مستقل^۷ و به Y_t متغیر وابسته^۸

1. Disturbance Term

2. Simple Linear Regression Model

3. Structural Model

4. Structural Parameters

5. Explanatory Variable

6. Explained Variable

7. Independent Variable

8. Dependent Variable

می‌گویند؛ زیرا مقادیری که X_t می‌گیرد به طور مستقل تعیین می‌شود؛ یعنی تابعی از مدل ۱-۸ نیست. به عبارت دیگر مقادیر X_t داده شده است، در حالی که مقادیر Y_t دقیقاً تابعی از مقادیر X_t و کیفیت مدل رگرسیون مفروض است و به همین دلیل Y_t یک متغیر تابع یا وابسته محسوب می‌شود. اصطلاحات متغیر برون‌زا^۱ و متغیر درون‌زا^۲ به ترتیب برای X_t و Y_t نیز به کار می‌رود؛ زیرا فرض بر این است که مقادیر X_t خارج از مدل رگرسیون مفروض تعیین شده و در نتیجه برون‌زا است، در حالی که مقادیر Y_t در داخل و براساس قانونمندی مدل رگرسیون تعیین می‌شود و به همین دلیل درون‌زا خواهد بود. همچنین متغیرهای X_t و Y_t را می‌توان متغیرهای «غیر تصادفی»^۳ و «تصادفی»^۴ نامید؛ زیرا مقادیری که X_t می‌گیرد به طور مستقل تعیین شده و تابعی از جمله اختلال مدل (U_t) نیست که خود یک متغیر تصادفی است، در حالی که Y_t دقیقاً تابعی از U_t بوده و یک متغیر تصادفی خواهد بود. در این کتاب بیشتر از واژه متغیرهای درون‌زا و برون‌زا استفاده می‌شود، هر چند در بعضی موارد و برحسب ضرورت از اصطلاحات دیگر نیز استفاده خواهد شد.

در مدل رگرسیون فوق، α و β پارامترهای مدل هستند. مقادیر واقعی این پارامترها را نمی‌دانیم، بنابراین سعی می‌کنیم مقادیر آنها را تخمین^۵ بزنیم. یکی از وظایف اصلی اقتصادسنجی تخمین پارامترهای مدل رگرسیون است. با اینکه مقادیر واقعی پارامترهای α و β را نمی‌دانیم اما مطمئن هستیم این پارامترها غیر تصادفی هستند؛ زیرا فرض بر این است که تابعی از جمله خطای مدل نیستند. با توجه به اینکه X_t نیز غیر تصادفی است به $(\alpha + \beta X_t)$ در اصطلاح قسمت معین^۶ مدل رگرسیون می‌گویند. همان گونه که دیدیم U_t نیز قسمت نامعین یا تصادفی مدل است. بنابراین نتیجه می‌گیریم که ساختار ریاضی مدل رگرسیون مفروض به گونه‌ای است که برای تفسیر تغییرات متغیر درون‌زای Y_t ، تغییرات آن را به دو قسمت تقسیم می‌کند: تغییرات معین،

1. Exogenous

2. Endogenous

3. Non-stochastic

4. Stochastic

5. Estimation

6. Deterministic Part

یعنی $(\alpha + \beta X_t)$ که در واقع همان مدل ریاضی است و تغییرات تصادفی، یعنی U_t که جمله اختلال مدل محسوب می شود. به مجموعه این دو تغییرات مدل رگرسیون گفته می شود.

تحلیلهای مدل رگرسیون می تواند در حوزه های مختلف صورت پذیرد. در این قسمت فقط به سه قلمرو اشاره می کنیم:

۱. آثار مختلف سیاستگذاریهای مرتبط با متغیرهای برونزا را می توان بر تغییرات متغیر درونزا تحلیل کرد. که در اصطلاح به آن تحلیل «سیاستگذاری» می گویند؛ برای مثال، در بحث مدل مصرف ۱-۷ می توان آثار تغییر در سطوح مختلف درآمد را بر کیفیت ایجاد تغییر در سطوح مختلف مصرف تجزیه و تحلیل کرد.

۲. به ازای مقادیر معینی از متغیر یا متغیرهای برونزا می توان مقادیر متغیر درونزا را پیش بینی کرد. برای مثال، می توان به ازای مقدار معینی از سطح درآمد در آینده مقدار مصرف را پیش بینی کرد.

۳. به کمک مدل های رگرسیون می توان تحلیلهای آماری، تأثیر و یا بهتر بگوییم معنی دار بودن^۲ تأثیر متغیرهای برونزا را بر متغیر درونزا مطالعه کرد. برای مثال، می توان به این سؤال پاسخ داد که آیا از نظر آماری تأثیر تغییرات درآمد بر تغییرات مصرف معتبر است؟ به عبارت دیگر آیا می توان تغییرات مصرف را به کمک تغییرات درآمد توضیح داد؟

با بررسی مدل رگرسیون ۱-۸ به سهولت ملاحظه می شود که هر گونه پیشرفت در تحلیلهای رگرسیونی متوقف به شناخت بیشتر از جمله اختلال مدل (U_t) است. در واقع در یک مدل رگرسیون، U_t با اینکه نقش مهمی ایفا می کند اما بنا بر تعریف ناشناخته است. هرگاه کوشش کنیم اجزایی از U_t را بشناسیم و آنها را اندازه گیری کنیم، این اجزای شناخته شده، در قسمت معین مدل قرار می گیرد و مجموعه عوامل مجهولی که باقی می ماند U_t را شکل می دهد؛ بنابراین U_t هیچگاه قابل مشاهده و اندازه گیری نیست.

1. Policy Analysis

2. Forecast or Prediction

3. Statistical Significance

در نتیجه راه خروج از این تنگنای نظری این است که یک سری فرضهای منطقی در مورد U_t مطرح کنیم تا بر آن اساس بتوان به تحلیلهای رگرسیونی ادامه داد. این فرضها با یک فرض در مورد متغیر برونزا، در ذیل با عنوان فرضهای کلاسیک مدلهای رگرسیون مطرح می شود.

۱. فرضهای کلاسیک مدلهای رگرسیون خطی ساده

در قسمت قبل، تعریف U_t و عوامل ایجادکننده آن را به طور مختصر بررسی کردیم. دیدیم که U_t را در اصطلاح جمله اختلال می گویند؛ زیرا وجود U_t موجب می شود که قسمت متین مدل که در واقع همان مدل ریاضی $(Y_t = \alpha + \beta X_t)$ است مختل شود. اصطلاحات دیگری نیز در متون اقتصادسنجی برای U_t به کار می رود؛ برای مثال، U_t را خطای معادله یا خطای مدل نیز می گویند؛ زیرا وجود U_t مبین خطای موجود در بیان دقیق عوامل و متغیرهایی است که می تواند تغییرات متغیر درونزا یعنی Y_t را در مدل ۱-۸ توضیح دهد. اصطلاح دیگری که برای U_t به کار می رود «جمله تصادفی» است؛ زیرا تأثیری که U_t بر تغییرات Y_t می گذارد، کاملاً تصادفی است. در گذشته بیشتر اصطلاح جمله اختلال در کتابهای اقتصادسنجی رایج بود، اما به نظر می رسد امروز گرایش به کاربرد جمله خطای مدل بیشتر شده است. با توجه به اینکه هنوز در کتابهای مرجع اقتصادسنجی، از واژه جمله اختلال استفاده می شود، ما نیز در این کتاب همین اصطلاح را برگزیده ایم.

مهمترین نکته در مورد U_t ، خصوصیت تصادفی بودن آن است. با توجه به تعریفی که از U_t ارائه شد، بدیهی است این فرض قابل قبول است و خلاف آن را نمی توان تصور نمود. U_t یک متغیر تصادفی است و مانند همه متغیرهای تصادفی دارای یک تابع توزیع احتمال است. می دانیم هر تابع توزیع احتمال دارای میانگین و واریانس است. سؤال مهمی که می توان مطرح کرد این است که خصوصیات آماری و شکل تابع

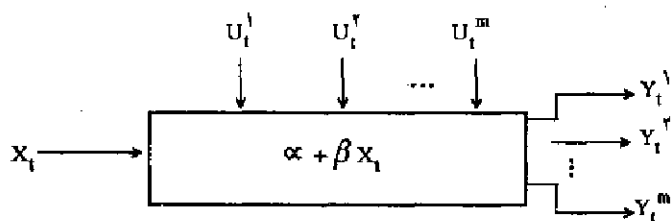
توزیع احتمال متغیر تصادفی U_t چیست؟ پاسخ به این سؤال با عنوان فرضهای کلاسیک چنین است:

اولین فرضی که در مورد U_t مطرح می شود این است که میانگین یا امید ریاضی آن صفر است. می دانیم امید ریاضی را با عملگر E نشان می دهند؛ بنابراین در یک نمونه n تایی داریم،

$$E(U_t) = 0, t = 1, 2, 3, \dots, n. \quad (1-9)$$

مفهوم این فرض اهمیت بسیار دارد. این فرض در واقع بدین معنی است که به ازای هر مقدار معین از X_t ، میانگین تمام مقادیر ممکن U_t برابر صفر است. ابتدا باید به این سؤال پاسخ داد که به ازای یک مقدار ثابت و معین X_t چگونه U_t می تواند مقادیر مختلفی داشته باشد؟ بدیهی است دقت در میانگین مقادیر مختلف U_t که آیا صفر است یا خیر بعد از پاسخ به این سؤال مطرح می شود.

ظهور مقادیر مختلف U_t به اعتبار فرض آزمایشهای فرضی تکراری^۲ به ازای مقدار معین و ثابت X_t است. برای تبیین این مطلب به نمودار ۱-۱ توجه می کنیم.



نمودار ۱-۱ مفهوم امید ریاضی U_t

مقدار X_t را ثابت فرض کرده، با توجه به مقادیر ثابت α و β نتیجه می گیریم که $(\alpha + \beta X_t)$ نیز مقدار معینی خواهد داشت. مقدار Y_t چیست؟ فرض می کنیم در دوره زمانی t هستیم. در طی این دوره کلیه عواملی که در تعیین U_t مؤثر هستند فعال می شود و مقدار U_t را شکل می دهند. می دانیم U_t یک متغیر تصادفی است، یعنی هر مقداری که

U_1 بدین ترتیب اتخاذ کند در واقع می توانست مقدار دیگری باشد. اگر بخواهیم این نکته را به زبانی دیگر بیان کنیم، کافی است بگوییم که در آزمایش فرضی اول، U_1 مقدار U_1 را می گیرد. این مقدار با $(\alpha + \beta X_1)$ جمع می شود و Y_1 را نتیجه می دهد. حال می گوییم اگر این آزمایش را دوباره انجام دهیم، با توجه به تصادفی بودن U_1 ، هیچ دلیلی وجود ندارد که U_1 درست همان مقدار قبلی را بگیرد؛ بنابراین باید گفت U_1 در آزمایش فرضی دوم، مقدار U_1 را خواهد داشت و وقتی با مقدار ثابت $(\alpha + \beta X_1)$ جمع شود، Y_1 را تشکیل می دهد. به همین ترتیب می توان m آزمایش فرضی را تصور کرد که مقادیر U_1^m, \dots, U_1^1 و U_1 ، با فرض ثابت $(\alpha + \beta X_1)$ به ترتیب مقادیر Y_1^m, \dots, Y_1^1 و Y_1 را نتیجه می دهند. انتظار این است که به ازای مقادیر بزرگ m ، میانگین یا امید ریاضی U_1 صفر باشد؛ زیرا مقادیر مختلف U_1^m, \dots, U_1^1 و U_1 در حول محور ثابت $(\alpha + \beta X_1)$ نوسان دارند؛ به عبارت دیگر بعضی از مقادیر فوق U_1 در بالای این خط ثابت قرار می گیرد و بعضی دیگر در پایین و در مجموع مقادیر مثبت و منفی، همدیگر را خنثی می کند؛ به گونه ای که در نهایت امید ریاضی یا میانگین U_1 برابر صفر خواهد شد. با وجود این، باید در نظر داشت که هیچیک از مقادیر U_1 در آزمایشهای تکراری قابل مشاهده نیست و آنچه مطرح شد در واقع دقت در خصوصیات آماری این متغیر تصادفی است.

دومین فرضی که برای U_1 مطرح می کنیم، ثابت بودن واریانس آن به ازای مقادیر مختلف X_1 است،

$$\text{Var}(U_i) = E(U_i)^2 = \sigma^2 \quad \text{و} \quad i = 1, 2, \dots, n. \quad (1.10)$$

اگر مقدار X_i در دوره i ، یعنی X_i را ملاحظه کنیم و نمودار ۱-۱ را برای آن بسازیم، m مقدار مختلف برای U_i به دست خواهیم آورد. واریانس این مقادیر مختلف U_i را با σ_i^2 نشان می دهیم؛ یعنی $\text{Var}(U_i) = \sigma_i^2$. به همین ترتیب اگر مقدار X_i را در دوره i ملاحظه کرده، واریانس U_i را به ترتیبی مشابه تعریف کنیم، خواهیم داشت: $\text{Var}(U_i) = \sigma_i^2$. فرض ثابت بودن واریانس U_i این است که برای تمام مقادیر i و زاین واریانسها با یکدیگر برابر باشد،

$$\text{Var}(U_i) = \text{Var}(U_j) = \sigma_i^2 = \sigma_j^2, \quad i, j = 1, 2, \dots, n.$$

با توجه به تعریف واریانس،

$$\text{Var}(U_i) = E[U_i - E(U_i)]^2,$$

و با در نظر گرفتن فرض اول، یعنی $E(U_i) = 0$ ، خواهیم داشت

$$\text{Var}(U_i) = E(U_i)^2 = \sigma^2.$$

هرگاه واریانس جمله اختلال ثابت باشد، می‌گوییم U_i واریانس همسانی^۱ دارد؛ در غیر این صورت U_i واریانس ناهمسان^۲ است.

سومین فرض U_i این است که U_i و U_j به ازای تمامی مقادیر $j \neq i$ از یکدیگر مستقلند؛ یعنی کوواریانس آنها صفر است.

$$\text{cov}(U_i, U_j) = E(U_i U_j) = 0, \quad i, j = 1, 2, \dots, n, \quad i \neq j \quad (1-11)$$

به عبارت دیگر، هرگاه دو مقدار X_i ، برای مثال، X_i و X_j را در نظر بگیریم، فرض بر این است که جمله‌های اختلال متناظر با آنها، یعنی U_i و U_j از یکدیگر مستقلند. همچنین با توجه به تعریف کوواریانس داریم

$$\text{cov}(U_i, U_j) = E[U_i - E(U_i)][U_j - E(U_j)] = 0,$$

و با در نظر گرفتن فرض اول، یعنی $E(U_i) = 0$ ، رابطه مذکور را می‌توان به صورت زیر نوشت،

$$\text{cov}(U_i, U_j) = E(U_i U_j) = 0.$$

در چنین حالتی می‌گوییم که جمله‌های اختلال خود همبستگی^۳ ندارد و در غیر این صورت (هنگامی که $E(U_i, U_j) \neq 0$) جمله‌های اختلال خود همبستگی خواهد داشت. اگر جمله‌های اختلال، واریانس همسان و خود همبستگی نداشته باشند، می‌گوییم U_i

«کروی»^۱ است.

چهارمین فرض U_i این است که تابع توزیع احتمال آن را نرمال بدانیم. بنابراین با توجه به فرضهای اول، دوم و سوم می‌توان گفت که U_i دارای توزیع مستقل نرمال با میانگین صفر و واریانس ثابت σ^2 است. اگر علامت $-$ را برای توزیع و IN را برای توزیع مستقل نرمال^۲ به کار ببریم، آنگاه خواهیم داشت

$$U_i \sim IN(0, \sigma^2). \quad (1-12)$$

در بعضی از منابع، به جای علامت IN از NID ^۳ استفاده می‌کنند،

$$U_i \sim NID(0, \sigma^2).$$

پنجمین و آخرین فرض از فرضهای کلاسیک این است که X_i یک متغیر غیر تصادفی است. این فرض بیشتر برای سهولت در استنتاج قضایا و نیز رسیدن به نتایج جالبتر در تخمین پارامترهاست. بدیهی است که می‌توان این فرض را نقض کرد و X_i را به صورت یک متغیر تصادفی در نظر گرفت. به هر حال فرض غیر تصادفی بودن X_i بدین معنی است که X_i از متغیر تصادفی U_i مستقل است (در این جلد تمام نتایج به صورت شرطی و با فرض غیر تصادفی بودن X_i بیان شده است).

در بعضی از کتابهای اقتصادسنجی، تصادفی بودن U_i و متساوی نبودن مقادیر مختلف X_i در یک نمونه را نیز فرضهای کلاسیک می‌نامند و در نتیجه هفت فرض را با این عنوان مطرح می‌کنند. در معادله ۱-۲۵ خواهیم دید که در صورت تساوی مقادیر X_i با یکدیگر، تخمین پارامترها ممکن نخواهد بود.

۲. خصوصیات آماری متغیر درون‌زا

یک بار دیگر مدل رگرسیون ۱-۸ را در نظر می‌گیریم،

1. Spherical Disturbances
3. Normally and Independently Distributed

2. Independent Normal

$$Y_i = \alpha + \beta X_i + U_i.$$

ملاحظه می‌شود که چون Y_i تابعی از جمله اختلال (U_i) است یک متغیر تصادفی خواهد بود؛ زیرا U_i یک متغیر تصادفی است. Y_i مانند هر متغیر تصادفی دیگر دارای تابع توزیع احتمال است و یک میانگین و یک واریانس خواهد داشت.

با توجه به اینکه در مدل رگرسیون خطی ۱-۸ متغیر درون‌زای Y_i تابعی از U_i است و علاوه بر این چون فرض کرده‌ایم U_i دارای توزیع نرمال است بنابراین متغیر تصادفی Y_i نیز دارای توزیع احتمال نرمال خواهد بود. برای یافتن میانگین آن کافی است از دو طرف رابطه ۱-۸ امید ریاضی بگیریم خواهیم داشت

$$E(Y_i) = E(\alpha) + E(\beta X_i) + E(U_i).$$

با توجه به اینکه α ثابت است پس $E(\alpha) = \alpha$. همچنین فرض کرده‌ایم X_i یک متغیر غیر تصادفی بوده و در آزمایشهای تکراری ثابت است؛ بنابراین با توجه به ثابت بودن β داریم

$$E(\beta X_i) = \beta E(X_i) = \beta X_i.$$

اگر اولین فرض مربوط به U_i ، یعنی $E(U_i) = 0$ را نیز در نظر بگیریم، آنگاه خواهیم داشت

$$E(Y_i) = \alpha + \beta X_i. \quad (1-13)$$

به رابطه ۱-۱۳ معمولاً تابع رگرسیون جامعه^۱ می‌گویند. همان گونه که بعداً خواهیم دید اگر به جای α و β تخمین آنها، یعنی $\hat{\alpha}$ و $\hat{\beta}$ را بگذاریم در آن صورت تابع رگرسیون نمونه^۲ حاصل خواهد شد^۳. تابع خطی ۱-۱۳ در حقیقت مکان هندسی نقاطی است که برای آنها $E(U_i) = 0$ است.

برای محاسبه واریانس Y_i کافی است به تعریف واریانس مراجعه کنیم،

$$\text{Var}(Y_i) = E[Y_i - E(Y_i)]^2.$$

1. Population Regression Function

2. Sample Regression Function

۳. معمولاً تخمین پارامترها را با علامت $\hat{\cdot}$ نشان می‌دهیم و آن را کلاه (Hat) می‌خوانیم.

با جایگزینی مقادیر Y_1 و $E(Y_1)$ به ترتیب از رابطه‌های ۱-۸ و ۱-۱۳، مقدار $\text{Var}(Y_1)$ به ترتیب زیر محاسبه می‌شود،

$$\begin{aligned} \text{Var}(Y_1) &= E[\alpha + \beta X_1 + U_1 - (\alpha + \beta X_1)]^2 \\ &= E(U_1)^2 = \sigma^2. \end{aligned} \quad (1-14)$$

واریانس متغیر درون‌زای Y_1 دقیقاً برابر واریانس جمله اختلال مدل رگرسیون است. نتیجه می‌گیریم که Y_1 یک متغیر درون‌زای نرمال با میانگین $(\alpha + \beta X_1)$ و واریانس σ^2 است. چون فرض کرده بودیم مقادیر مختلف U_1 دارای توابع توزیع نرمال مستقل از یکدیگر هستند بنابراین Y_1 نیز دقیقاً چنین خصوصیتی خواهد داشت. این خصوصیات آماری Y_1 را می‌توان به صورت زیر نشان داد،

$$Y_1 \sim \text{IN}[(\alpha + \beta X_1), \sigma^2]. \quad (1-15)$$

۳. بیان هندسی خصوصیات آماری مدل‌های رگرسیون خطی ساده آنچه را که تا به حال در مورد خصوصیات آماری یک مدل رگرسیون خطی گفته‌ایم، در یک نمودار نشان خواهیم داد. در یک مدل رگرسیون خطی ساده سه متغیر اصلی داریم: Y_1 ، X_1 و U_1 . برای سهولت کافی است در یک نمودار سه بعدی مقادیر Y_1 ، X_1 و $P(U_1)$ را منعکس سازیم.

ابتدا باید نمودارهای تابع رگرسیون جامعه و تابع رگرسیون نمونه را رسم کرد. تابع رگرسیون جامعه، یعنی مدل ۱-۱۳ را یک بار دیگر همراه با مدل رگرسیون خطی ۱-۸ می‌نویسیم،

$$Y_1 = \alpha + \beta X_1 + U_1,$$

$$E(Y_1) = \alpha + \beta X_1.$$

یکی از مباحث اصلی تحلیلهای رگرسیونی تخمین مدل‌های رگرسیون است. منظور از تخمین یک مدل رگرسیونی در واقع تخمین پارامترهای آن است؛ بنابراین اگر تخمینهای

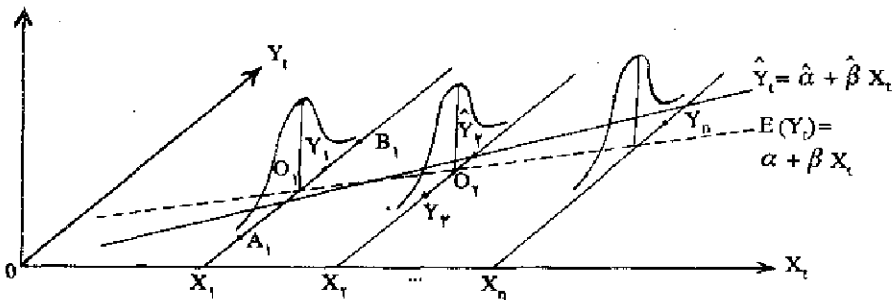
α ، β و Y_i را به ترتیب با $\hat{\alpha}$ ، $\hat{\beta}$ و \hat{Y}_i نشان دهیم، تابع یا مدل رگرسیون نمونه - که در واقع تخمین مدل رگرسیون ۱-۸ است - عبارت خواهد بود از

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i. \quad (1-16)$$

رسیدن به معادله ۱-۱۶ یکی از اهداف مهم تخمین در تحلیل‌های رگرسیونی است. در مباحث آینده خواهیم دید که $\hat{\alpha}$ و $\hat{\beta}$ قابل اندازه‌گیری است و در نتیجه می‌توان مدل رگرسیون نمونه ۱-۱۶ را رسم کرد. اما پارامترهای واقعی جامعه (α و β) هیچگاه قابل مشاهده و اندازه‌گیری نیستند؛ زیرا U_i اساساً قابل مشاهده نیست؛ بنابراین نمودار هندسی مدل رگرسیون جامعه ۱-۱۳ در عمل امکان‌پذیر نیست.

در نمودار ۱-۲ ابتدا منحنی تغییرات مدل رگرسیون جامعه را به صورت یک خط مستقیم و فرضی با نقطه چین رسم می‌کنیم؛ آنگاه منحنی تغییرات مدل رگرسیون نمونه را به صورت یک خط مستقیم نشان می‌دهیم. از این به بعد به مدل رگرسیون نمونه «تخمین مدل رگرسیون جامعه» یا به طور خلاصه «تخمین مدل رگرسیون» می‌گوییم.

$$P(U_i) = P(Y_i | X_i)$$



نمودار ۱-۲ خصوصیات آماری مدل رگرسیون خطی

X_i در آزمایش‌های فرضی تکراری ثابت فرض می‌شود. برای n دوره مشاهده، یعنی $n = 1, 2, 3, \dots$ فرض کنیم X_1 (مقدار X) را در اولین دوره مشاهده ثابت گرفته‌ایم، این مقدار ثابت در نمودار ۱-۲ با OX_1 مشخص شده است. می‌دانیم به ازای X_1 معین، Y_1 می‌تواند مقادیر متفاوتی داشته باشد. اما این مقادیر دقیقاً تابعی از این نکته

است که U_1 چه مقادیری کسب می‌کند. فرض بر این است که تابع توزیع احتمال U_1 نرمال است؛ بنابراین قلمرو تغییرات U_1 مثلاً محدود به A_1, B_1 است. حال با فرض ثبات X_1 می‌گوییم که در آزمایش اول فرض می‌کنیم Y مقدار Y_1, X_1 را گرفته است. اگر مدل رگرسیون ۱-۸ را برای دوره اول، یعنی $t = 1$ بنویسیم، خواهیم داشت

$$Y_1 = \alpha + \beta X_1 + U_1 .$$

این واقعیت که در آزمایش اول، Y در سال اول مقدار Y_1 را گرفته است، این نکته را بیان می‌کند که جمله اختلال (U) در این آزمایش، مقدار مثبت U_1, Y_1 را داشته است؛ زیرا می‌دانیم $E(Y_1) = \alpha + \beta X_1$ در واقع همان $(\alpha + \beta X_1)$ است. البته در آزمایشهای فرضی دیگر، U می‌تواند هر مقدار دیگری را در قلمرو A_1, B_1 بگیرد و بدین ترتیب مقادیر مثبت یا منفی دیگری به مقدار ثابت $\alpha + \beta X_1$ که در واقع همان $E(Y_1)$ است اضافه کند. نتیجه می‌گیریم که به ازای X_1 ثابت و بعد از ظهور مقدار Y_1 در عینیت، می‌توان مقدار U_1 را به صورت زیر نوشت،

$$\begin{aligned} U_1 &= X_1 Y_1 - (\alpha + \beta X_1) , \\ &= Y_1 - E(Y_1) . \end{aligned}$$

بنابراین با توجه به اینکه مقدار مشاهده شده Y_1 در عینیت بالای خط فرضی رگرسیون جامعه $(E(Y_1) = \alpha + \beta X_1)$ یا پایین آن قرار می‌گیرد، می‌توان نتیجه گرفت که مقدار U_1 مثبت یا منفی خواهد بود.

همین بحث را دقیقاً برای دوره‌های $n, \dots, 3, 2, t$ تکرار می‌کنیم و به ازای X_t, X_n تا X_2 و مشاهده مقادیر Y_2, Y_3, \dots, Y_n تا Y_2 مقادیر U را استخراج می‌کنیم. در حالت کلی می‌توان چنین نوشت

$$\begin{aligned} U_t &= X_t Y_t - (\alpha + \beta X_t) , \\ &= Y_t - E(Y_t) . \end{aligned} \tag{1.17}$$

در نمودار ۱-۲ فرض $E(U_t) = 0$ را می‌توان بدین صورت منعکس کرد که به ازای

هر مقدار X_t ثابت، منحنی تابع توزیع احتمال U_t در حول محور $E(Y_t) = \alpha + \beta X_t$ متمرکز بوده و دارای میانگین صفر است. به همین دلیل است که منحنی $E(Y_t) = \alpha + \beta X_t$ را می توان مکان هندسی نقاطی دانست که برای آنها $E(U_t)$ برابر صفر است فرض ثابت بودن واریانس U_t به ازای مقادیر مختلف t (فرض واریانس همسانی) چیزی نیست جز اینکه بگوییم منحنی نمایش هندسی توابع توزیع احتمال U_t به ازای مقادیر ثابت X_1, X_2, \dots, X_n دارای پراکندگی یکسانی است. فرض عدم وجود خودهمبستگی (صفر بودن کوواریانس U_i و U_j) به این معنی است که مقدار U در دوره i مستقل از مقداری است که U در دوره j گرفته است یا خواهد گرفت و برعکس. سرانجام تعبیر هندسی فرض غیر تصادفی بودن X_t بدین گونه است که مقدار X_t مثلاً در دوره دوم، $t=2$ ، به طور مطلق تابعی از مقادیر جمله اختلال، (U_t) در این دوره، دوره قبل یا دوره بعد نیست؛ به عبارت دیگر، مقدار X در دوره i اساساً تابعی از این نیست که جمله اختلال در دوره $i+1$ چه مقداری خواهد داشت یا در دوره $i-1$ چه مقداری داشته است.

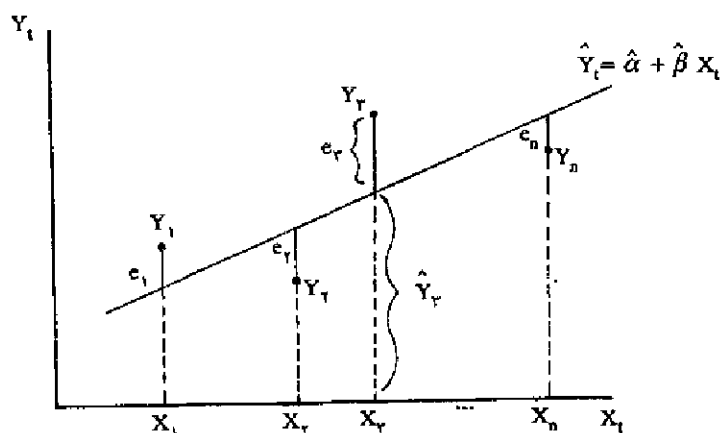
۱-۳ تخمین مدل رگرسیون خطی: روش حداقل مربعات معمولی (OLS)

مدل رگرسیون خطی ساده ۱-۸ را یک بار دیگر می نویسیم،

$$Y_t = \alpha + \beta X_t + U_t.$$

می خواهیم پارامترهای این مدل (α و β) را تخمین بزنیم. برای این منظور مشاهدات متغیرهای برونزا و درونزا (Y_t و X_t) را برای n دوره، $t = 1, 2, 3, \dots, n$ ، جمع آوری کرده ایم. در نمودار ۱-۲ این مشاهدات را با $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ نشان داده ایم. برای سهولت بحث، صفحه مختصات XOY از نمودار سه بعدی ۱-۲ را در نمودار ۱-۳ نشان می دهیم. برای تخمین پارامترهای α و β ضرورتاً نیاز به یک روش تخمین داریم. در اقتصادسنجی روشهای مختلفی برای تخمین پارامترها وجود دارد. ساده ترین و بهترین روش تخمین، روش حداقل مربعات معمولی است. طرح اولیه این

روش را که معمولاً با OLS نشان داده می‌شود گاس^۱ ریاضیدان معروف آلمانی در قرن هیجدهم مطرح کرده است.



نمودار ۱.۳ روش تخمین OLS

مدل رگرسیون ۱.۸ در واقع یک مدل خطی است که می‌خواهد تغییرات متغیر درون‌زای Y_i را در خلال مشاهدات n دوره توضیح دهد. زیربنای فکری روش حداقل مربعات معمولی (OLS) این است که بتوانیم این مدل خطی را با یک مدل خطی دیگر از نوع

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i,$$

چنان تخمین بزنیم که در بالاترین سطح ممکن، بتواند n مشاهده ما را توضیح دهد. مدل فوق که همان مدل ۱.۱۶ است تخمین مدل رگرسیون ۱.۸ است. در حقیقت باید $\hat{\alpha}$ و $\hat{\beta}$ چنان مقادیری داشته باشند که مدل ۱.۱۶ نسبت به هر مدل خطی دیگر، بیشترین نزدیکی را به مشاهدات Y_1, Y_r, \dots, Y_n داشته باشد؛ به عبارت دیگر کمترین انحراف را از مشاهدات فوق نشان دهد.

اگر معادله خط تخمین مدل رگرسیون (۱.۱۶) را در نمودار ۱.۳ رسم کنیم

۱. Carl Friedrich Gauss 1777-1855؛ با اینکه روش حداقل مربعات ابتدا در کتاب لژاندر (A. M. Legendre) ریاضیدان معروف فرانسوی در سال ۱۸۰۵ به چاپ رسید، اما گاس از سال ۱۷۹۵ از این روش در حل بسیاری از مسائل استفاده کرده است. برای توضیحات بیشتر به پیوست «د» مراجعه کنید.

بنا بر اقتضای بحث این خط باید کمترین فاصله را با مشاهدات ما داشته باشد. به ازای مقادیر X_1, X_2, \dots, X_n ، تخمین مدل رگرسیون، یعنی معادله ۱-۱۶، مقادیر متناظر Y_1 را به ترتیب به صورت $\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_n$ تخمین می‌زند. بدیهی است تخمین ما از مدل رگرسیون نمی‌تواند مقداری برای \hat{Y}_1 بدهد که دقیقاً برابر Y_1 باشد. به همین ترتیب در حالت کلی نمی‌توان انتظار داشت \hat{Y}_i برابر Y_i شود. اختلاف بین مشاهده (Y_i) و تخمین (\hat{Y}_i) را اصطلاحاً «پسماند»^۱ یا «انحراف»^۲ گفته و با e_i نشان می‌دهیم،

$$e_i = Y_i - \hat{Y}_i. \quad (1-18)$$

اصطلاح «خطای تخمین»^۳ نیز در مواردی برای این مفهوم به کار رفته است. ناگفته نماند که e_i را با \hat{U}_i نیز می‌توان نشان داد.

معیار روش حداقل مربعات معمولی این است که α و β را باید چنان تخمین زد که مجموع مربعات پسماند به حداقل برسد؛ به عبارت دیگر باید $\hat{\alpha}$ و $\hat{\beta}$ را چنان محاسبه کرد که $\sum e_i^2$ به حداقل برسد. علت اینکه به جای مجموع ساده پسماند، مجموع مربع آنها را در نظر می‌گیریم، این است که می‌خواهیم نه تنها انحرافات مثبت و منفی یکدیگر را خنثی نکنند، بلکه انحرافهای بزرگ نسبت به انحرافهای کوچک از اهمیت بیشتری برخوردار شود. می‌دانیم، $\sum e_i^2$ برابر است با

$$\begin{aligned} \sum e_i^2 &= \sum (Y_i - \hat{Y}_i)^2, \\ &= \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2. \end{aligned} \quad (1-19)$$

بنابراین مسأله تخمین α و β به مسأله حداقل سازی رابطه ۱-۱۹ نسبت به $\hat{\alpha}$ و $\hat{\beta}$ تبدیل می‌شود. تخمین پارامترها با روش حداقل مربعات معمولی (OLS) همواره می‌تواند به حداقل سازی یک تابع تبدیل شود. بدین ترتیب باید $\sum e_i^2$ را در رابطه ۱-۱۹

نسبت به $\hat{\alpha}$ و $\hat{\beta}$ حداقل کنیم. کافی است از $\sum e_i^2$ نسبت به $\hat{\alpha}$ و $\hat{\beta}$ مشتق بگیریم و آنها را مساوی صفر قرار داده و $\hat{\alpha}$ و $\hat{\beta}$ مطلوب را محاسبه کنیم.

$$\begin{aligned} \frac{\partial \sum e_i^2}{\partial \hat{\alpha}} &= \frac{\partial \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2}{\partial \hat{\alpha}} = 0, \\ &= 2 \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) (-1) = 0, \\ &= \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0. \end{aligned}$$

$$\begin{aligned} \frac{\partial \sum e_i^2}{\partial \hat{\beta}} &= \frac{\partial \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2}{\partial \hat{\beta}} = 0, \\ &= 2 \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) (-X_i) = 0, \\ &= \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) (X_i) = 0. \end{aligned}$$

نتایج این مشتق‌گیریهای جزئی را «معادله‌های اول و دوم نرمال»^۱ می‌گویند. این دو معادله را به صورت زیر می‌نویسیم،

$$\begin{cases} \sum Y_i = \sum \hat{\alpha} + \sum \hat{\beta} X_i, \\ \sum X_i Y_i = \sum \hat{\alpha} X_i + \sum \hat{\beta} X_i^2, \end{cases}$$

با توجه به اینکه n مشاهده داریم و مقادیر $\hat{\alpha}$ و $\hat{\beta}$ نیز ثابت است، خواهیم داشت

$$\begin{cases} \sum Y_i = n \hat{\alpha} + \hat{\beta} \sum X_i, \\ \sum X_i Y_i = \hat{\alpha} \sum X_i + \hat{\beta} \sum X_i^2. \end{cases} \quad (1-20)$$

از حل دستگاه دو معادله دو مجهولی ۱-۲۰ می‌توان $\hat{\alpha}$ و $\hat{\beta}$ را به دست آورد،

$$\hat{\alpha}_{OLS} = \frac{\sum X_i^T \sum Y_i - \sum X_i \sum X_i Y_i}{n \sum X_i^T - (\sum X_i)^T}, \quad (1.21)$$

$$\hat{\beta}_{OLS} = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n \sum X_i^T - (\sum X_i)^T}. \quad (1.22)$$

منظور از $\hat{\alpha}_{OLS}$ و $\hat{\beta}_{OLS}$ در فرمولهای فوق این است که تخمینهای α و β با استفاده از روش حداقل مربعات معمولی به دست آمده است. به فرمولهای ۱-۲۱ و ۱-۲۲ که در واقع فرمولهای تخمین هستند، «تخمین زنده»^۱ می‌گویند.

می‌توان به فرمولهای محاسباتی ساده‌تری برای $\hat{\alpha}_{OLS}$ و $\hat{\beta}_{OLS}$ رسید. اما قبل از آن دو نتیجه بسیار مهم را که از معادله‌های نرمال حاصل می‌شوند بررسی می‌کنیم. معادله اول نرمال را یک بار دیگر می‌نویسیم،

$$\frac{\partial \sum e_i^2}{\partial \hat{\alpha}} = \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0.$$

با توجه به معادله ۱-۱۸ می‌توان از معادله فوق نتیجه گرفت که مجموع پسماندها و میانگین آنها صفر است،

$$\bar{e} = 0 \quad \text{یا} \quad \sum e_i = 0. \quad (1.23)$$

معادله دوم نرمال را نیز یک بار دیگر می‌نویسیم،

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}} = \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) (X_i) = 0,$$

با توجه به معادله ۱-۱۸ می‌توان گفت که مجموع حاصلضرب پسماندها در متغیر توضیحی برابر صفر است،

$$\sum e_i X_i = 0.$$

همچنین می‌توان نشان داد که $\sum e_i X_i = \sum e_i x_i = 0$ ؛ زیرا با توجه به معادله ۱-۲۳

داریم:

$$\sum e_i x_i = \sum e_i (x_i + \bar{X}) = \sum e_i x_i - \bar{X} \sum e_i = \sum e_i x_i = 0.$$

بنابراین جمله پسماند و متغیر توضیحی از یکدیگر مستقل هستند؛ زیرا

$$\sum (e_i - \bar{e})(X_i - \bar{X}) = \sum e_i x_i = 0. \quad (1.24)$$

با استفاده از تعاریف $x_i = X_i - \bar{X}$ و $y_i = Y_i - \bar{Y}$ که در آن \bar{X} و \bar{Y} به ترتیب میانگین X_i و Y_i هستند، می‌توان به فرمولهای ساده‌تری برای $\hat{\alpha}$ و $\hat{\beta}$ نیز رسید. می‌دانیم،

$$\sum x_i y_i = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n} \quad \text{و} \quad \sum x_i^2 = \frac{n \sum X_i^2 - (\sum X_i)^2}{n}.$$

اگر صورت و مخرج تخمین‌زننده $\hat{\beta}_{OLS}$ در رابطه ۱.۲۲ را بر n تقسیم کرده، از رابطه‌های فوق استفاده کنیم، خواهیم داشت

$$\hat{\beta}_{OLS} = \frac{\sum x_i y_i}{\sum x_i^2}. \quad (1.25)$$

دو طرف معادله اول از سیستم معادله‌های ۱.۲۰ را بر n تقسیم می‌کنیم،

$$\frac{\sum Y_i}{n} = \hat{\alpha} + \hat{\beta} \frac{\sum X_i}{n}, \quad \text{یا}$$

$$\bar{Y} = \hat{\alpha} + \hat{\beta} \bar{X}. \quad (1.26)$$

رابطه فوق بر این دلالت می‌کند که مختصات \bar{Y} و \bar{X} در تخمین مدل رگرسیون (معادله ۱-۱۶)، صدق می‌کند؛ به عبارت دیگر می‌توان گفت که مدل $\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i$ از نقطه میانگین X_i و میانگین Y_i می‌گذرد. به سهولت ملاحظه می‌شود که

$$\hat{\alpha}_{OLS} = \bar{Y} - \hat{\beta} \bar{X}.$$

برای تبیین کیفیت استفاده از تخمین‌زننده‌های $\hat{\beta}_{OLS}$ و $\hat{\alpha}_{OLS}$ در حل مسائل،

به ذکر چند مثال می‌پردازیم. (در ادامه بحث هر گاه تخمین زنده‌ای بدون اندیس نوشته شود، منظور تخمین زنده حداقل مربعات معمولی (OLS) است).

مثال ۱-۱ مدل رگرسیون مصرف زیر مفروض است،

$$C_t = \alpha + \beta Y_t + U_t.$$

با استفاده از آمار موجود در جدول ۱-۱ پارامترهای α و β را با استفاده از روش حداقل مربعات معمولی تخمین بزنید. مقادیر پسماند را نیز حساب کنید و مجموع مربع آنها را به دست آورید. اگر در دوره $t = 6$ مقدار درآمد ۱۲ میلیارد تومان باشد، سطح مصرف را پیش‌بینی کنید.

بهرتر است نام متغیرهای موجود در مسأله را به X_t و Y_t تغییر دهیم تا بهتر بتوان از فرمولها یا تخمین زنده‌های ۱-۲۵ و ۱-۲۶ استفاده کرد.

$$Y_t = \text{مصرف} \quad \text{و} \quad X_t = \text{درآمد}$$

محاسبات لازم را در جدولی به شرح زیر تنظیم می‌کنیم.

جدول ۱-۲

t	Y_t	X_t	$y_t = Y_t - \bar{Y}$	$x_t = X_t - \bar{X}$	$x_t y_t$	x_t^2	\hat{Y}_t	$e_t = Y_t - \hat{Y}_t$	e_t^2
۱	۲	۳	۲ - ۵ = -۳	۳ - ۸ = -۵	۱۵	۲۵	۱/۶۵	۲ - ۱/۶۵ = ۰/۳۵	۰/۱۲۲۵
۲	۳	۵	۳ - ۵ = -۲	۵ - ۸ = -۳	۶	۹	۲/۹۹	۳ - ۲/۹۹ = ۰/۰۱	۰/۰۰۰۱
۳	۵	۹	۵ - ۵ = ۰	۹ - ۸ = ۱	۰	۱	۵/۶۷	۵ - ۵/۶۷ = ۰/۶۷	۰/۴۴۸۹
۴	۶	۱۰	۶ - ۵ = ۱	۱۰ - ۸ = ۲	۲	۴	۶/۳۴	۶ - ۶/۳۴ = ۰/۳۴	۰/۱۱۵۶
۵	۹	۱۳	۹ - ۵ = ۴	۱۳ - ۸ = ۵	۲۰	۲۵	۸/۳۵	۹ - ۸/۳۵ = ۰/۶۵	۰/۴۲۲۵
Σ	۲۵	۴۰	.	.	۴۰	۶۴	۲۵	.	۱/۰۹۶

$$\hat{\beta} = \frac{\Sigma x_t y_t}{\Sigma x_t^2}, \quad \bar{Y} = \frac{\Sigma Y_t}{n} = \frac{۲۵}{۵} = ۵, \quad \bar{X} = \frac{\Sigma X_t}{n} = \frac{۴۰}{۵} = ۸,$$

$$= \frac{۴۳}{۶۴} = ۰/۶۷,$$

مفاهیم و تخمین مدل رگرسیون ... ۳۵

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X},$$

$$= 0 - 0/67(8) = -0/36,$$

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i,$$

$$\hat{Y}_i = -0/36 + 0/67 X_i.$$

بر اساس نامگذاریهای موجود در مسأله داریم

$$\hat{C}_i = -0/36 + 0/67 Y_i.$$

میل نهایی به مصرف، $mpc = 0/67$ است. تخمینهایی که از پارامترها زده‌ایم، بر این دلالت می‌کند که اگر درآمد صفر شود، مصرف برابر $0/36$ میلیارد تومان است. بدیهی است چنین نتیجه‌ای مطلقاً بی‌معنی است. بنابراین، محاسبات متغیر درون‌زا (متغیر موجود در سمت چپ مدل رگرسیون) فقط باید به ازای قلمرو محدودی از تغییرات متغیر برون‌زا صورت گیرد و در خارج از این قلمرو معمولاً نتایج، بی‌معنی خواهد بود.

با توجه به $\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i$ ، مقادیر \hat{Y}_i را به ترتیب زیر محاسبه می‌کنیم.

$$\hat{Y}_1 = -0/36 + 0/67(3) = 1/60,$$

$$\hat{Y}_2 = -0/36 + 0/67(5) = 2/99,$$

$$\hat{Y}_3 = -0/36 + 0/67(9) = 0/67,$$

$$\hat{Y}_4 = -0/36 + 0/67(10) = 7/34,$$

$$\hat{Y}_5 = -0/36 + 0/67(13) = 8/35.$$

با استفاده از رابطه $e_i = Y_i - \hat{Y}_i$ ، به راحتی می‌توان پسماندها را برای هر یک از

دوره‌های مورد مطالعه محاسبه کرد. این مقادیر در ستون ۹ جدول منعکس شده است. مجموع مربعات پسماند $(\sum e_i^2 = 1/1096)$ نیز به راحتی قابل محاسبه است. مقدار مجموع مربعات پسماند به این معنی است که با روش حداقل مربعات معمولی مقادیر $\hat{\alpha}$ و $\hat{\beta}$ چنان تعیین شده است که مجموع مربعات پسماند کمترین مقدار ممکن را دارد و این کمترین مقدار برابر $1/1096$ میلیارد تومان شده است. برای پیش‌بینی مقدار مصرف در $t=6$ کافی است که در مدل

$$\hat{C}_t = -0.36 + 0.67 Y_t,$$

به جای Y_t ، مقدار $Y_t = 12$ را قرار داده، در آن صورت \hat{C}_t برابر $7/68$ محاسبه می‌شود. همچنین می‌توان نشان داد که $\sum e_i X_i$ یا $\sum e_i x_i$ تقریباً برابر صفر است.

مثال ۱-۲ برای بررسی رابطه تولید با ساعتهای نیروی کار از مدل رگرسیون زیر استفاده کرده‌ایم،

$$Q_i = \alpha + \beta L_i + U_i,$$

که در آن Q_i تولید و L_i ساعت کار است.

برای تخمین α و β با استفاده از روش حداقل مربعات معمولی این مشاهدات را داریم

جدول ۱-۳

Q_i	۱۱	۱۰	۱۲	۶	۱۰	۷	۹	۱۰	۱۱	۱۰
L_i	۱۰	۷	۱۰	۵	۸	۸	۶	۷	۹	۱۰

$\hat{\alpha}_{OLS}$ و $\hat{\beta}_{OLS}$ را بدون استفاده از محاسبات انحراف از میانگین به دست آورید.

در مسأله قبل با استفاده از انحراف از میانگین، یعنی $x_i = X_i - \bar{X}$ و $y_i = Y_i - \bar{Y}$

پارامترها را تخمین زدیم. می‌توان به طور مستقیم و به کمک مشاهدات اصلی نیز پارامترها را تخمین زد. می‌دانیم،

$$\sum x_i y_i = \sum (X_i - \bar{X})(Y_i - \bar{Y}) = \sum X_i Y_i - n \bar{X} \bar{Y},$$

$$\sum x_i^2 = \sum (X_i - \bar{X})^2 = \sum X_i^2 - n \bar{X}^2.$$

بنابراین:

$$\begin{aligned} \hat{\beta} &= \frac{\sum x_i y_i}{\sum x_i^2}, \\ &= \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sum X_i^2 - n \bar{X}^2}. \end{aligned}$$

با تغییر در نامگذاری متغیرها خواهیم داشت:

$$X_i = \text{ساعت کار} \quad \text{و} \quad Y_i = \text{تولید}$$

بنابراین مدل رگرسیون مفروض را می‌توان به صورت زیر نوشت،

$$Y_i = \alpha + \beta X_i + U_i.$$

با استفاده از محاسبات جدول ۱-۴ کمیت‌های لازم برای $\hat{\beta}_{OLS}$ را به دست می‌آوریم،

$$\begin{aligned} \sum x_i y_i &= \sum X_i Y_i - n \bar{X} \bar{Y}, \\ &= 789 - 10(8)(9/6) = 789 - 768 = 21, \end{aligned}$$

$$\begin{aligned} \sum x_i^2 &= \sum X_i^2 - n \bar{X}^2, \\ &= 668 - 10(8)^2 = 668 - 640 = 28. \end{aligned}$$

بنابراین

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{21}{28} = 0.75.$$

برای تخمین $\hat{\alpha}_{OLS}$ کافی است از رابطه زیر استفاده کنیم،

$$\begin{aligned} \hat{\alpha} &= \bar{Y} - \hat{\beta} \bar{X}, \\ &= 9/6 - 0.75(8) = 3/6. \end{aligned}$$

جدول ۱.۴

مشاهدات	Y_i	X_i	$X_i Y_i$	X_i^2
۱	۱۱	۱۰	۱۱۰	۱۰۰
۲	۱۰	۷	۷۰	۴۹
۳	۱۲	۱۰	۱۲۰	۱۰۰
۴	۶	۵	۳۰	۲۵
۵	۱۰	۸	۸۰	۶۴
۶	۷	۸	۵۶	۶۴
۷	۹	۶	۵۴	۳۶
۸	۱۰	۷	۷۰	۴۹
۹	۱۱	۹	۹۹	۸۱
۱۰	۱۰	۱۰	۱۰۰	۱۰۰
Σ	۹۶	۸۰	۷۸۹	۶۶۸

$$\bar{X} = \frac{\sum X_i}{n} = \frac{80}{10} = 8,$$

$$\bar{Y} = \frac{\sum Y_i}{n} = \frac{96}{10} = 9.6.$$

بدین ترتیب تخمین مدل رگرسیون مفروض عبارت است از

$$\hat{Q}_i = 3/6 + 0/70 L_i.$$

ملاحظه می شود که بهره وری نهایی نیروی کار^۱ برابر ۰/۷۰ است. تخمین مدل رگرسیون بر این دلالت می کند که اگر ساعت - کار، صفر باشد؛ یعنی اساساً هیچ کاری انجام نشود، تولید ۳/۶ واحد خواهد بود. علت رسیدن به چنین نتیجه بی معنایی است که مقدار متغیر درونزا را به ازای آن مقدار از متغیر برونزا حساب کرده ایم که فاصله بسیاری از قلمرو تغییرات متغیر برونزا دارد. در جدول ۱-۴ ملاحظه می کنیم که قلمرو تغییرات متغیر برونزا بین ۵ تا ۱۰ است. هر مقدار متغیر برونزا از این فاصله دورتر شود، محاسبه و پیش بینی مقدار متناظر متغیر برونزا معنای خود را از دست خواهد داد (اثبات ریاضی این نکته در مبحث پیش بینیها در ۳-۲ آمده است).

1. Marginal Productivity of Labour

۱-۴ ضریب تعیین

مدل رگرسیون ۱-۸ را یک بار دیگر می‌نویسیم،

$$Y_t = \alpha + \beta X_t + U_t.$$

می‌دانیم تخمین این مدل با استفاده از روش حداقل مربعات معمولی عبارت است از

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t.$$

سؤال این است که چگونه \hat{Y}_t توانسته است تغییرات Y_t را توضیح دهد؟ در واقع می‌دانیم یکی از اهداف اساسی در تحلیل‌های رگرسیونی این است که بتوان تغییرات متغیر درون‌زا را توضیح داد. متغیر درون‌زای ما در اینجا Y_t است. به کمک X_t و با استفاده از مدل ۱-۸ کوشش کرده‌ایم تغییرات Y_t را توضیح دهیم و سرانجام بعد از به دست آوردن $\hat{\alpha}$ و $\hat{\beta}$ ، مقادیر \hat{Y}_t به دست آمده است. بنابراین می‌توان در این قسمت این سؤال را مطرح کرد که مسیری از تغییرات Y_t - که ما تخمین زده‌ایم - (\hat{Y}_t)، چقدر با مسیر اصلی تغییرات Y_t تناسب دارد. برای این منظور کمیتی را محاسبه می‌کنیم که بتواند درجهٔ برازندگی یا تناسب \hat{Y}_t را نسبت به Y_t تعیین کند. به این کمیت «ضریب تعیین درجهٔ تناسب»^۱ تخمین مدل رگرسیون گفته می‌شود. از این به بعد، برای سهولت، این کمیت را «ضریب تعیین» می‌خوانیم. مسلماً هر چه تخمین ما از تغییرات (\hat{Y}_t) به Y_t نزدیکتر باشد، نشان‌دهندهٔ بهتر بودن تناسب تخمین با مقادیر واقعی است.

توجه به این نکته اهمیت دارد که وقتی از تغییرات یک متغیر صحبت می‌شود، باید بتوانیم کل تغییرات ملاحظه شده در یک دورهٔ زمانی را با یک کمیت واحد نشان دهیم؛ برای مثال، وقتی می‌گوییم تغییرات Y_t یا \hat{Y}_t در فاصلهٔ زمانی $n, \dots, 2, 1, t$ است، به دو کمیت واحد نیاز داریم تا از طریق مقایسه آنها یا یکدیگر بتوان نسبت به تناسب دو مسیر از تغییرات Y_t و \hat{Y}_t اظهار نظر کرد. راه‌های مختلفی برای بیان میزان تغییرات یک متغیر وجود دارد؛ یکی از بهترین روشها این است که ابتدا میانگین تغییرات

1. Coefficient of Determination of Goodness of Fit

را مبنای سنجش قرار دهیم و انحراف مشاهدات را از میانگین به دست آوریم، سپس مربع انحرافات را حساب کرده، مجموع آنها را به عنوان شاخص کل تغییرات معرفی کنیم. بدین ترتیب شاخص تغییرات Y_t و \hat{Y}_t به ترتیب عبارتند از

$$\sum (Y_t - \bar{Y})^2 = \sum y_t^2, \quad (1.27)$$

$$\sum (\hat{Y}_t - \bar{Y})^2 = \sum \hat{y}_t^2. \quad (1.28)$$

$\sum y_t^2$ در حقیقت کل تغییرات مشاهده شده Y_t است که معمولاً آن را با TSS یا TSSQ نشان می‌دهند. $\sum \hat{y}_t^2$ نیز تغییرات توضیح داده شده توسط تخمین مدل رگرسیون است که می‌توان آن را با ESS یا ESSQ نشان داد. نزدیکی هر چه بیشتر ESS به TSS مبین این واقعیت است که تناسب تخمین مدل رگرسیون با مدل واقعی رگرسیون بیشتر است. بنابراین، بهتر است کسری را در نظر بگیریم که صورت آن تغییرات توضیح داده شده و مخرج آن کل تغییرات باشد. این کسر همان ضریب تعیین خواهد بود،

$$\begin{aligned} \text{ضریب تعیین} &= \frac{\text{تغییرات توضیح داده شده}}{\text{کل تغییرات}}, \\ &= \frac{\text{ESS}}{\text{TSS}} = \frac{\sum (\hat{Y}_t - \bar{Y})^2}{\sum (Y_t - \bar{Y})^2} = \frac{\sum \hat{y}_t^2}{\sum y_t^2}. \end{aligned} \quad (1.29)$$

برای اینکه بتوان به فرمولهای ساده‌تری رسید، ابتدا دو نکته را مطرح می‌کنیم:
نکته اول: ثابت می‌کنیم که میانگین Y_t با میانگین \hat{Y}_t برابر است یعنی $\bar{Y} = \bar{\hat{Y}}$.
برای اثبات کافی است معادله اول نرمال ۱-۲۰ را یک بار دیگر بنویسیم،

$$\sum Y_t = n \hat{\alpha} + \hat{\beta} \sum X_t.$$

دو طرف را بر n تقسیم می‌کنیم،

$$\bar{Y} = \hat{\alpha} + \hat{\beta} \bar{X}.$$

سپس تخمین مدل رگرسیون، یعنی معادله ۱-۱۶ را نیز دوباره می‌نویسیم،

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i.$$

دو طرف رابطه فوق را برای $i = 1, 2, \dots, n$ جمع می‌کنیم،

$$\sum \hat{Y}_i = n \hat{\alpha} + \hat{\beta} \sum X_i.$$

و یا تقسیم دو طرف بر n خواهیم داشت

$$\bar{\hat{Y}} = \hat{\alpha} + \hat{\beta} \bar{X}.$$

پس نتیجه می‌گیریم که

$$\bar{Y} = \bar{\hat{Y}}. \quad (1.30)$$

نکته دوم: نشان می‌دهیم که $\hat{y}_i = \hat{\beta} \hat{x}_i$ است. اثبات این نکته بسیار ساده است.

مدل تخمین رگرسیون (معادله ۱-۱۶) را یک بار دیگر می‌نویسیم،

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i.$$

در نکته اول دیدیم که

$$\bar{Y} = \hat{\alpha} + \hat{\beta} \bar{X}.$$

کافی است دو رابطه فوق را از یکدیگر کم کنیم،

$$\hat{Y}_i - \bar{Y} = \hat{\beta} (X_i - \bar{X}).$$

با استفاده از رابطه ۱-۳۰، معادله فوق را می‌توان چنین نوشت،

$$\hat{Y}_i - \bar{Y} = \hat{\beta} (X_i - \bar{X}),$$

یا

$$\hat{y}_i = \hat{\beta} x_i. \quad (1.31)$$

بعد از بیان این دو نکته به بحث ضریب تعیین برمی‌گردیم. رابطه ۱-۲۹ را یک بار

دیگر می‌نویسیم،

$$\text{ضریب تعیین} = \frac{\sum \hat{y}_i^2}{\sum y_i^2}$$

کافی است به جای \hat{y}_i مقدار آن را از رابطه ۱-۳۱ قرار دهیم،

$$\text{ضریب تعیین} = \frac{\sum (\hat{\beta} x_i)^2}{\sum y_i^2} = \frac{\hat{\beta}^2 \sum x_i^2}{\sum y_i^2}$$

با توجه به رابطه ۱-۲۵ می‌دانیم

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

در نتیجه خواهیم داشت

$$\begin{aligned} \text{ضریب تعیین} &= \frac{(\sum x_i y_i)^2 \cdot \sum x_i^2}{(\sum x_i^2)^2 \cdot \sum y_i^2} \\ &= \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2} \end{aligned}$$

اما در آمار دیده‌ایم که ضریب همبستگی بین دو متغیر X_i و Y_i از رابطه زیر به دست می‌آید،

$$r_{xy} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \cdot \sum y_i^2}} \quad (1-32)$$

بنابراین، خواهیم داشت

$$\text{ضریب تعیین} = \frac{(\sum x_i y_i)^2}{\sum x_i^2 \cdot \sum y_i^2} = r^2$$

در واقع نشان دادیم که ضریب تعیین برابر مجذور ضریب همبستگی است.

به عبارت دیگر، r^2 که معمولاً در کتابهای اقتصادسنجی برای نشان دادن ضریب تعیین به کار می‌رود یک علامتگذاری یا قرارداد نیست، بلکه ثابت می‌شود که ضریب تعیین برابر r^2 است، بنابراین

$$r^2 = \frac{\sum \hat{y}_t^2}{\sum y_t^2} \quad (1.33)$$

$$r^2 = \frac{(\sum x_t y_t)^2}{\sum x_t^2 \cdot \sum y_t^2} \quad (1.34)$$

می‌توان فرمول سومی نیز برای r^2 استخراج کرده که از نظر محاسباتی بسیار ساده‌تر است. فرمول ۱.۳۴ را یک بار دیگر می‌نویسیم و آن را تفکیک می‌کنیم،

$$r^2 = \frac{(\sum x_t y_t)^2}{\sum x_t^2 \cdot \sum y_t^2} = \frac{\sum x_t y_t}{\sum x_t^2} \cdot \frac{\sum x_t y_t}{\sum y_t^2}$$

اما با توجه به رابطه ۱.۲۵ می‌دانیم $\hat{\beta} = \frac{\sum (x_t y_t)}{\sum x_t^2}$ در نتیجه خواهیم داشت

$$r^2 = \hat{\beta}^2 \frac{\sum x_t^2}{\sum y_t^2} \quad \text{یا} \quad r^2 = \hat{\beta} \frac{\sum x_t y_t}{\sum y_t^2} \quad (1.35)$$

برای رسیدن به فرمول دیگری برای r^2 که جنبه مفهومی بسیار خوبی دارد - ابتدا دو نکته را مطرح کرده، سپس یک قضیه را ثابت می‌کنیم.
نکته اول: با استفاده از نمودار ۱-۳ ملاحظه می‌شود که

$$Y_t = \hat{Y}_t + e_t \quad (1.36)$$

یعنی مقدار مشاهده شده متغیر درون‌زا برابر است با مقدار تخمین زده شده به علاوه مقدار پسماند.

این نکته را به بیان ساده‌تر می‌توان چنین نوشت.

$$\text{پسماند} + \text{تخمین} = \text{مشاهده}$$

از دو طرف رابطه ۱-۳۶ مقدار \bar{Y} را کم می‌کنیم،

$$(Y_t - \bar{Y}) = (\hat{Y}_t - \bar{Y}) + e_t.$$

با استفاده از رابطه ۱-۳۰ و جایگزینی $\bar{Y} = \bar{Y}$ در طرف راست آن خواهیم داشت

$$(Y_t - \bar{Y}) = (\bar{\hat{Y}}_t - \hat{Y}_t) + e_t,$$

یا:

$$y_t = \hat{y}_t + e_t. \quad (1.37)$$

در نتیجه نشان دادیم، رابطه ۱-۳۶ که برای مشاهدات اصلی صادق است برای حالت تفاوت از میانگین نیز برقرار است.

نکته دوم: در معادله ۱-۲۳ دیدیم که می‌توان به کمک معادلات نرمال نشان داد که $\bar{e} = 0$ است. در اینجا می‌خواهیم با روش دیگر به همین نتیجه برسیم. از تعریف میانگین

e_t شروع می‌کنیم،

$$\bar{e} = \frac{\sum e_t}{n}.$$

با جایگزینی رابطه ۱-۱۸، یعنی $e_t = Y_t - \hat{Y}_t$ خواهیم داشت

$$\begin{aligned} \bar{e} &= \frac{\sum (Y_t - \hat{Y}_t)}{n} = \frac{\sum Y_t}{n} - \frac{\sum \hat{Y}_t}{n}, \\ &= \bar{Y} - \bar{\hat{Y}}. \end{aligned}$$

با توجه به رابطه ۱-۳۰، یعنی $\bar{Y} = \bar{\hat{Y}}$ به سهولت ملاحظه می‌شود که

$$\bar{e} = 0. \quad (1.38)$$

قضیه در مدل رگرسیون $Y_t = \alpha + \beta X_t + U_t$ ، اگر پارامترها را با روش حداقل مربعات معمولی تخمین بزنیم، رابطه زیر برقرار است.

تغییرات توضیح داده نشده + تغییرات توضیح داده شده = کل تغییرات متغیر درون‌زا

با توجه به اینکه مقادیر توضیح داده نشده، همان پسماندهاست؛ تغییرات توضیح داده

نشده نیز برابر تغییرات پسماندها خواهد بود. اگر تغییرات پسماندها را نیز به صورت مجموع مربعات انحراف از میانگین پسماندها تعریف کنیم، خواهیم داشت

$$\text{تغییرات پسماندها} = \sum e_i^2 = \sum (e_i - \bar{e})^2.$$

رابطه ۱-۳۸ را در رابطه فوق جایگزین می‌کنیم،

$$\text{تغییرات پسماندها} = \sum e_i^2. \quad (1-39)$$

با توجه به تعریف $e_i = Y_i - \hat{Y}_i$ ، می‌توان $\sum e_i^2$ را مجموع مربعات پسماندها (RSS) نامید. با توجه به رابطه‌های ۱-۲۷، ۱-۲۸ و ۱-۳۹، قضیه فوق را می‌توان چنین نوشت،

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2, \quad (1-40)$$

یا

$$TSS = ESS + RSS.$$

اثبات این قضیه بسیار ساده است. رابطه ۱-۳۷ را دوباره می‌نویسیم،

$$y_i = \hat{y}_i + e_i.$$

دو طرف آن را مجذور کرده و برای تمام مشاهدات $i = 1, 2, \dots, n$ جمع می‌کنیم،

$$y_i^2 = (\hat{y}_i + e_i)^2,$$

$$\sum y_i^2 = \sum (\hat{y}_i^2 + e_i^2).$$

طرف راست را بسط می‌دهیم،

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2 + 2 \sum e_i \hat{y}_i.$$

برای اثبات رابطه ۱-۴۰ کافی است ثابت کنیم $\sum e_i \hat{y}_i$ برابر صفر است. ابتدا معادله

۱-۳۱ را دوباره می‌نویسیم،

$$\hat{y}_i = \hat{\beta} x_i.$$

دو طرف را در e_i ضرب کرده و نسبت به تمام مشاهدات جمع می‌کنیم،

$$e_i \hat{y}_i = \hat{\beta} x_i e_i,$$

$$\sum e_i \hat{y}_i = \hat{\beta} \sum x_i e_i. \quad (1-41)$$

با استفاده از رابطه‌های ۱-۳۷ و ۱-۳۱ داریم

$$e_i = y_i - \hat{y}_i,$$

$$= y_i - \hat{\beta} x_i,$$

و با جایگزینی در طرف راست رابطه ۱-۴۱ خواهیم داشت

$$\sum e_i \hat{y}_i = \hat{\beta} \sum x_i (y_i - \hat{\beta} x_i),$$

$$= \hat{\beta} [\sum x_i y_i - \hat{\beta} \sum x_i^2].$$

مقدار $\hat{\beta}$ از رابطه ۱-۲۵ را در این رابطه جایگزین می‌کنیم،

$$\sum e_i \hat{y}_i = \hat{\beta} \left[\sum x_i y_i - \frac{\sum x_i y_i}{\sum x_i^2} \cdot \sum x_i^2 \right],$$

در نتیجه:

$$\sum e_i \hat{y}_i = 0. \quad (1-42)$$

با استفاده از معادله ۱-۲۴ نیز به طور مستقیم می‌توان صفر بودن $\sum e_i \hat{y}_i$ را از معادله ۱-۴۱ نتیجه گرفت. بدیهی است با اثبات صفر شدن $\sum e_i \hat{y}_i$ ، رابطه ۱-۴۰ ثابت می‌شود.^۱

۱. در فصل دوم در مبحث آنالیز واریانس، به کمک معادلات ۲-۴۵ و ۲-۴۶ اثبات دیگری برای این قضیه ارائه شده است.

حال که این قضیه بسیار مهم به اثبات رسید به استخراج یک فرمول دیگر برای r^2 می پردازیم. رابطه ۱-۴۰ را به صورت زیر می نویسیم،

$$\sum \hat{y}_i^2 = \sum y_i^2 - \sum e_i^2,$$

و دو طرف آن را بر $\sum y_i^2$ تقسیم می کنیم،

$$\frac{\sum \hat{y}_i^2}{\sum y_i^2} = \frac{\sum y_i^2}{\sum y_i^2} - \frac{\sum e_i^2}{\sum y_i^2}.$$

با استفاده از رابطه ۱-۳۳ خواهیم داشت

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2} = 1 - \frac{RSS}{TSS}. \quad (1-43)$$

فرمول ۱-۴۳ به دو دلیل بسیار مهم است:

اولاً، این فرمول می تواند قلمرو تغییرات r^2 را مشخص کند. در واقع حد بالای r^2 موقعی است که مطلقاً پسماندی در تخمین نداشته باشیم و آنچه تخمین زده ایم دقیقاً برابر مشاهدات باشد. در این حالت $\sum e_i^2$ برابر صفر و r^2 مساوی یک می شود. حد پایین r^2 هنگامی است که تخمین مدل رگرسیون ما مطلقاً قدرت توضیحی ندارد؛ بنابراین پسماندها به قدری بزرگند که برابر مشاهدات می باشند. در چنین حالتی $\sum e_i^2$ برابر $\sum y_i^2$ و مقدار r^2 برابر صفر خواهد بود. بنابراین قلمرو تغییرات r^2 را می توان به صورت زیر خلاصه کرد،

$$0 \leq r^2 \leq 1. \quad (1-44)$$

ثانیاً، با استفاده از این فرمول به راحتی می توان مجموع مربعات پسماند را محاسبه کرد. یادآوری می کنیم که فرمول ۱-۴۳ فرمول مناسبی برای محاسبه r^2 نیست؛ زیرا مستلزم محاسبه جمله های پسماند است. از بین فرمولهایی که تا اینجا برای محاسبه r^2 مطرح کردیم، به نظر می رسد فرمول ۱-۳۵ مناسبتر باشد،

$$r^2 = \hat{\beta} \frac{\sum x_i y_i}{\sum y_i^2}.$$

دلیل این امر کاملاً روشن است. می‌دانیم r^1 را معمولاً بعد از تخمین پارامترهای یک مدل محاسبه می‌کنند؛ بنابراین در موقع محاسبه r^2 مقدار $\hat{\beta}$ مشخص است. برای تخمین $\hat{\beta}$ نیز مقدار $\sum x_i y_i$ باید ضرورتاً محاسبه شود. بنابراین وقتی شروع به محاسبه r^2 می‌کنیم، $\hat{\beta}$ و $\sum x_i y_i$ از قبل محاسبه شده است. فقط کافی است $\sum y_i^2$ محاسبه شود تا r^2 به دست آید. حال اهمیت فرمول ۱-۴۳،

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2},$$

در این است که اگر r^1 را از فرمول ۱-۳۵ حساب کنیم و در آن قرار دهیم، e_i^2 را به راحتی محاسبه می‌شود. محاسبه e_i^2 با روشهای دیگر، ممکن است مستلزم تقریب و محاسبات طولانی در مورد هر یک از پسماندها باشد؛ در حالی که با استفاده از فرمول فوق مقدار e_i^2 به طور مستقیم قابل محاسبه خواهد بود. در فصل دوم نه تنها فرمولهای محاسباتی دقیقتری مانند ۲-۴۶ را برای $\sum e_i^2$ به دست خواهیم آورد؛ بلکه نشان خواهیم داد که محاسبه e_i^2 را قدم اساسی در تخمین واریانس پارامترهاست.

برای تبیین نحوه محاسبه ضریب تعیین چند مثال می‌آوریم:

مثال ۱-۳ مدل مصرف. موضوع مثال ۱-۱ را در نظر بگیرید. با استفاده از جدول ۱-۲ ضریب تعیین r^1 را برای این مدل محاسبه کنید و مجموع مربعات پسماند، یعنی $\sum e_i^2$ را به دست آورید.

با استفاده از فرمول ۱-۳۵ داریم

$$r^1 = \hat{\beta} \frac{\sum x_i y_i}{\sum y_i^2}.$$

در مثال ۱-۱ دیدیم که $\hat{\beta} = 0/67$ و $\sum x_i y_i = 43$. بنابراین باید فقط $\sum y_i^2$ محاسبه شود.

با استفاده از جدول ۱-۲ خواهیم داشت

$$\sum y_i^2 = 9 + 4 + 1 + 16 = 30,$$

به این ترتیب

$$r^2 = \frac{0.67(43)}{30} = \frac{28.81}{30} = 0.96,$$

یعنی متغیر درآمد، با توجه به مدل رگرسیون خطی مصرف، ۹۶ درصد تغییرات مصرف را توضیح داده است.

برای محاسبه مجموع مربعات پسماند، می‌توان از فرمول ۱-۴۳ استفاده کرد،

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2},$$

$$\frac{0.67(43)}{30} = 1 - \frac{\sum e_i^2}{30}, \quad 28.81 - 30 = -\sum e_i^2,$$

$$\sum e_i^2 = 1.19.$$

مثال ۱-۴ مدل رابطه تولید با ساعتهای نیروی کار. موضوع مثال ۱-۲ را ملاحظه کرده، با استفاده از جدول ۱-۴ ضریب تعیین r^2 را برای این مدل محاسبه کنید و مجموع مربعات پسماند $\sum e_i^2$ را به دست آورید. با توجه به معادله ۱-۳۵ می‌دانیم

$$r^2 = \hat{\beta}^2 \frac{\sum x_i y_i}{\sum y_i^2}.$$

در مثال ۱-۲ دیدیم که $\hat{\beta} = 0.75$ و $\sum x_i y_i = 21$ ؛ بنابراین باید فقط $\sum y_i^2$ محاسبه شود. با توجه به اینکه محاسبات در جدول ۱-۴ بر اساس مشاهدات واقعی و نه انحراف از میانگین ارائه شده است، ابتدا $\sum Y_i^2$ را محاسبه می‌کنیم،

$$\sum Y_i^2 = (11)^2 + (10)^2 + (12)^2 + (6)^2 + (10)^2 + (7)^2 + (9)^2 + (10)^2 + (11)^2 + (10)^2 = 952.$$

$\sum y_i^2$ به راحتی و با استفاده از $\sum Y_i^2$ به دست می‌آید. می‌دانیم

$$\sum y_i^2 = \sum Y_i^2 - n(\bar{Y})^2,$$

$$= 952 - 10(9/6)^2 = 952 - 921/6 = 30/4,$$

بدین ترتیب r^2 محاسبه می شود،

$$r^2 = \hat{\beta} \frac{\sum x_i y_i}{\sum y_i^2},$$

$$= \frac{0.70(21)}{30/4} = 0.92,$$

یعنی متغیر برونزا (ساعاتی کار)، فقط ۹۲ درصد از تغییرات متغیر درونزا (تولید) را توضیح داده است.

به عنوان تمرین، r^2 را از فرمول ضریب همبستگی نیز به دست می آوریم. با توجه به معادله ۱.۳۲ داریم

$$r_{xy} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \cdot \sum y_i^2}},$$

$$= \frac{21}{\sqrt{28(30/4)}} = \frac{21}{29} = 0.724,$$

یا

$$r^2 = 0.92.$$

برای محاسبه $\sum e_i^2$ ، از فرمول ۱.۴۳ استفاده می کنیم.

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2},$$

$$\frac{0.92(21)}{30/4} = 1 - \frac{\sum e_i^2}{30/4},$$

$$\sum e_i^2 = 14/60.$$

۱.۵ خطای معیار تخمین

در این قسمت به بررسی روش دیگری برای ارزیابی تخمین یک مدل رگرسیون

می پردازیم. - همان گونه که اشاره خواهد شد - با اینکه این روش جامعیت معیار ضریب تعیین را ندارد، در مواردی می تواند بسیار مفید باشد.

سؤالی را که در قسمت ۱-۴ مطرح کردیم، یک بار دیگر مطرح می کنیم. مدل رگرسیون مفروض، عبارت است از

$$Y_i = \alpha + \beta X_i + U_i ,$$

که با روش حداقل مربعات معمولی آن را به صورت زیر تخمین زده ایم،

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i .$$

می دانیم مبنای روش حداقل مربعات معمولی این است که مجموع مربعات پسماند باید حداقل شود؛ پس این سؤال که آیا \hat{Y}_i تخمین مناسبی از Y_i است یا خیر، می تواند به این سؤال تبدیل شود که آیا مقدار عددی مجموع مربعات پسماند رضایت بخش است یا خیر؟ به عبارت دیگر $\sum e_i^2$ می تواند بالقوه اساس ارزیابی ما در تخمین یک مدل رگرسیون باشد. حال به بررسی جواب این مسأله می پردازیم.

مهمترین اشکالی که به $\sum e_i^2$ به منزله شاخص ارزیابی به تخمین مدل رگرسیون وارد می شود، این است که مقدار آن به طور مستقیم، تابعی از حجم مشاهدات است. با افزایش مشاهدات، $\sum e_i^2$ به طور مرتب زیاد می شود و اساساً حدی ندارد. این اشکال را می توان بدین گونه رفع کرد که مجموع مربعات پسماند را بر تعداد مشاهدات تقسیم کرده و نتیجه را به عنوان شاخصی در ارزیابی خوبی تخمین مدل رگرسیون معرفی کنیم،

$$\text{معیار ارزیابی تخمین مدل رگرسیون} = \frac{\sum e_i^2}{n} = \frac{RSS}{n} ,$$

که در آن n حجم مشاهدات است.

با توجه به مباحث آماری، می توان گفت که اگر $\sum e_i^2$ را به جای اینکه بر n تقسیم کنیم، بر درجات آزادی آن تقسیم نماییم، معیار بهتری خواهیم داشت. این نکته را در

فصل دوم در بحث واریانس جمله اختلال و در معادله ۲-۲۸ ثابت خواهیم کرد. به هر حال می‌دانیم که در محاسبه تغییرات توضیح داده نشده (RSS) یا $\sum e_i^2$ دو درجه آزادی را از دست می‌دهیم. یادآوری می‌شود که تعداد درجات آزادی که از دست می‌دهیم برابر با تعداد پارامترهایی است که باید برای محاسبه مجموع مربعات پسماند تخمین بزنیم. واضح است که برای محاسبه تغییرات توضیح داده نشده باید α و β تخمین زده شود؛ بنابراین می‌گوییم که در محاسبه تغییرات توضیح داده نشده دو درجه آزادی را از دست می‌دهیم. اگر درجات آزادی را - که از دست داده‌ایم - از تعداد کل مشاهدات کم کنیم، آنگاه به درجات آزادی باقیمانده یا به طور خلاصه درجات آزادی می‌رسیم. بنابراین معیار بهتر، عبارت خواهد بود از

$$\text{معیار ارزیابی تخمین مدل رگرسیون} = \frac{\sum e_i^2}{n-2} = \frac{RSS}{n-2}$$

با وجود این، هنوز می‌توان اشکال دیگری به این معیار وارد کرد. با این معیار می‌خواهیم توانایی \hat{Y}_i در تفسیر تغییرات Y_i را اندازه‌گیری کنیم، در حالی که بُعد این معیار با بُعد Y_i متفاوت است؛ به عبارت دیگر، از رتبه دو است در حالی که Y_i از رتبه یک است. بدیهی است این اشکال را نیز می‌توان به سهولت رفع کرد. کافی است از $\frac{\sum e_i^2}{n-2}$ جذر بگیریم تا هم‌رتبه Y_i یا \hat{Y}_i شود. به معیاری که به این ترتیب به دست می‌آید «خطای معیار تخمین»^۱ می‌گویند و آن را «SEE» نشان می‌دهند،

$$SEE = \sqrt{\frac{\sum e_i^2}{n-2}} \quad (1-45)$$

می‌توان برای خطای معیار تخمین (SEE) فرمول ساده‌تری نیز به دست آورد،

1. Dimension

2. Standard Error of Estimation

برای توضیح بیشتر این فرمول می‌توان به بحث واریانس جمله اختلال در قسمت ۲-۲ و فرمول ۲-۲۸ مراجعه کرد.

به گونه‌ای که با داشتن مشاهدات اصلی قابل محاسبه باشد. با استفاده از معادله‌های ۱-۳۷ و ۱-۳۱ و جایگزینی آنها در فرمول ۱-۴۵ داریم

$$\begin{aligned} \text{SEE} &= \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} \\ &= \sqrt{\frac{\sum (y_i - \hat{\beta} x_i)^2}{n-2}} \end{aligned}$$

می‌دانیم

$$\sum (y_i - \hat{\beta} x_i)^2 = \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2}$$

زیرا کافی است طرف چپ رابطه فوق را بسط داده، به جای $\hat{\beta}$ مقدار آن را از فرمول ۱-۲۵ قرار دهیم،

$$\begin{aligned} \sum (y_i - \hat{\beta} x_i)^2 &= \sum y_i^2 + \hat{\beta} \sum x_i^2 - 2\hat{\beta} \sum x_i y_i \\ &= \sum y_i^2 + \frac{(\sum x_i y_i)^2}{(\sum x_i^2)^2} \cdot \sum x_i^2 - 2 \frac{(\sum x_i y_i)^2}{\sum x_i^2} \\ &= \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} \end{aligned}$$

بدین ترتیب خطای معیار تخمین (SEE) عبارت خواهد بود از

$$\text{SEE} = \sqrt{\left[\sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} \right] / (n-2)} \quad (1-46)$$

قبل از اینکه به نکته‌های مثبت و منفی خطای معیار تخمین (SEE) بپردازیم، برای تبیین بیشتر فرمول ۱-۴۶ مثالی می‌زنیم.

مثال ۱-۵ مدل رابطه تولید با ساعتهای نیروی کار. موضوع مثال ۱-۲ را در نظر بگیرید. با استفاده از جدول ۱-۴ خطای معیار تخمین را حساب کنید.

می‌دانیم $\sum y_i^2 = 30/4$ ، $\sum x_i^2 = 28$ و $\sum x_i y_i = 21$ است، بنابراین با توجه به $n=10$ و با استفاده از فرمول ۱-۴۶ خواهیم داشت

$$SEE = \sqrt{\left[30/4 - \frac{(21)^2}{28} \right] / (10 - 2)} ,$$

$$= \sqrt{14/60/8} = 1/303 .$$

با استفاده از فرمول ۱-۴۵ نیز می‌توان خطای معیار تخمین (SEE) را حساب کرد. در مثال ۱-۴ دیدیم که $\sum e_i^2 = 14/60$ ؛ بنابراین مقدار خطای معیار تخمین به صورت زیر محاسبه می‌شود،

$$SEE = \sqrt{\frac{14/60}{10 - 2}} = 1/303 .$$

حال که با کیفیت محاسبه خطای معیار تخمین آشنا شدیم، به موارد مثبت و منفی آن اشاره می‌کنیم. برای مواردی که می‌خواهیم حالت‌های مختلف یک مدل رگرسیون مفروض را مقایسه کنیم، خطای معیار تخمین می‌تواند بسیار مفید باشد؛ برای مثال، اگر بخواهیم به این سؤال پاسخ دهیم که آیا در مدل تولید و نیروی کار مذکور می‌توان به جای ۱۰ مشاهده فقط از ۸ مشاهده استفاده کرد یا خیر، کافی است خطای معیار تخمین را یک بار برای ۱۰ مشاهده و بار دیگر برای ۸ مشاهده محاسبه می‌کنیم. اگر نتایج به دست آمده تفاوت چندانی با یکدیگر نداشته باشد، می‌توان نتیجه گرفت که ضرورتی برای ۱۰ مشاهده نیست و ۸ مشاهده کافی است. بدیهی است با معیار r^2 نیز می‌توان به این سؤال پاسخ داد. اما اشکال اساسی روش خطای معیار تخمین این است که نمی‌تواند معیار مناسبی برای مقایسه تخمین‌های دو یا چند مدل رگرسیون باشد؛ زیرا خطای معیار تخمین تابعی از مقیاس اندازه‌گیری متغیر درون‌زا است. به عبارت دیگر، اگر بخواهیم برای چند متغیر درون‌زای متفاوت، مدل‌های رگرسیون بسازیم، برای مثال، برای تابع مصرف و شاخص قیمت‌ها، بر همین اساس، جمله‌های پسماند، مقیاس‌های متفاوتی خواهد داشت و در نتیجه نمی‌توان به طور مستقیم آنها را با یکدیگر مقایسه کرد. در نتیجه به منظور رسیدن به

معیاری برای ارزیابی تخمین یک مدل رگرسیون باید سعی کنیم معیار ما تابعی از مقیاس اندازه گیری متغیرها نباشد و دیدیم که معیار ضریب تعیین (r^2) این خصوصیت مطلوب را دارد. بنابراین r^2 نسبت به خطای معیار تخمین مرجع است، هر چند در مواردی نیز می توان از خطای معیار تخمین استفاده کرد. کاربرد اساسی خطای معیار تخمین در محاسبه واریانس تخمین پارامترها و آزمون فرضیه ها موضوع فصل آینده است.

مسائل فصل اول

۱-۱ مدل رگرسیون خطی

$$Y_i = \alpha + \beta X_i + U_i ,$$

مفروض است. U_i جمله اختلال مدل می باشد.

۱. معمولاً در اقتصادسنجی چه دلایلی برای ضرورت وجود U_i در یک مدل رگرسیون مطرح می شود؟

۲. اگر بخواهیم α و β را با روش حداقل مربعات معمولی تخمین بزنیم، چه فرضهایی برای U_i لازم است؟

۳. فرض کنید X_i برای تمام مشاهدات $n, \dots, 3, 2, 1, t$ مقدار ثابتی دارد. تأثیر این فرض را بر تخمین پارامترهای α و β ارزیابی کنید.

۴. فرض $E(U_i) = 0$ را با استفاده از نمودار، دقیقاً توضیح دهید.

۵. نشان دهید که واریانس متغیر درون‌زا برابر واریانس جمله اختلال مدل است.

۶. امید ریاضی متغیر درون‌زا را به دست آورید.

۷. با استفاده از نمودار، مفاهیم جمله اختلال مدل و جمله پسماند را دقیقاً توضیح دهید و با یکدیگر مقایسه کنید. می دانیم

$$U_i = Y_i - E(Y_i) \quad , \quad e_i = Y_i - \hat{Y}_i ,$$

آیا می توان e_i را به صورت تخمینی از U_i به کار برد. چرا؟

۸. با استفاده از نمودار، فرض واریانس همسانی U_i را توضیح دهید.

۹. با استفاده از نمودار، فرض استقلال X_i از U_i را توضیح دهید.

۱۰. برای هر یک از X_i و Y_i ، e اصطلاح رایج در اقتصادسنجی را نام ببرید.

۱-۲ مدل رگرسیون خطی زیر مفروض است،

$$Y_i = \alpha + \beta X_i + U_i .$$

۱. نشان دهید $\bar{Y} = \bar{Y}$.

۲. نشان دهید $\hat{y}_i = \hat{\beta} x_i$ ، که y و x به ترتیب عبارتند از $y_i = Y_i - \bar{Y}$ و

$$x_i = X_i - \bar{X}$$

۳. نشان دهید $\bar{e} = 0$.

۴. ثابت کنید $\sum y_i' = \sum \hat{y}_i' + \sum e_i'$ است.

۵. ثابت کنید $\sum e_i \hat{y}_i = 0$.

که در آن، $e_i = y_i - \hat{y}_i$.

۶. ثابت کنید $\sum e_i x_i = 0$.

۷. چرا ضریب تعیین r^2 بین صفر و یک است؟

۸. چرا r^2 به خطای معیار تخمین (SEE) برتری دارد؟

۱-۳ مدل رگرسیون

$$y_i = \beta x_i + u_i ,$$

مفروض است. x_i و y_i به ترتیب انحراف X_i و Y_i از میانگین آنهاست. با استفاده از روش حداقل مربعات معمولی (OLS) پارامتر β را به طور مستقیم از این مدل تخمین بزنید. آیا $\hat{\beta}$ به دست آمده از این مدل با $\hat{\beta}$ به دست آمده از مدل

$$Y_i = \alpha + \beta X_i + U_i$$

متفاوت است؟ چرا؟

۱-۴ مدل رگرسیون

$$CF_i = \alpha + \beta DI_i + U_i ,$$

مفروض است. CF_i و DI_i به ترتیب هزینه مصرف مواد غذایی و درآمد قابل تصرف خانوار i است. $i = 1, 2, \dots, 10$. مشاهدات زیر را داریم

DI_i	۲۰	۳۰	۲۲	۴۰	۱۵	۱۲	۲۶	۲۸	۳۵	۴۳
CF_i	۷	۹	۸	۱۱	۵	۴	۸	۱۰	۹	۱۰

۱. با روش حداقل مربعات معمولی پارامترهای α و β را تخمین بزنید.
۲. $e_i = CF_i - \hat{CF}_i$ را برای تمام مقادیر i حساب کنید.
۳. $\sum e_i^2$ را به طور مستقیم محاسبه کنید.
۴. ضریب تعیین r^2 را حساب کنید.
۵. با استفاده از r^2 به دست آمده در بند ۴، مقدار e_i^2 را حساب کنید.
۶. خطای معیار تخمین را به دست آورید.
۷. نشان دهید که $\sum cf_i^2 = \sum \hat{cf}_i^2 + \sum e_i^2$
۸. اگر خانوار یازدهم، درآمد قابل تصرفی برابر ۵۰ هزار تومان داشته باشد، مصرف مواد غذایی آن خانوار را تخمین بزنید.
۹. اگر DI_i را برابر صفر فرض کنیم، مقدار مصرف چقدر خواهد بود؟ این نتیجه را چگونه توجیه می کنید؟

۱-۵ مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i .$$

مفروض است. برای تخمین پارامترهای این مدل، ۱۶ مشاهده روی X و Y داشته و این کمیتها را محاسبه کرده ایم،

$$\sum Y_i^2 = ۵۲۶ , \quad \sum X_i^2 = ۶۵۷ , \quad \sum X_i Y_i = ۴۹۲ .$$

$$\sum Y_i = ۶۴ , \quad \sum X_i = ۹۶ ,$$

۱. با روش حداقل مربعات معمولی پارامتر β را تخمین بزنید.
۲. r^2 را محاسبه کنید.
۳. خطای معیار تخمین را به دست آورید.

۱-۶ برای تخمین یک مدل رگرسیون خطی بین دو متغیر X و Y که Y درون زاست - تعداد ۲۷ مشاهده روی X و Y انجام داده‌ایم و این نتایج را به دست آورده‌ایم،

$$\bar{X} = 100, \quad \sum_{i=1}^{27} (X_i - \bar{X})^2 = 100, \quad \sum_{i=1}^{27} (X_i - \bar{X})(Y_i - \bar{Y}) = 200.$$

$$\bar{Y} = 150, \quad \sum_{i=1}^{27} (Y_i - \bar{Y})^2 = 500.$$

۱. با روش حداقل مربعات معمولی پارامترهای α و β را تخمین بزنید.

۲. r^2 را محاسبه کنید.

۳. $\sum e_i^2$ را به دست آورید.

۴. خطای معیار تخمین را حساب کنید.

۱-۷ مدل زیر مفروض است،

$$Y_i = a_1 + b_1 X_i + U_i. \quad (1)$$

۱. فرض می‌شود جمله اختلال U_i ، تابعی از متغیر توضیحی X_i است، یعنی

$$U_i = a_2 + b_2 X_i + \varepsilon_i, \quad (2)$$

که ε_i جمله اختلال مدل (۲) بوده است و تمام فرضهای کلاسیک در مورد آن صادق است. همچنین فرض می‌کنیم که $b_2 > 0$ است. نشان دهید که پارامتر b_1 در مدل (۱) تأثیر X_i بر Y_i را کمتر از مقدار واقعی آن منعکس می‌کند.

۲. اگر معادله (۲) را به صورت زیر بنویسیم،

$$U_i = a_2 + b_2 X_i + \varepsilon_i, \quad (3)$$

در این صورت آیا همه فرضهای کلاسیک جمله اختلال در مدل (۱) به قوت خود باقی می‌ماند؟

۱-۸ مدل رگرسیون

$$Y_i = a + b X_i + U_i. \quad (1)$$

را در نظر می‌گیریم. فرض کنید در اندازه‌گیری متغیر توضیحی X_i خطای اندازه‌گیری داریم؛ یعنی به جای اینکه به طور مستقیم X_i را اندازه‌گیری کنیم، مقدار X_i^* را محاسبه کرده‌ایم، به گونه‌ای که

$$X_i^* = X_i + \varepsilon_i, \quad (2)$$

ε_i جملهٔ اختلال مدل (۲) و مستقل از X_i است و تمام فرضهای کلاسیک نیز در مورد آن صادق است. اگر فرض کنیم ε_i و U_i از یکدیگر مستقل باشند:

۱. نشان دهید که X_i^* از U_i مستقل است.

۲. یک مدل رگرسیون بسازید که Y_i را با X_i^* مرتبط سازد.

۳. در این مدل رگرسیون که ساخته‌اید آیا همهٔ فرضهای کلاسیک در مورد جملهٔ اختلال آن صادق است؟

۱-۹ مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i.$$

مفروض است. نشان دهید که هیچگونه تضمینی وجود ندارد که تخمینهای به دست آمده با روش حداقل مربعات ($\hat{\alpha}$ و $\hat{\beta}$)، دقیقاً با مقادیر واقعی پارامترها (α و β) برابر باشد.

۱-۱۰ در مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i,$$

۱. نشان دهید که اگر تمام مقادیر X_i با یکدیگر برابر باشند تخمین $\hat{\beta}$ ممکن نیست (بند ۳ مسأله ۱-۱).

۲. آیا می‌توان یک مدل رگرسیون داشت که علی‌رغم ثابت بودن مقادیر متغیر توضیحی، تخمین $\hat{\beta}$ ممکن باشد؟

۱-۱۱ برای مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i,$$

۱. نشان دهید که مجموع پسماندها، برابر صفر است؛ یعنی $\sum e_i = 0$.

۲. آیا به نظر شما می‌توان گفت $\sum U_i = 0$ ؟

۱-۱۲ مدل رگرسیون

$$Y_t = \alpha + \beta X_t + U_t .$$

مفروض است. فرض کنید مشاهدات مربوط به Y_t را به دو قسمت تقسیم کرده‌ایم: Y_{1t} و Y_{2t} و بدیهی است که $Y_t = Y_{1t} + Y_{2t}$. به کمک Y_{1t} و Y_{2t} دو مدل رگرسیون به شرح زیر می‌سازیم،

$$Y_{1t} = \alpha_1 + \beta_1 X_t + U_{1t} ,$$

$$Y_{2t} = \alpha_2 + \beta_2 X_t + U_{2t} .$$

نشان دهید که $\hat{\alpha} = \hat{\alpha}_1 + \hat{\alpha}_2$ و $\hat{\beta} = \hat{\beta}_1 + \hat{\beta}_2$.

۱-۱۳ متغیر تصادفی X_t مفروض است. این متغیر را به k قسمت تقسیم می‌کنیم، به گونه‌ای که برای تمام مقادیر t داشته باشیم

$$\sum X_{it} = X_t , \quad i = 1, 2, 3, \dots, k .$$

X_{it} را تابعی از X_t گرفته و مدل‌های زیر را می‌سازیم.

$$X_{1t} = \alpha_1 + \beta_1 X_t + U_{1t} ,$$

$$X_{2t} = \alpha_2 + \beta_2 X_t + U_{2t} ,$$

⋮

$$X_{kt} = \alpha_k + \beta_k X_t + U_{kt} ,$$

که در آن، $t = 1, 2, \dots, n$. نشان دهید $\sum_{i=1}^k \hat{\beta}_i = 1$ و $\sum_{i=1}^k \hat{\alpha}_i = 0$.

۱-۱۴ مدل رگرسیون

$$Y_t = \alpha + \beta X_t + U_t .$$

مفروض است. این مدل را تخمین زده و \hat{Y}_t را به دست می آوریم. نشان دهید که اگر مدل رگرسیون Y_t روی \hat{Y}_t را بسازیم،

$$Y_t = a + b \hat{Y}_t + \varepsilon_t ,$$

آنگاه $\hat{a} = 0$ ، $\hat{b} = 1$.

۱-۱۵ مدل رگرسیون

$$Y_t = \beta X_t + U_t ,$$

مفروض است. آیا تخمین این مدل با روش حداقل مربعات معمولی ($\hat{Y} = \hat{\beta} X_t$) از نقطه $(\bar{X}$ و $\bar{Y})$ می گذرد؟

۱-۱۶ می دانیم مدل واقعی برای تخمین پارامتر β عبارت است از:

$$Y_t = \beta X_t + U_t . \quad (1)$$

حال فرض کنید β را از این مدل تخمین می زنیم،

$$Y_t = \alpha + \beta X_t + U_t .$$

نشان دهید هنگامی تخمین β از مدل (۲) دقیقاً به همان نتیجه تخمین β از مدل (۱) می رسد که مدل (۲) را مقید تخمین یزنیم؛ یعنی شرط کرده باشیم که تخمین ضریب ثابت در مدل دوم صفر است.

۱-۱۷ مدل رگرسیون

$$Y_t = \beta X_t + U_t .$$

مفروض است. بدون استفاده از مشتق گیری نشان دهید که

$$\hat{\beta} = \frac{\sum X_t Y_t}{\sum X_t^2} .$$

۱-۱۸ مدل رگرسیون

$$Y_t = \alpha + \beta X_t + U_t .$$

مفروض است. مشاهدات X_t و Y_t را برحسب انحراف از مقادیر معلوم و اختیاری X^* و Y^* می نویسم،

$$Y_t^* = Y_t - Y^* \quad , \quad X_t^* = X_t - X^*$$

حال مدل

$$Y_t^* = a + b X_t^* + U_t \quad (۲)$$

را تخمین می زنیم. به نظر شما چه رابطه ای بین تخمین پارامترها در مدل اولیه و مدل جدید وجود دارد؟

۱-۱۹ در مدل رگرسیون

$$Y_t = \alpha + \beta X_t + U_t \quad ,$$

نشان دهید که ضریب تعیین r^2 را می توان به کمک هر یک از فرمولهای زیر محاسبه کرد.

$$۱) \quad r^2 = \frac{\sum (x_t y_t)^2}{\sum x_t^2 \sum y_t^2} \quad , \quad ۳) \quad r^2 = \frac{(\sum y_t \hat{y}_t)^2}{\sum y_t^2 \sum \hat{y}_t^2} \quad ,$$

$$۲) \quad r^2 = \hat{\beta} \frac{\sum x_t y_t}{\sum y_t^2} \quad , \quad ۴) \quad r^2 = 1 - \frac{\sum e_t^2}{\sum y_t^2} \quad .$$

۱-۲۰ در مدل رگرسیون

$$Y_t = \beta X_t + U_t \quad ,$$

نشان دهید که ضریب تعیین r^2 را می توان به کمک هر یک از فرمولهای زیر محاسبه کرد:

$$۱) \quad r^2 = \frac{(\sum X_t Y_t)^2}{\sum X_t^2 \sum Y_t^2} \quad , \quad ۳) \quad r^2 = 1 - \frac{\sum e_t^2}{\sum Y_t^2} \quad .$$

$$۲) \quad r^2 = \frac{\sum (Y_t \hat{Y}_t)^2}{\sum Y_t^2 \sum \hat{Y}_t^2} \quad ,$$

۱-۲۱ در مدل رگرسیون

$$Y_t = \beta X_t + U_t \quad ,$$

نشان دهید که اگر بخواهیم r^2 را با استفاده از فرمول

$$r^2 = 1 - \frac{RSS}{TSS} ,$$

محاسبه کنیم که در آن $RSS = \sum e_i^2$ و $TSS = \sum y_i^2$ ، چه بسا r^2 ممکن است منفی شود.
۱-۲۲ مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i . \quad (1)$$

مفروض است. از دو طرف مدل (۱) مقدار X_i را کم می‌کنیم:

$$Y_i - X_i = \alpha + (\beta - 1) X_i + U_i . \quad (2)$$

۱. نشان دهید که تخمین هر یک از معادلات (۱) و (۲) برای $\hat{\alpha}$ و $\hat{\beta}$ به نتایج

یکسانی می‌رسد.

۲. آیا تغییرات توضیح داده نشده (RSS) و r^2 هر یک از دو مدل فوق نیز با

یکدیگر برابر است؟

حل مسائل فصل اول

مسائل ۱-۱ تا ۱-۳ در واقع سؤالاتی است که از فصل اول کتاب، انتخاب شده و پاسخ آنها در متن کتاب آمده است.

مسائل ۱-۴ تا ۱-۶ دقیقاً مشابه مثالهایی است که در متن فصل اول کتاب حل شده است و اساساً شامل نکته‌های جدید نیست.

۱-۷ برای پاسخ به قسمت اول این مسأله کافی است معادله (۲) را در مدل رگرسیون مفروض قرار دهیم. خواهیم داشت

$$Y_i = (a_1 + a_2) + (b_1 + b_2) X_i + \varepsilon_i .$$

ملاحظه می‌شود که تأثیر X_i بر Y_i در مدل فوق - که در واقع صورت صحیح مدل (۱) است - یا $(b_1 + b_2)$ برابر است. چون فرض بر مثبت بودن b_2 است؛ بنابراین

$$b_1 < b_1 + b_2 ,$$

و در نتیجه b_1 تأثیر X_i بر Y_i در مدل (۱) را کمتر از میزان واقعی آن، یعنی $(b_1 + b_2)$ منعکس می‌سازد.

پاسخ قسمت دوم این مسأله منفی است؛ یعنی با فرض معادله (۳)، بیشتر فرضهای کلاسیک جمله اختلال در معادله (۱) برقرار نخواهد بود. برای تبیین این نکته معادله (۳) را در معادله (۱) قرار می‌دهیم،

$$Y_i = (a_1 + a_2) + b_1 X_i + (b_2 X_i^2 + \varepsilon_i) .$$

جمله اختلال در مدل فوق که صورت واقعی مدل رگرسیون (۱) است برابر با

$$w_i = b_2 X_i^2 + \varepsilon_i .$$

است. در اینجا فرضهای کلاسیک در مورد w_i را بررسی می‌کنیم. ملاحظه می‌شود که

امید ریاضی w_1 برابر صفر نیست؛ زیرا

$$E(w_1) = b_1 E(X_1^*) + E(\varepsilon_1) ,$$

است. با فرض $E(\varepsilon_1) = 0$ ، خواهیم داشت

$$E(w_1) = b_1 X_1^* \neq 0 .$$

چون مقادیر مختلف X_1 در زمان از یکدیگر مستقل نیست؛ بنابراین مقادیر w_1 نیز در دوره‌های زمانی مختلف از هم مستقل نبوده و در نتیجه فرض عدم خودهمبستگی برقرار نیست

$$\text{Cov}(w_1, w_2) \neq 0 .$$

با توجه به اینکه X_1 و X_1^* از یکدیگر مستقل نیستند، w_1 که تابعی از X_1^* است تابعی از X_1 نیز خواهد بود و می‌دانیم X_1 متغیر توضیحی در مدل رگرسیون است؛ بنابراین فرض استقلال متغیر توضیحی از جمله اختلال نیز باطل می‌شود،

$$E(w_1 | X_1) \neq X_1 E(w_1) \neq 0 .$$

همچنین چون w_1 تابعی از X_1 است و X_1 در زمان تغییر می‌کند؛ در نتیجه واریانس w_1 نیز در طی زمان ثابت نخواهد بود؛ در نتیجه به جای فرض کلاسیک واریانس همسانی، باید فرض واریانس ناهمسانی را بپذیریم. تنها جزئی که از فرضهای کلاسیک تغیر نخواهد کرد فرض توزیع نرمال جمله اختلال است.

۱-۸. چون X_1^* تابعی از ε_1 و ε_1 بنا بر فرض از U_1 مستقل است؛ بنابراین X_1^* از U_1 مستقل خواهد بود.

۲. با استفاده از معادله (۲) داریم

$$X_1 = X_1^* - \varepsilon_1 ,$$

که با جایگزینی در معادله (۱) خواهیم داشت

$$Y_1 = a + b X_1^* + (U_1 - b \varepsilon_1) ,$$

یا

$$Y_1 = a + b X_1^* + U_1^* .$$

۳. خیر؛ زیرا در این مدل جدید، متغیر توضیحی از جمله خطای مدل مستقل نیست. صحت این امر با ملاحظه معادله (۲) واضح است. X_i^* تابعی از ε_i است و چون ε_i قسمتی از جمله اختلال مدل جدید را تشکیل می‌دهد؛ یعنی $U_i^* = U_i - b \varepsilon_i$ ، پس X_i^* از U_i^* مستقل نخواهد بود. نتیجه می‌گیریم که فرض غیر تصادفی بودن متغیر توضیحی دیگر برقرار نیست.

۱-۹ ابتدا درباره $\hat{\beta}$ و β بحث می‌کنیم. می‌دانیم $\hat{\beta}_{OLS}$ با تخمین زتنده زیر به دست می‌آید،

$$\hat{\beta}_{OLS} = \frac{\sum x_i y_i}{\sum x_i^2}$$

با جایگزینی $y_i = \beta x_i + U_i$ خواهیم داشت

$$\begin{aligned} \hat{\beta}_{OLS} &= \frac{\sum x_i (\beta x_i + U_i)}{\sum x_i^2} \\ &= \beta + \frac{\sum x_i U_i}{\sum x_i^2} \end{aligned}$$

برای اینکه $\hat{\beta} = \beta$ باشد، باید $\frac{\sum x_i U_i}{\sum x_i^2}$ ضرورتاً برابر صفر شود. اما می‌دانیم $\sum x_i^2$ همواره مثبت است، بجز موردی که تمام مقادیر x_i با یکدیگر برابر باشند. بدیهی است در چنین صورتی محاسبه $\hat{\beta}$ ممکن نیست؛ چون منخرج کسر $\hat{\beta}_{OLS}$ برابر صفر است؛ بنابراین برای $\hat{\beta} = \beta$ کافی است $\sum x_i U_i$ برابر صفر شود. با توجه به اینکه مقادیر U_i ، هیچگاه شناخته نخواهند شد، نمی‌توان صحت رابطه $\sum x_i U_i = 0$ را ارزیابی کرد؛ بدین ترتیب، هیچ حالتی وجود نخواهد داشت که مطمئن باشیم $\hat{\beta} = \beta$.
در مورد امکان $\hat{\alpha} = \alpha$ ، ابتدا می‌گوییم

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

با توجه به رابطه

$$\bar{Y} = \alpha + \beta \bar{X} + \bar{U}$$

خواهیم داشت

$$\begin{aligned}\hat{\alpha} &= \alpha + \beta \bar{X} + \bar{U} - \hat{\beta} \bar{X} , \\ &= \alpha + \beta \bar{X} + \bar{U} - \left(\beta + \frac{\sum x_t u_t}{\sum x_t^2} \right) \bar{X} , \\ &= \alpha + \bar{U} - \frac{\sum x_t u_t}{\sum x_t^2} \bar{X} .\end{aligned}$$

برای اینکه $\hat{\alpha}$ برابر α شود، نه تنها باید شرط $\sum x_t u_t = 0$ برقرار باشد، بلکه \bar{U} نیز باید برابر صفر شود. دوباره این استدلال تکرار می‌شود که چون مقادیر U_t بنا بر تعریف، نامعلوم است، هیچ راهی وجود ندارد که بتوان، برای مجموعه معینی از مشاهدات داده‌ای، نشان داد که آیا شرایط $\sum x_t u_t = 0$ یا $\bar{U} = 0$ برقرار است یا خیر.

۱-۱۰. این سؤال دقیقاً موضوع بند ۳ مسأله (۱-۱) است. می‌دانیم

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2} ,$$

بنابراین، اگر مقادیر X_t ثابت باشد، آنگاه $\bar{X} = X_t$ و در نتیجه $x_t = X_t - \bar{X} = 0$ یا $\sum x_t^2 = 0$. در اینصورت محاسبه $\hat{\beta}$ ممکن نخواهد بود.

۲. در مدل $Y_t = \beta X_t + U_t$ داریم

$$\hat{\beta} = \frac{\sum X_t Y_t}{\sum X_t^2} ,$$

و اگر فرض کنیم مقادیر X_t ثابت باشد، خواهیم داشت

$$\hat{\beta} = \frac{\sum X_t Y_t}{\sum X_t^2} = \frac{X_t \sum Y_t}{n X_t^2} = \frac{\bar{Y}}{\bar{X}} .$$

نتیجه فوق را می‌توان تعمیم داد. مدل‌های رگرسیون را که فاقد متغیر توضیحی بوده اما ضریب ثابت دارند، می‌توان به صورت زیر نوشت،

$$Y_t = \alpha + U_t .$$

اگر α را با روش حداقل مربعات معمولی تخمین بزنیم، خواهیم داشت

$$\hat{\alpha}_{OLS} = \bar{Y} .$$

اثبات این نکته بسیار ساده است. مدل فوق را به صورت زیر می‌نویسیم،

$$Y_i = \alpha X_i + U_i ,$$

که در آن برای تمام مقادیر i داریم $X_i = 1$. می‌دانیم

$$\hat{\alpha}_{OLS} = \frac{\sum X_i Y_i}{\sum X_i^2} ,$$

با توجه به $X_i = 1$ ، خواهیم داشت $\hat{\alpha}_{OLS} = \bar{Y}$.

۱-۱۱. می‌دانیم $e_i = Y_i - \hat{Y}_i$ است؛ بنابراین

$$e_i = Y_i - \hat{\alpha} - \hat{\beta} X_i .$$

می‌توان نوشت،

$$\sum e_i = \sum Y_i - n \hat{\alpha} - \hat{\beta} \sum X_i ,$$

یا

$$\sum e_i = \sum Y_i - n (\bar{Y} - \hat{\beta} \bar{X}) - \hat{\beta} \sum X_i ,$$

$$= \hat{\beta} \sum X_i - \hat{\beta} \sum X_i ,$$

$$= 0 .$$

۲. با توجه به مفهوم U_i ، ملاحظه می‌شود که اگر مجموع پسماندها برابر صفر باشد، هیچ دلیلی وجود ندارد که مجموع جملات اختلال نیز صفر شود. U_i و e_i مفاهیم کاملاً متفاوتی هستند و چون مقادیر U_i مطلقاً ناشناخته است نمی‌توان در مورد مجموع مقادیر آنها هیچگونه اظهار نظری کرد.

۱-۱۲ از مدل اولیه داریم

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} , \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} .$$

از دو مدل دوم داریم

$$\hat{\beta}_i = \frac{\sum x_{it} y_{it}}{\sum x_{it}^2}, \quad i = 1, 2$$

$$\hat{\alpha}_i = \bar{Y}_i - \hat{\beta}_i \bar{X} \quad i = 1, 2$$

می توان نوشت،

$$\begin{aligned} \sum \hat{\beta}_i &= \hat{\beta}_1 + \hat{\beta}_2, \\ &= \frac{\sum (y_{1t} + y_{2t}) x_{1t}}{\sum x_{1t}^2} = \frac{\sum x_{1t} y_{1t}}{\sum x_{1t}^2} = \hat{\beta}. \end{aligned}$$

حال برای نشان دادن صحت $\hat{\alpha} = \hat{\alpha}_1 + \hat{\alpha}_2$ چنین می نویسیم،

$$\begin{aligned} \sum \hat{\alpha}_i &= \hat{\alpha}_1 + \hat{\alpha}_2 = \sum \bar{Y}_i - \sum \hat{\beta}_i \bar{X}, \\ &= \bar{Y}_1 + \bar{Y}_2 - (\hat{\beta}_1 + \hat{\beta}_2) \bar{X}, \\ &= \bar{Y} - \hat{\beta} \bar{X}, \\ &= \hat{\alpha}. \end{aligned}$$

البته باید توجه کرد که این قاعده در مورد متغیر توضیحی صادق نیست؛ یعنی اگر X_{it} را به k قسمت تقسیم کنیم و به کمک Y_i برای هر سری از مشاهدات X_i ، $i = 1, 2, \dots, k$ یک مدل رگرسیون بسازیم، مشابه نتیجه فوق حاصل نخواهد شد.

۱-۱۳ می دانیم

$$\hat{\beta}_i = \frac{\sum_t x_{it} x_{it}}{\sum_t x_{it}^2}, \quad \hat{\alpha}_i = \bar{X}_i - \hat{\beta}_i \bar{X},$$

اگر برای تمام مقادیر i جمع کنیم، خواهیم داشت

$$\sum_{i=1}^k \hat{\beta}_i = \frac{\sum_i \sum_t x_{it} x_{it}}{\sum_t x_{it}^2} = \frac{\sum_t x_{it}^2}{\sum_t x_{it}^2} = 1,$$

$$\sum_i \hat{\alpha}_i = \sum_i \bar{x}_i - \sum_i \hat{\beta}_i \bar{x}$$

می‌دانیم $\sum_i \bar{x}_i = \bar{x}$ و همچنین $\sum_i \hat{\beta}_i \bar{x} = \bar{x}$ ؛ زیرا

$$\sum_i \hat{\beta}_i \bar{x} = \bar{x} \sum_i \hat{\beta}_i = \bar{x} (1) = \bar{x}$$

بنابراین خواهیم داشت

$$\sum_i \hat{\alpha}_i = \bar{x} - \bar{x} = 0$$

۱-۱۴ می‌دانیم

$$\hat{b} = \frac{\sum y_t \hat{y}_t}{\sum \hat{y}_t}$$

با توجه به $\hat{y}_t = \hat{\beta} x_t$ داریم

$$\hat{b} = \frac{\sum y_t (\hat{\beta} x_t)}{\sum (\hat{\beta} x_t)} = \frac{\hat{\beta} \sum x_t y_t}{\hat{\beta} \sum x_t} = 1$$

از طرف دیگر می‌دانیم

$$\hat{a} = \bar{Y} - \hat{b} \bar{Y} \quad (1)$$

اما:

$$\bar{Y} = \frac{\sum (\hat{\alpha} + \hat{\beta} x_t)}{n} = \hat{\alpha} + \hat{\beta} \bar{x}$$

با جایگزینی $\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{x}$ در رابطه فوق داریم.

$$\bar{Y} = \bar{Y}$$

اگر در نتیجه فوق یعنی $\hat{b} = 1$ و $\bar{Y} = \bar{Y}$ را در فرمول (۱) قرار دهیم خواهیم داشت

$$\hat{a} = 0$$

۱-۱۵ اگر $\hat{Y} = \hat{\beta} x_t$ بخواند از نقطه (\bar{x}, \bar{Y}) بگذرد، باشد داشته باشیم

$$\bar{Y} = \hat{\beta} \bar{X} .$$

و چون می‌دانیم

$$\hat{\beta} = \frac{\sum X_t Y_t}{\sum X_t^2} ,$$

بنابراین رابطه زیر باید همواره برقرار باشد،

$$\bar{Y} = \frac{\bar{X} \sum X_t Y_t}{\sum X_t^2} .$$

ملاحظه می‌شود که هیچ دلیلی بر صحت رابطه فوق به ازای جمیع مقادیر t وجود ندارد. با اینکه مدل‌های رگرسیون با ضریب ثابت، همواره از نقطه (\bar{X}, \bar{Y}) می‌گذرند، این خصوصیت برای مدل‌های رگرسیون فاقد ضریب ثابت معمولاً برقرار نیست.

۱-۱۶ هر چند در فصل هشتم، این گونه تخمین‌های مقید به طور گسترده مطرح خواهد شد، در این فصل به طور خلاصه به آن اشاره می‌شود. همان گونه که مسأله تخمین در اقتصادسنجی به مسأله حداقل سازی توابع تبدیل می‌شود، تخمین‌های مقید نیز به حداقل سازی مقید تبدیل می‌گردد. در چنین مواردی برای اینکه مجموع مربعات پسماندها را با توجه به شرط یا قید مفروض حداقل کنیم، دو راه وجود دارد: یا باید قید را در معادله‌ای که قرار است حداقل شود جایگزین کنیم یا با روش ضریب لاگرانژ عمل کنیم. راه حل دوم بیشتر معمول است. اگر λ را ضریب لاگرانژ بنامیم، با توجه به اینکه قید ما صفر بودن $\hat{\alpha}$ است، باید مجموع مربعات پسماند را با توجه به این قید حداقل کنیم. در واقع باید عبارت زیر حداقل شود،

$$s = \sum e_i + \lambda \hat{\alpha} .$$

با توجه به اینکه $e_i = Y_i - \hat{Y}_i$ یا $e_i = Y_i - \hat{\alpha} - \hat{\beta} X_i$ ؛ بنابراین

$$s = \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2 + \lambda \hat{\alpha} .$$

از S نسبت به $\hat{\alpha}$ ، $\hat{\beta}$ و λ مشتق می‌گیریم،

$$\frac{\partial s}{\partial \hat{\alpha}} = -2 \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i) + \lambda = 0 ,$$

$$\frac{\partial s}{\partial \hat{\beta}} = -2 \sum X_i (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0 ,$$

$$\frac{\partial s}{\partial \lambda} = \hat{\alpha} = 0 .$$

با جایگزینی معادله سوم در معادله دوم، به معادله نرمال برای مدل رگرسیون بدون ضریب ثابت می‌رسیم. اگر آن را برای $\hat{\beta}$ حل کنیم، خواهیم داشت

$$\hat{\beta} = \frac{\sum X_i Y_i}{\sum X_i^2} . \quad (1)$$

همچنین از معادله اول نیز داریم

$$\lambda = 2 \sum (Y_i - \hat{\beta} X_i) .$$

اگر λ صفر شود، در آن صورت قید مفروض بی‌تأثیر است؛ یعنی اساساً فرقی نمی‌کند که تخمین با توجه به وجود قید صورت پذیرد یا اینکه یک تخمین آزاد باشد. حال باید دید که در چه صورتی λ برابر صفر است. بدیهی است موقعی λ برابر صفر است که

$$\sum (Y_i - \hat{\beta} X_i) = 0 ,$$

یعنی $\bar{Y} = \hat{\beta} \bar{X}$ رابطه اخیر موقعی صحیح است که تخمین مدل رگرسیون $(\hat{Y}_i = \hat{\beta} X_i)$ از نقطه (\bar{X}, \bar{Y}) بگذرد.

با جایگزینی $\hat{\beta}$ در معادله $\bar{Y} = \hat{\beta} \bar{X}$ خواهیم داشت

$$\bar{Y} = \frac{\bar{X} \sum X_i Y_i}{\sum X_i^2} . \quad (2)$$

قبلاً در مسأله ۱-۱۵ نیز دقیقاً به معادله (۲) رسیده بودیم. از طرف دیگر می‌دانیم که تخمین β از مدل $Y_i = \alpha + \beta X_i + U_i$ برابر است با

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2} \quad (3)$$

با یک جایگزینی ساده می‌توان نشان داد که موقعی تخمین β از دو مدل فوق به یک جواب می‌رسد - یعنی (۱) و (۳) برابر است - که رابطه (۲) برقرار باشد.

۱-۱۷ این مسأله نکته خاصی در اقتصادسنجی ندارد و بیشتر کاربرد روش حداقل سازی معادله‌های، بدون استفاده از مشتق‌گیری است. می‌دانیم

$$\begin{aligned} \sum e_t^2 &= \sum (Y_t - \hat{Y}_t)^2 = \sum (Y_t - \hat{\beta} X_t)^2, \\ &= \sum Y_t^2 + \hat{\beta}^2 \sum X_t^2 - 2\hat{\beta} \sum X_t Y_t. \end{aligned}$$

متغیر ما، $\hat{\beta}$ است. معادله فوق را برحسب توانهای $\hat{\beta}$ مرتب می‌کنیم،

$$\sum e_t^2 = A_0 \hat{\beta}^2 + A_1 \hat{\beta} + A_2, \quad (1)$$

که در آن $A_0 = \sum X_t^2$ و $A_1 = -2 \sum X_t Y_t$ و $A_2 = \sum Y_t^2$. سعی می‌کنیم یک مربع کامل برحسب $\hat{\beta}$ از رابطه (۱) به دست آوریم،

$$\begin{aligned} \sum e_t^2 &= A_0 \left(\hat{\beta}^2 + \frac{A_1}{A_0} \hat{\beta} + \frac{A_1^2}{4A_0^2} \right) - \left(\frac{A_1^2}{4A_0} - A_2 \right), \\ &= A_0 \left(\hat{\beta} + \frac{A_1}{2A_0} \right)^2 - \left(\frac{A_1^2}{4A_0} - A_2 \right). \end{aligned}$$

فقط جمله اول، تابعی از $\hat{\beta}$ است؛ بنابراین $\sum e_t^2$ هنگامی حداقل است که این جمله حداقل باشد. با توجه به اینکه این جمله مربع کامل است، حداقل آن برابر صفر است.

$$\hat{\beta} + \frac{A_1}{2A_0} = 0,$$

$$\hat{\beta} = -\frac{A_1}{2A_0} = -\frac{-2 \sum X_t Y_t}{2 \sum X_t^2} = \frac{\sum X_t Y_t}{\sum X_t^2}.$$

۱-۱۸ هر یک از دو مدل را یک بار دیگر می‌نویسیم،

$$Y_i = \alpha + \beta X_i + U_i, \quad (1)$$

$$Y_i^* = \alpha + \beta X_i^* + U_i. \quad (2)$$

توجه داریم که U_i تغییر نمی‌کند؛ زیرا فقط مقیاس اندازه‌گیری هر یک از دو متغیر را تغییر داده‌ایم که قاعدتاً در ساختار دینامیکی مدل تأثیری ندارد. بدیهی است که

$$\hat{b} = \frac{\sum x_i^* y_i^*}{\sum x_i^{*2}} \quad (3)$$

حال می‌گوییم $y_i^* = y_i$ ؛ برای اثبات کافی است در دو طرف این معادله، مقادیر هر یک را جایگزین کنیم خواهیم داشت

$$Y_i^* - \bar{Y}^* = Y_i - \bar{Y}.$$

با جایگزینی $Y_i^* = Y_i - Y^*$ داریم:

$$Y_i - Y^* - \frac{\sum (Y_i - Y^*)}{n} = Y_i - \bar{Y},$$

$$Y_i - Y^* - \bar{Y} + \frac{n Y^*}{n} = Y_i - \bar{Y}.$$

ملاحظه می‌شود که رابطه فوق همواره برقرار است. به همین ترتیب می‌توان نشان داد که $x_i^* = x_i$. با جایگزینی مقادیر $y_i^* = y_i$ و $x_i^* = x_i$ در معادله (۳) خواهیم داشت

$$\hat{b} = \frac{\sum x_i y_i}{\sum x_i^2} = \hat{\beta}.$$

نتیجه می‌گیریم که تخمین شیب معادله رگرسیونی در هر دو مدل دقیقاً با یکدیگر برابر است. اما در مورد ضریب ثابت به ترتیب زیر عمل می‌کنیم. با استفاده از مدل (۲) داریم

$$\hat{a} = \bar{Y}^* - \hat{b} \bar{X}^*,$$

$$= \frac{\sum (Y_i - Y^*)}{n} - \hat{b} \frac{\sum (X_i - X^*)}{n},$$

$$= \bar{Y} - Y^* - \hat{b} \bar{X} + \hat{b} X^*.$$

با توجه به معادله (۱) می‌دانیم

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X},$$

در نتیجه، خواهیم داشت

$$\hat{a} = \hat{\alpha} - Y^* + \hat{\beta} X^*.$$

۱-۱۹ تساوی معادله‌های (۱)، (۲) و (۴) را در متن دیدیم. در اینجا تنها به اثبات

تساوی دو معادله (۱) و (۳) می‌پردازیم. می‌دانیم

$$\hat{y}_i = \hat{\beta} x_i.$$

با جایگزینی در معادله (۳) خواهیم داشت

$$r^2 = \frac{[\sum y_i (\hat{\beta} x_i)]^2}{\sum \hat{y}_i^2 \sum (\hat{\beta} x_i)^2},$$

$$= \frac{\hat{\beta}^2 (\sum x_i y_i)^2}{\hat{\beta}^2 \sum x_i^2 \sum y_i^2} = \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2},$$

که در واقع همان معادله (۱) است.

۱-۲۰ می‌دانیم، $\hat{Y}_i = \hat{\beta} X_i$. با جایگزینی در معادله (۲) داریم

$$r^2 = \frac{(\sum Y_i \hat{\beta} X_i)^2}{\sum Y_i^2 \sum (\hat{\beta} X_i)^2} = \frac{\hat{\beta}^2 (\sum X_i Y_i)^2}{\hat{\beta}^2 \sum X_i^2 \sum Y_i^2},$$

$$= \frac{(\sum X_i Y_i)^2}{\sum X_i^2 \sum Y_i^2},$$

که در واقع همان معادله (۱) است.

برای نشان دادن تساوی معادله‌های (۳) و (۱) کافی است معادله (۳) را به صورت

زیر بنویسیم،

$$r^2 = 1 - \frac{\sum (Y_i - \hat{\beta} X_i)^2}{\sum Y_i^2}$$

با جایگزینی $\hat{\beta} = \frac{\sum X_i Y_i}{\sum X_i^2}$ خواهیم داشت

$$\begin{aligned} r^2 &= 1 - \frac{\left[\sum Y_i^2 - \frac{\sum (X_i Y_i)^2}{\sum X_i^2} \right]}{\sum Y_i^2} \\ &= \frac{(\sum X_i Y_i)^2}{\sum X_i^2 \sum Y_i^2} \end{aligned}$$

که دقیقاً همان معادله (۱) است. و معادله (۱) نیز تعریف r^2 در مدلهایی است که جمله ثابت ندارند.

۱-۲۱ راه حل اول: برای اثبات، ابتدا به ذکر یک مقدمه ریاضی می‌پردازیم. نامساوی

$$\sum a_i^2 \sum b_i^2 \geq (\sum a_i b_i)^2$$

به نامساوی کوچی - شوارتز معروف است. با توجه به اینکه اثبات دقیق این نامساوی مستلزم دقت‌های ریاضی است - که ضرورتاً کاربردی در تحلیلهای اقتصادسنجی ندارد - در اینجا به طرح روش بسیار ساده اثبات این نامساوی می‌پردازیم. ابتدا نامساوی را برای حالتی می‌نویسیم که $t = 1$ و 2 ، خواهیم داشت

$$\sum_{i=1}^t a_i^2 \sum_{i=1}^t b_i^2 \geq \left(\sum_{i=1}^t a_i b_i \right)^2$$

$$(a_1^2 + a_2^2) (b_1^2 + b_2^2) \geq (a_1 b_1 + a_2 b_2)^2$$

دو طرف نامساوی را بسط می دهیم و ساده می کنیم،

$$a_1^2 b_2^2 + a_2^2 b_1^2 - 2 a_1 a_2 b_1 b_2 \geq 0 ,$$

یا

$$(a_1 b_2 - a_2 b_1)^2 \geq 0 ,$$

که همواره مثبت است، مگر در حالتی که $\frac{a_1}{b_1} = \frac{a_2}{b_2}$ که در آن صورت برابر صفر می شود. در حالت کلی، اگر جمله های $a_1^2 b_2^2$ را برای تمام مقادیر t از دو طرف نامساوی حذف کنیم، خواهیم داشت

$$\sum_{s=1}^n \sum_{v=1}^n (a_s b_v - a_v b_s)^2 \geq 0 , \quad s \neq v$$

که همواره مثبت است، مگر در حالتی که $\frac{a_s}{b_s} = \frac{a_v}{b_v}$ که در آن صورت برابر صفر می شود. بعد از یادآوری این نکته ریاضی به بحث خود برمی گردیم. می دانیم

$$r^2 = 1 - \frac{RSS}{TSS}$$

در مدل $Y_t = \beta X_t + U_t$ داریم

$$RSS = \sum e_t^2 = \sum (Y_t - \hat{Y}_t)^2 ,$$

$$= \sum (Y_t - \hat{\beta} X_t)^2 ,$$

$$= \sum Y_t^2 + \hat{\beta}^2 \sum X_t^2 - 2 \hat{\beta} \sum X_t Y_t ,$$

که با جایگزینی $\hat{\beta} = \frac{\sum X_t Y_t}{\sum X_t^2}$ خواهیم داشت

$$RSS = \sum Y_t^2 - \frac{(\sum X_t Y_t)^2}{\sum X_t^2} .$$

همچنین می دانیم

$$TSS = \sum y_t^2 = \sum (Y_t - \bar{Y})^2 .$$

بنابراین

$$\begin{aligned} r^2 &= 1 - \frac{RSS}{TSS} , \\ &= 1 - \frac{\sum Y_i' - \frac{(\sum X_i Y_i)'}{\sum X_i'}}{\sum (Y_i - \bar{Y})'} , \\ &= 1 - \left[\frac{\sum Y_i'}{\sum (Y_i - \bar{Y})'} - \frac{(\sum X_i Y_i)'}{\sum (Y_i - \bar{Y})' \sum X_i'} \right] , \\ &= 1 - \left[\frac{\sum Y_i'}{\sum (Y_i - \bar{Y})'} - \frac{(\sum X_i Y_i)'}{\sum X_i' \sum Y_i'} \cdot \frac{\sum Y_i'}{\sum (Y_i - \bar{Y})'} \right] . \end{aligned}$$

اگر فرض کنیم

$$q = \frac{\sum Y_i'}{\sum (Y_i - \bar{Y})'} , \quad r^{*2} = \frac{(\sum X_i Y_i)'}{\sum X_i' \sum Y_i'} ,$$

آنگاه خواهیم داشت

$$\begin{aligned} r^2 &= 1 - [q - r^{*2} q] , \\ &= 1 - (1 - r^{*2}) q . \end{aligned} \quad (1)$$

در اینجا به بررسی q و r^{*2} می پردازیم. می دانیم

$$q = \frac{\sum_{i=1}^n Y_i'}{\sum_{i=1}^n (Y_i - \bar{Y})'} = \frac{Y_1' + Y_2' + \dots + Y_n'}{(Y_1 - \bar{Y})' + (Y_2 - \bar{Y})' + \dots + (Y_n - \bar{Y})'}$$

صورت و مخرج کسر q را بر \bar{Y}' تقسیم می کنیم،

$$q = \frac{\frac{Y_1^2}{\bar{Y}^2} + \frac{Y_2^2}{\bar{Y}^2} + \dots + \frac{Y_n^2}{\bar{Y}^2}}{\frac{Y_1^2 + \bar{Y}^2 - 2Y_1\bar{Y}}{\bar{Y}^2} + \frac{Y_2^2 + \bar{Y}^2 - 2Y_2\bar{Y}}{\bar{Y}^2} + \dots + \frac{Y_n^2 + \bar{Y}^2 - 2Y_n\bar{Y}}{\bar{Y}^2}}$$

$$= \frac{\sum \left(\frac{Y_t}{\bar{Y}}\right)^2}{\sum \left(\frac{Y_t}{\bar{Y}}\right)^2 - n}$$

هنگامی که تمام مقادیر Y_t با یکدیگر برابر شوند حد بالای q برابر $+\infty$ است و وقتی که تغییرات Y_t بسیار زیاد باشد، حد پایین آن برابر یک است. اما با توجه به اینکه r^{*2} بنا بر قضیه کوچکی - شوارتز بین صفر و یک نوسان می‌کند، می‌توان اینگونه نتیجه گرفت که به ازای مقادیر کوچک r^{*2} و مقادیر بزرگ q ، چه بسا مقدار r^2 بتواند منفی شود. البته توجه داریم که مقادیر r^{*2} و q مستقل از یکدیگر تعیین می‌شوند.

راه حل دوم: راه حل ساده تری نیز برای این مسأله وجود دارد؛ با وجود این، به این نکته اشاره می‌کنیم که مزیت راه حل فوق این است که می‌تواند رابطه بین ضریب تعیین معمولی r^2 و ضریب تعیین برای مدل‌های بدون جمله ثابت r^{*2} را نشان دهد. به هر حال راه حل دوم بسیار خلاصه و به شرح زیر است.

در مدل $Y_t = \beta X_t + U_t$ ، می‌دانیم

$$Y_t = \hat{Y}_t + e_t .$$

با توجه به $\hat{Y}_t = \hat{\beta} X_t$ ، خواهیم داشت

$$\hat{Y}_t = \hat{\beta} X_t + e_t .$$

دو طرف رابطه فوق را مجذور کرده و برای تمام مقادیر t جمع می‌کنیم،

$$\sum Y_t^2 = \hat{\beta}^2 \sum X_t^2 + \sum e_t^2 + 2\hat{\beta} \sum X_t e_t .$$

می‌دانیم $\sum X_i e_i$ بنا بر معادلهٔ نرمال صفر است؛ زیرا

$$\begin{aligned} \frac{d \sum e_i^2}{d \hat{\beta}} &= \frac{d \sum (Y_i - \hat{Y}_i)^2}{d \hat{\beta}} = \frac{d \sum (Y_i - \hat{\beta} X_i)^2}{d \hat{\beta}} \\ &= -2 \sum X_i (Y_i - \hat{\beta} X_i) = -2 \sum X_i e_i = 0 \end{aligned}$$

بنابراین خواهیم داشت:

$$\sum Y_i^2 = \sum \hat{Y}_i^2 + \sum e_i^2$$

اگر r^2 را برای اینگونه مدلها که فاقد جملهٔ ثابت هستند به صورت

$$r^2 = 1 - \frac{\sum e_i^2}{\sum Y_i^2}$$

تعریف کنیم، آنگاه $0 < r^2 < 1$ در غیر این صورت اگر r^2 را به روال معمول تعریف نماییم، خواهیم داشت

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2} = 1 - \frac{\sum e_i^2}{\sum Y_i^2 - n \bar{Y}^2} = \frac{\sum Y_i^2 - n \bar{Y}^2 - (\sum Y_i^2 - \hat{\beta}^2 \sum X_i^2)}{\sum Y_i^2 - n \bar{Y}^2}$$

در نتیجه داریم

$$r^2 = \frac{\hat{\beta}^2 \sum X_i^2 - n \bar{Y}^2}{\sum y_i^2}$$

مخرج کسر همواره مثبت است، اما در مواردی که $n \bar{Y}^2 > \hat{\beta}^2 \sum X_i^2$ مقدار r^2 منفی خواهد شد. یادآوری می‌کنیم در مدلهایی که ضریب ثابت ندارند $\bar{Y} \neq \bar{y}$ ؛ بنابراین $\hat{\beta}^2 \sum X_i^2 = \sum \hat{Y}_i^2 \neq \sum \hat{y}_i^2 + n \bar{Y}^2$

۱-۲۲. مدل (۲) را به صورت زیر می‌نویسیم،

$$Y_i^* = a + b X_i + U_i \quad (3)$$

که در آن، $Y_i^* = Y_i - X_i$ و $a = \alpha$ و $b = (\beta - 1)$ به سهولت ملاحظه می شود که با تخمین b از مدل (۲) یا (۳) می توان β را به دست آورد؛ زیرا

$$\hat{b} = \frac{\sum x_i y_i}{\sum x_i^2} ,$$

$$= \frac{\sum x_i (y_i - x_i)}{\sum x_i^2} = \hat{\beta} - 1 .$$

همچنین ملاحظه می شود که $\hat{a} = \hat{\alpha}$ ؛ به عبارت دقیقتر از معادله (۳) داریم

$$\hat{a} = \bar{Y}^* - \hat{b} \bar{X} .$$

با توجه به $\hat{b} = \hat{\beta} - 1$ داریم

$$\begin{aligned} \hat{a} &= \bar{Y} - \bar{X} - (\hat{\beta} - 1) \bar{X} , \\ &= \bar{Y} - \bar{X} - \hat{\beta} \bar{X} + \bar{X} , \\ &= \bar{Y} - \hat{\beta} \bar{X} , \\ &= \hat{\alpha} . \end{aligned}$$

۲. با استفاده از فرمول $r^2 = 1 - \frac{RSS}{\sum y_i^2}$ ، برای مدل (۱) داریم

$$\frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2} = 1 - \frac{RSS}{\sum y_i^2} ,$$

یا:

$$RSS = \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} .$$

اگر RSS (تغییرات توضیح داده نشده) را برای مدل (۲) به صورت RSS^* تعریف کنیم، خواهیم داشت

$$RSS^* = \sum y_i^{*2} - \frac{(\sum x_i y_i^*)^2}{\sum x_i^2} .$$

می‌دانیم $y_i^* = y_i - x_i$ در نتیجه داریم

$$RSS^* = \sum (y_i - x_i)^2 - \frac{[\sum x_i (y_i - x_i)]^2}{\sum x_i^2}$$

که بعد از ساده کردن، عبارت است از

$$\begin{aligned} RSS^* &= \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} \\ &= RSS \end{aligned}$$

نتیجه می‌گیریم که مجموع مربعات پسماند در هر دو مدل با یکدیگر برابر است. در اینجا به بررسی r^2 در هر یک از دو مدل (۱) و (۲) می‌پردازیم. اگر r^2 را برای مدل دوم به صورت r^{*2} تعریف کنیم، خواهیم داشت

$$r^{*2} = \frac{(\sum x_i y_i^*)^2}{\sum x_i^2 \sum y_i^{*2}} \quad (4)$$

رابطه $y_i^* = y_i - x_i$ را جایگزین می‌کنیم،

$$r^{*2} = \frac{(\sum x_i y_i)^2 + M}{\sum x_i^2 \sum y_i^2 + M} \quad (5)$$

که در آن:

$$M = (\sum x_i^2)^2 - 2(\sum x_i y_i)(\sum x_i^2)$$

با ملاحظه معادله (۴) و بنا بر نامساوی کوچی - شوارتز، می‌دانیم

$$\sum x_i^2 \sum y_i^{*2} \geq (\sum x_i y_i^*)^2$$

در معادله (۴) مخرج کسر معمولاً از صورت آن بزرگتر است. می‌دانیم اگر به صورت و مخرج کسر $\frac{A}{B}$ ، که در آن $A < B$ ، یک عدد مثبت اضافه کنیم، مقدار کسر بزرگتر می‌شود

و در صورتی که آن عدد منفی باشد، مقدار کسر کوچکتر خواهد شد،

$$\frac{A}{B} < \frac{A+M}{B+M} .$$

برای اثبات کافی است که مخرجها را به یک مخرج مشترک تبدیل کنیم. با توجه به اینکه در معادله (۵)، M نمی تواند مقداری باشد که r^{*2} منفی شود، داریم

$$M > 0 \quad \text{هرگاه} \quad r^* > r^2 ,$$

$$M < 0 \quad \text{هرگاه} \quad r^* < r^2 .$$

اما $M > 0$ بدین معنی است که

$$(\sum x_i^2)^2 > 2 (\sum x_i y_i) (\sum x_i^2) ,$$

یا

$$\frac{\sum x_i y_i}{\sum x_i^2} < \frac{1}{2} ,$$

یعنی

$$\hat{\beta} < \frac{1}{2} .$$

نتیجه می گیریم که اگر در مدل اولیه $\hat{\beta} < \frac{1}{2}$ ، می توان مدل را چنان تبدیل کرد که ضریب تعیین جدید آن بزرگتر از حالت قبلی باشد. باید توجه داشت که این قاعده عمومی است و برای حالت زیر نیز صادق خواهد بود،

$$Y_i - k X_i = \alpha + (\beta - k) X_i + U_i .$$

آزمونهای آماری و خصوصیات مطلوب تخمین زنده‌ها در مدل رگرسیون خطی ساده

۲-۱ مقدمه

اگر مدل رگرسیون و تخمین آن به ترتیب عبارت باشند از:

$$Y_t = \alpha + \beta X_t + U_t ,$$

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t ,$$

آنگاه مقادیر $\hat{\alpha}_{OLS}$ و $\hat{\beta}_{OLS}$ را می‌توان از طریق فرمولهای زیر به دست آورد،

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2} , \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} .$$

$\hat{\beta}$ یک متغیر تصادفی است؛ زیرا تابعی از y_t می‌باشد که خود تابعی از U_t است. همچنین $\hat{\alpha}$ که تابعی از $\hat{\beta}$ است تصادفی خواهد شد. $\hat{\beta}$ و $\hat{\alpha}$ مانند هر متغیر تصادفی دارای تابع توزیع احتمال با میانگین و واریانس معین هستند.

یکی از هدفهای اصلی در محاسبه $\hat{\alpha}$ و $\hat{\beta}$ به دست آوردن اطلاعات درباره α و β واقعی است؛ به عبارت دیگر، با داشتن $\hat{\alpha}$ و $\hat{\beta}$ باید بتوان فرضیه‌های مختلفی را در مورد α و β آزمون کرد. در قسمت ۲-۳ آزمون فرضیه‌های مربوط به هر یک از پارامترها را مطالعه خواهیم کرد. در این آزمونها از آماره‌های Z و t استفاده خواهد شد. تعیین فاصله‌های اطمینان برای پارامترهای α و β و رابطه بین فاصله اطمینان و دقت تخمین پارامترها از موضوعات دیگر این قسمت است. آزمون واریانس جمله اختلال، که به کمک آماره F انجام می‌شود، نیز در همین قسمت

بررسی خواهد شد.

کاربرد مسأله آنالیز واریانس در آزمون معنی دار بودن کل مدل رگرسیون موضوع قسمت ۲-۴ است. بعد از معرفی اجمالی آنالیز واریانس و ارائه تعاریف توزیع و آزمون F ، به چگونگی اجرای آزمون معنی دار بودن مدل رگرسیون می پردازیم. رابطه بین آزمونهای F و t و نیز رابطه بین آزمونهای معنی دار بودن ضریب تعیین و معنی دار بودن کل رگرسیون از اهمیت خاصی برخوردار است که در همین قسمت به طور خلاصه بررسی می شود.

خصوصیات مطلوب تخمین زنده‌ها در حالت کلی، موضوع قسمت ۲-۵ را تشکیل می دهد. این سؤال از لحاظ نظری اهمیت بسیاری دارد که اگر $\hat{\theta}$ تخمینی از پارامتر θ باشد، آنگاه خصوصیات مطلوب $\hat{\theta}$ چیست؟ خصوصیات مطلوب $\hat{\theta}$ را در دو بخش بررسی خواهیم کرد. ابتدا این خصوصیات را در نمونه‌هایی مطالعه می کنیم که حجم آنها محدود است، سپس همین مسأله، با عنوان خصوصیات حدی تخمین زنده‌ها، برای نمونه‌های بسیار بزرگ بررسی خواهد شد. بدیهی است، مطالعه این خصوصیات، مستقل از روشی است که به کمک آن پارامتر θ را تخمین زده‌ایم. در قسمت ۲-۶ به تحلیل این مسأله می پردازیم که اگر پارامتر مفروض را با روش حداقل مربعات معمولی (OLS) تخمین بزنیم، آیا خصوصیات مطلوب تخمین زنده‌ها را خواهد داشت؟ خواهیم دید که بنا بر قضیه گاس - مارکف پاسخ این سؤال مثبت است.

۲-۲ خصوصیات آماری $\hat{\alpha}$ و $\hat{\beta}$

مدل رگرسیون ۱-۸ و تخمین آن، معادله ۱-۱۶ را یک بار دیگر می نویسیم،

$$Y_i = \alpha + \beta X_i + U_i ,$$

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i .$$

بحث را روی خصوصیات آماری $\hat{\beta}$ متمرکز می کنیم. برای $\hat{\alpha}$ می توان دقیقاً استدلال مشابهی عرضه کرد.

در معادله ۱-۲۵ دیدیم که^۱

$$\hat{\beta}_{OLS} = \frac{\sum x_i y_i}{\sum x_i^2}$$

می دانیم

$$y_i = \beta x_i + u_i \quad (2-1)$$

زیرا کافی است از دو طرف مدل ۱-۸ میانگین گرفته و معادله به دست آمده را از ۱-۸ کم کنیم، خواهیم داشت

$$(Y_i - \bar{Y}) = (\alpha - \alpha) + \beta (X_i - \bar{X}) + (U_i - \bar{U})$$

یا:

$$y_i = \beta x_i + u_i$$

معادله ۲-۱ را در ۱-۲۵ جایگزین می کنیم

$$\hat{\beta} = \frac{\sum x_i (\beta x_i + u_i)}{\sum x_i^2}$$

$$\hat{\beta} = \beta + \frac{\sum x_i u_i}{\sum x_i^2} \quad (2-2)$$

$\hat{\beta}$ تابعی از u_i است، در نتیجه یک متغیر تصادفی است.

۱. در فرمول $\hat{\beta}$ می توان در صورت کسر به جای $\sum x_i y_i$ مقدار $\sum x_i Y_i$ را نوشت؛ زیرا

$$\sum x_i y_i = \sum x_i (Y_i - \bar{Y}) = \sum x_i Y_i - \bar{Y} \sum x_i = \sum x_i Y_i$$

۲. در معادله ۲-۲ می توان در صورت کسر به جای $\sum x_i u_i$ مقدار $\sum x_i U_i$ را نوشت؛ زیرا

$$\sum x_i u_i = \sum x_i (U_i - \bar{U}) = \sum x_i U_i - \bar{U} \sum x_i = \sum x_i U_i$$

توجه داریم که $\sum x_i = 0$ ؛ زیرا

$$\sum x_i = \sum (X_i - \bar{X}) = \sum X_i - n \bar{X} = 0$$

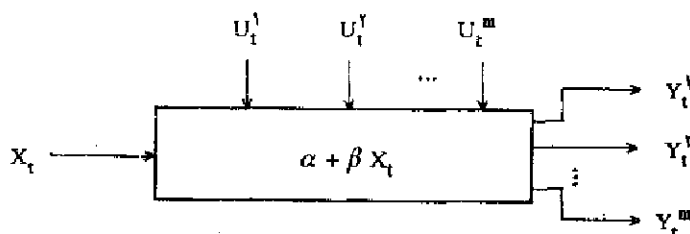
قبلاً فرض کرده ایم U_t دارای توزیع نرمال است؛ بنابراین با توجه به اینکه $\hat{\beta}$ یک تابع خطی از U_t است، $\hat{\beta}$ نیز دارای توزیع نرمال خواهد بود. سؤال این است که میانگین و واریانس $\hat{\beta}$ چیست؟

۱. میانگین $\hat{\beta}$

ابتدا سعی می‌کنیم مفهوم میانگین یا امید ریاضی $\hat{\beta}$ روشن شود. مدل رگرسیون ۱.۸ را دوباره می‌نویسیم:

$$Y_t = \alpha + \beta X_t + U_t$$

می‌دانیم X_t در آزمایشهای فرضی تکراری ثابت فرض می‌شود. برای تمام مقادیر t ، $t = 1, 2, 3, \dots, n$ ، مقدار X_t را ثابت نگاه داشته، m آزمایش انجام می‌دهیم. در هر آزمایش جمله اختلال، یک مقداری به خود می‌گیرد و با افزوده شدن به مقدار معین $(\alpha + \beta X_t)$ ، مقدار متغیر درون‌زا را مشخص می‌کند. این نکته‌ها را می‌توان در نمودار ۲-۱ ملاحظه کرد.



نمودار ۲-۱ مفهوم $E(\hat{\beta})$

به ازای X_t معین و در آزمایش اول، جمله اختلال مدل، مقدار U_t^1 را می‌گیرد. این مقدار به کمیت ثابت $(\alpha + \beta X_t)$ اضافه می‌شود و Y_t^1 را نتیجه می‌دهد. در آزمایش دوم، U_t مقدار U_t^2 را گرفته که بعد از ورود به سیستم، مقدار Y_t^2 به دست می‌آید. به همین ترتیب با تعیین U_t^m مقدار Y_t^m حاصل خواهد شد.

در اینجا به ازای سری مشاهدات X_t ، $t = 1, 2, \dots, n$ و سری مشاهدات Y_t برای تمام $t = 1, 2, \dots, n$ ، می‌توان یک رگرسیون ساخت و پارامتر β را تخمین زد. این

تخمین را $\hat{\beta}^1$ می‌نامیم. به همین ترتیب با استفاده از n مشاهده X_t و مشاهدات n تایی Y_t در آزمایش دوم (Y_t^2) می‌توان مدل رگرسیون را دوباره تخمین زد و این بار $\hat{\beta}^2$ را به دست آورد. این استدلال به همین روال ادامه پیدا می‌کند تا $\hat{\beta}^m$ حاصل شود. امید ریاضی $\hat{\beta}$ در واقع میانگین این m مقدار است که به طور فرضی برای $\hat{\beta}$ به دست آمده است. با افزایش m می‌توان امید ریاضی $\hat{\beta}$ را به صورت $E(\hat{\beta})$ تعریف کرد. بعد از آشنایی با مفهوم $E(\hat{\beta})$ به بررسی مقدار آن می‌پردازیم. معادله ۲-۲ را یک بار دیگر می‌نویسیم،

$$\hat{\beta} = \beta + \frac{\sum x_t U_t}{\sum x_t^2} .$$

می‌دانیم X_t در آزمایشهای تکراری، ثابت است؛ بنابراین x_t نیز همین خصوصیت را خواهد داشت. نتیجه می‌گیریم که $\sum x_t^2$ نیز برای تمام مقادیر t ($t = 1, 2, \dots, n$) در مجموع کمیت ثابتی است. این مقدار را با A نشان می‌دهیم؛ یعنی $\sum x_t^2 = A$. خواهیم داشت

$$\begin{aligned} \hat{\beta} &= \beta + \frac{x_1 U_1 + x_2 U_2 + \dots + x_n U_n}{A} , \\ &= \beta + \frac{x_1}{A} U_1 + \frac{x_2}{A} U_2 + \dots + \frac{x_n}{A} U_n . \end{aligned}$$

با توجه به اینکه مقادیر x_t برای $t = 1, 2, \dots, n$ ثابت است؛ بنابراین

$$\frac{x_t}{A} = \frac{x_t}{\sum x_t^2} = w_t = \text{ثابت} . \quad (2-3)$$

در نتیجه خواهیم داشت

$$\hat{\beta} = \beta + w_1 U_1 + w_2 U_2 + \dots + w_n U_n . \quad (2-4)$$

یعنی $\hat{\beta}$ یک تابع خطی از جمله‌های اختلال است.

از دو طرف معادله ۲-۴ امید ریاضی می‌گیریم. با توجه به ثابت بودن w_t داریم

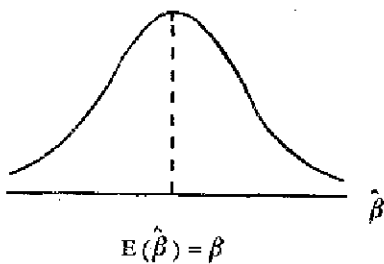
$$E(\hat{\beta}) = E(\beta) + w_1 E(U_1) + w_2 E(U_2) + \dots + w_n E(U_n) .$$

چون β ثابت است بنابراین $E(\hat{\beta}) = \beta$ ، با در نظر گرفتن $E(U) = 0$ برای تمام مقادیر $i = 1, 2, \dots, n$ خواهیم داشت

$$E(\hat{\beta}_{OLS}) = \beta \quad (2.5)$$

نتیجه می‌گیریم که میانگین یا امید ریاضی مقادیر مختلف $\hat{\beta}$ که در آزمایش‌های فرضی تکراری و با ثابت بودن X_i به دست می‌آید، برابر مقدار واقعی پارامتر جامعه، یعنی β است. یادآوری میشود که β واقعی قابل مشاهده نیست، همانگونه که انجام چنین آزمایش‌های فرضی نیز در عمل ممکن نخواهد بود. معادله ۲-۵ در واقع منعکس‌کننده یک خصوصیت آماری است و اساساً قابل محاسبه نیست.

در تمام مواردی که امید ریاضی یک تخمین‌زننده، برابر مقدار واقعی پارامتر باشد، می‌گویند آن تخمین‌زننده ناریب^۱ است. بنابراین $\hat{\beta}_{OLS}$ نیز یک تخمین‌زننده ناریب از β است. برای درک مفهوم خصوصیت ناریبی، کافی است به نمودار ۲-۲ توجه کنیم. می‌دانیم تابع توزیع احتمال $\hat{\beta}$ نرمال و میانگین این توزیع برابر $E(\hat{\beta})$ است. $\hat{\beta}$ موقعی ناریب است که این میانگین، یعنی $E(\hat{\beta})$ دقیقاً بر مقدار واقعی پارامتر، یعنی β منطبق باشد.



۲. میانگین $\hat{\alpha}$

در این قسمت، ثابت می‌کنیم که میانگین $\hat{\alpha}$ نیز برابر α است، یعنی $E(\hat{\alpha}) = \alpha$ ، که در واقع بیان دیگری از این واقعیت است که $\hat{\alpha}$ یک تخمین‌زننده ناریب از α است.

با توجه به معادله ۱-۲۶ می‌دانیم

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

می‌توان نوشت

$$\bar{Y} = \alpha + \beta \bar{X} + \bar{U} \quad (2.6)$$

معادله ۲-۶ را در ۱-۲۶ جایگزین می‌کنیم،

$$\hat{\alpha} = \alpha + \beta \bar{X} + \bar{U} - \hat{\beta} \bar{X} ,$$

$$\hat{\alpha} = \alpha + (\beta - \hat{\beta}) \bar{X} + \bar{U} . \quad (۲-۷)$$

از دو طرف رابطه فوق امید ریاضی می‌گیریم. با توجه به غیر تصادفی بودن X_i خواهیم داشت

$$E(\hat{\alpha}) = \alpha + \bar{X} E(\beta - \hat{\beta}) + E(\bar{U}) .$$

فرض صفر بودن امید ریاضی $E(U_i)$ ضرورتاً ایجاب می‌کند که $E(\bar{U})$ نیز صفر شود،

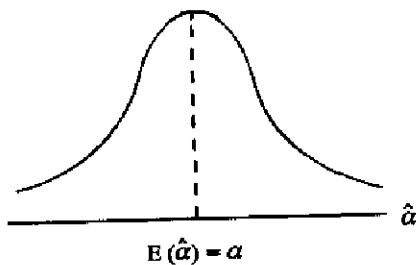
$$E(\hat{\alpha}) = \alpha + \bar{X} [E(\beta) - E(\hat{\beta})] + 0 ,$$

که با توجه به $E(\hat{\beta}) = \beta$ ، به صورت زیر نوشته می‌شود،

$$E(\hat{\alpha}) = \alpha + \bar{X} [\beta - \beta] + 0 ,$$

یا

$$E(\hat{\alpha}_{OLS}) = \alpha . \quad (۲-۸)$$



نمودار ۲-۳ خصوصیت نااریبی $\hat{\alpha}$

نمودار ۲-۳ نشان می‌دهد که

تخمین‌زننده $\hat{\alpha}$ نااریب است. ملاحظه می‌شود که مانند $\hat{\beta}$ ، تابع توزیع $\hat{\alpha}$ نیز نرمال است؛ زیرا با توجه به معادله ۱-۲۶ $\hat{\alpha}$ یک تابع خطی از $\hat{\beta}$ است، در نتیجه از توزیع $\hat{\beta}$ تبعیت می‌کند.

بعد از آشنایی با میانگین تخمین

زننده‌های $\hat{\alpha}$ و $\hat{\beta}$ ، به بررسی واریانس آنها می‌پردازیم. مطابق معمول بحث را با واریانس $\hat{\beta}$ آغاز می‌کنیم.

۳. واریانس $\hat{\beta}$

با استفاده از تعریف واریانس، می‌دانیم

$$\text{Var}(\hat{\beta}) = E[\hat{\beta} - E(\hat{\beta})]^2,$$

در معادله ۲-۵ دیدیم که، $E(\hat{\beta}) = \beta$ است. با جایگزینی ۲-۵ در معادله فوق خواهیم داشت

$$\text{Var}(\hat{\beta}) = E(\hat{\beta} - \beta)^2. \quad (2-9)$$

در اینجا باید مقدار $(\hat{\beta} - \beta)$ را به دست آوریم. در معادله ۲-۴ دیدیم که $\hat{\beta}$ یک تابع خطی از جمله‌های اختلال است. این معادله را یک بار دیگر می‌نویسیم،

$$\hat{\beta} = \beta + w_1 U_1 + w_2 U_2 + \dots + w_n U_n,$$

یا

$$(\hat{\beta} - \beta) = w_1 U_1 + w_2 U_2 + \dots + w_n U_n. \quad (2-10)$$

با جایگزینی ۲-۱۰ در ۲-۹ خواهیم داشت

$$\text{Var}(\hat{\beta}) = E(w_1 U_1 + w_2 U_2 + \dots + w_n U_n)^2.$$

طرف راست را بسط می‌دهیم،

$$\begin{aligned} \text{Var}(\hat{\beta}) = E & (w_1^2 U_1^2 + w_2^2 U_2^2 + \dots + w_n^2 U_n^2 + 2 w_1 w_2 U_1 U_2 + 2 w_1 w_3 U_1 U_3 \\ & + \dots + 2 w_1 w_n U_1 U_n + \dots) \end{aligned}$$

در معادله ۲-۳ دیدیم که w_i مقدار ثابتی است، بنابراین

$$\begin{aligned} \text{Var}(\hat{\beta}) = w_1^2 E(U_1)^2 + w_2^2 E(U_2)^2 + \dots + w_n^2 E(U_n)^2 + \\ 2 w_1 w_2 E(U_1 U_2) + 2 w_1 w_3 E(U_1 U_3) + \dots + \\ 2 w_1 w_n E(U_1 U_n) + \dots. \end{aligned} \quad (2-11)$$

به تعاریف واریانس و کوواریانس جمله اختلال برمی گردیم،

$$\text{Var}(U_t) = E[U_t - E(U_t)]^2 .$$

با توجه به $E(U_t) = 0$ ، نتیجه می گیریم که

$$\text{Var}(U_t) = E(U_t)^2 . \quad (2-12)$$

همچنین بنا بر تعریف می دانیم

$$\text{Cov}(U_t, U_s) = E[U_t - E(U_t)][U_s - E(U_s)] ,$$

و با توجه به $E(U_t) = 0$ ، برای تمام مقادیر t ، خواهیم داشت

$$\text{Cov}(U_t, U_s) = E(U_t U_s) . \quad (2-13)$$

معادله های ۲-۱۲ و ۲-۱۳ را در معادله ۲-۱۱ جایگزین می کنیم،

$$\begin{aligned} \text{Var}(\hat{\beta}) &= w_1^t \text{Var}(U_1) + w_1^t \text{Var}(U_2) + \dots + w_n^t \text{Var}(U_n) + \\ &+ 2 w_1 w_2 \text{Cov}(U_1, U_2) + 2 w_1 w_3 \text{Cov}(U_1, U_3) + \\ &\dots + 2 w_1 w_n \text{Cov}(U_1, U_n) + \dots . \end{aligned}$$

فرض بر این است که U_t کلاسیک است، یعنی از خصوصیات واریانس همسانی و عدم خودهمبستگی برخوردار است. به عبارت دیگر

$$\text{Var}(U_t) = \sigma^2 , \quad t = 1, 2, \dots, n ,$$

$$\text{Cov}(U_t, U_s) = 0 , \quad t, s = 1, 2, \dots, n , \quad t \neq s .$$

بنابراین خواهیم داشت

$$\text{Var}(\hat{\beta}) = w_1^t \sigma^2 + w_1^t \sigma^2 + \dots + w_n^t \sigma^2 + 2 w_1 w_1 (0) +$$

$$\begin{aligned} & \gamma w_1 w_r (\cdot) + \dots + \gamma w_1 w_n (\cdot) + \cdot + \cdot + \dots + \cdot \\ & = \sigma^2 \sum w_i^2 . \end{aligned} \quad (2-14)$$

از معادله ۲-۳ داریم

$$w_i^2 = \frac{x_i^2}{A^2} ,$$

یا

$$\sum w_i^2 = \sum \frac{1}{A^2} x_i^2 .$$

با توجه به تعریف $A = \sum x_i^2$ و نظر به اینکه A مقدار ثابتی است، خواهیم داشت

$$\sum w_i^2 = \frac{1}{(\sum x_i^2)^2} \sum x_i^2 = \frac{1}{\sum x_i^2} . \quad (2-15)$$

معادله ۲-۱۵ را در ۲-۱۴ جایگزین کرده، داریم

$$\text{Var}(\hat{\beta}_{OLS}) = \frac{\sigma^2}{\sum x_i^2} . \quad (2-16)$$

۴. واریانس $\hat{\alpha}$

از تعریف واریانس $\hat{\alpha}$ شروع می‌کنیم،

$$\text{Var}(\hat{\alpha}) = E [\hat{\alpha} - E(\hat{\alpha})]^2 ,$$

که با توجه به معادله ۲-۸ خواهیم داشت

$$\text{Var}(\hat{\alpha}) = E(\hat{\alpha} - \alpha)^2 .$$

باید به بررسی مقدار $(\hat{\alpha} - \alpha)$ بپردازیم. از معادله ۲-۷ داریم

$$(\hat{\alpha} - \alpha) = (\beta - \hat{\beta}) \bar{X} + \bar{U} ,$$

که در فرمول $\text{Var}(\hat{\alpha})$ به صورت زیر جایگزین می شود،

$$\begin{aligned}\text{Var}(\hat{\alpha}) &= E[(\beta - \hat{\beta}) \bar{X} + \bar{U}]^2, \\ &= E(\beta - \hat{\beta})^2 \bar{X}^2 + 2 E[(\beta - \hat{\beta}) \bar{X} \bar{U}] + E(\bar{U})^2, \\ &= \bar{X}^2 E(\beta - \hat{\beta})^2 + 2 \bar{X} E[(\beta - \hat{\beta}) \bar{U}] + E(\bar{U})^2.\end{aligned}$$

می دانیم بنا بر تعریف $\text{Var}(\hat{\beta}) = E(\hat{\beta} - \beta)^2$ ، بنابراین

$$\text{Var}(\hat{\alpha}) = \bar{X}^2 \text{Var}(\hat{\beta}) + 2 \bar{X} E[(\beta - \hat{\beta}) \bar{U}] + E(\bar{U})^2. \quad (2-17)$$

به بررسی مقدار $E[(\beta - \hat{\beta}) \bar{U}]$ می پردازیم، از معادله ۲-۴ داریم

$$(\hat{\beta} - \beta) = \sum w_i U_i,$$

که در آن $w_i = \frac{x_i}{\sum x_i^2}$ یادآوری می کنیم که چون ثابت $\sum x_i^2 = A$ ؛ بنابراین

$$\sum w_i = \frac{\sum x_i}{A} = \frac{\sum x_i}{\sum x_i^2},$$

و چون $\sum x_i = 0$ ، بنابراین

$$\sum w_i = 0. \quad (2-18)$$

در اینجا به خاصیت دیگری از w_i اشاره می شود. می توان نشان داد که $\sum w_i X_i = 1$. کافی است دو طرف معادله ۲-۳ را در X_i ضرب کرده،

$$w_i X_i = \frac{x_i X_i}{A},$$

یا

$$\sum w_i X_i = \frac{\sum x_i X_i}{A} = \frac{\sum x_i^2}{A} = \frac{\sum x_i^2}{\sum x_i^2},$$

$$\sum w_t X_t = 1 .$$

بعد از ذکر این مقدمه به بررسی مقدار $\bar{U} (\beta - \hat{\beta})$ می‌رسیم، می‌دانیم

$$\begin{aligned} (\beta - \hat{\beta}) \bar{U} &= \left(\sum_{t=1}^n w_t U_t \right) \frac{1}{n} \sum_{t=1}^n U_t = \frac{1}{n} \sum w_t U_t \sum U_t , \\ &= \frac{1}{n} [(w_1 U_1 + w_2 U_2 + \dots + w_n U_n) (U_1 + U_2 + \dots + U_n)] , \\ &= \frac{1}{n} [U_1 (w_1 U_1 + w_2 U_2 + \dots + w_n U_n) + U_2 (w_1 U_1 + \dots + w_n U_n) + \dots] . \end{aligned}$$

در نتیجه با گرفتن امید ریاضی خواهیم داشت

$$\begin{aligned} E(\beta - \hat{\beta}) \bar{U} &= \frac{1}{n} [w_1 E(U_1)^2 + w_2 E(U_1 U_2) + \dots + w_n E(U_1 U_n) + \\ &\quad w_1 E(U_1 U_2) + w_2 E(U_2)^2 + \dots + w_n E(U_2 U_n) + \dots] . \end{aligned}$$

با فرض واریانس همسانی $\sigma^2 = E(U_t)^2$ و عدم خودهمبستگی $E(U_t U_s) = 0$ داریم

$$\begin{aligned} E(\beta - \hat{\beta}) \bar{U} &= \frac{1}{n} [w_1 \sigma^2 + \dots + \dots + \dots + w_2 \sigma^2 + \dots + \dots + \dots \\ &\quad + w_n \sigma^2] , \\ &= \frac{1}{n} [\sum w_t \sigma^2] . \end{aligned}$$

با توجه به ثابت بودن σ^2 داریم

$$E(\beta - \hat{\beta}) \bar{U} = \frac{1}{n} \sigma^2 \sum w_t ,$$

و با استفاده از معادله ۲-۱۸ داریم،

$$E(\beta - \hat{\beta}) \bar{U} = 0 . \quad (2-19)$$

مقدار $E(\bar{U})^t$ را حساب می‌کنیم،

$$\begin{aligned} E(\bar{U})^t &= E\left(\frac{\sum U_i}{n}\right)^t = \frac{1}{n^t} E(\sum U_i)^t, \\ &= \frac{1}{n^t} E(U_1 + U_2 + \dots + U_n)^t, \\ &= \frac{1}{n^t} E(U_1^t + U_2^t + \dots + U_n^t + 2U_1 U_2 + 2U_1 U_3 + \dots). \end{aligned}$$

طبق فرض واریانس همسانی و عدم خودهمبستگی، داریم

$$= \frac{1}{n^t} (\sigma^t + \sigma^t + \dots + \sigma^t + 0 + 0 + \dots + 0).$$

در نتیجه

$$E(\bar{U})^t = \frac{\sigma^t}{n}. \quad (2.20)$$

معادله‌های ۲-۱۹ و ۲-۲۰ را در ۲-۱۷ جایگزین می‌کنیم،

$$\text{Var}(\hat{\alpha}) = \bar{X}^t \text{Var}(\hat{\beta}) + 0 + \frac{\sigma^t}{n}.$$

معادله ۲-۱۶ را در معادله فوق جایگزین می‌کنیم خواهیم داشت

$$\text{Var}(\hat{\alpha}) = \bar{X}^t \frac{\sigma^t}{\sum x_i^t} + \frac{\sigma^t}{n},$$

در نتیجه

$$\text{Var}(\hat{\alpha}) = \sigma^t \left[\frac{1}{n} + \frac{\bar{X}^t}{\sum x_i^t} \right]. \quad (2.21)$$

در بعضی موارد، واریانس $\hat{\alpha}$ را به صورت دیگری نیز ارائه می‌کنند. از معادله

۲-۲۱ داریم

$$\text{Var}(\hat{\alpha}) = \frac{n \bar{X}' \sigma^2 + \sigma^2 \sum x_i^2}{n \sum x_i^2} = \frac{\sigma^2}{n \sum x_i^2} [n \bar{X}' + \sum x_i^2] .$$

نشان می‌دهیم که $n \bar{X}' + \sum x_i^2 = \sum X_i^2$. برای این منظور $\sum x_i^2$ را چنین می‌نویسیم،

$$\begin{aligned} \sum x_i^2 &= \sum (X_i - \bar{X})^2 = \sum X_i^2 + n \bar{X}' - 2 \bar{X} \sum X_i , \\ &= \sum X_i^2 - n \bar{X}' . \end{aligned}$$

بدین ترتیب

$$\begin{aligned} n \bar{X}' + \sum x_i^2 &= n \bar{X}' + \sum X_i^2 - n \bar{X}' , \\ &= \sum X_i^2 . \end{aligned}$$

رابطه فوق را در فرمول $\text{Var}(\hat{\alpha})$ جایگزین می‌کنیم،

$$\text{Var}(\hat{\alpha}) = \frac{\sigma^2}{n \sum x_i^2} [\sum X_i^2] ,$$

یا

$$\text{Var}(\hat{\alpha}) = \sigma^2 \left(\frac{\sum X_i^2}{n \sum x_i^2} \right) . \quad (\text{a ۲-۲۱})$$

۵. روش دیگری برای استخراج واریانسهای $\hat{\alpha}$ و $\hat{\beta}$

با اینکه روش فوق در به دست آوردن واریانس $\hat{\beta}$ و $\hat{\alpha}$ ، یعنی فرمولهای ۲-۱۶ و ۲-۲۱ تا حدی طولانی است، از نظر مفهومی بسیار جالب و هماهنگ با روش استخراج واریانس پارامترها در مدل رگرسیون چند متغیره است. البته روش ساده تری نیز وجود دارد که در زیر به آن اشاره می‌کنیم. این روش به این بستگی دارد که ابتدا ثابت کنیم $\hat{\alpha}$ و $\hat{\beta}$ یک ترکیب خطی از مشاهدات Y_i است. این قضیه را ابتدا برای $\hat{\beta}$ و سپس برای $\hat{\alpha}$ ثابت می‌کنیم.

قضیه ۱. ثابت کنید $\hat{\beta}_{OLS}$ یک ترکیب خطی از Y_t است.
معادله ۱-۲۵ را می‌نویسیم،

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2} .$$

رابطه $y_t = Y_t - \bar{Y}$ را جایگزین می‌کنیم،

$$\hat{\beta} = \frac{\sum x_t (Y_t - \bar{Y})}{\sum x_t^2} = \frac{\sum x_t Y_t - \bar{Y} \sum x_t}{\sum x_t^2} .$$

با توجه به $\sum x_t = 0$ ، خواهیم داشت

$$\hat{\beta} = \frac{\sum x_t Y_t}{\sum x_t^2} . \quad (۲-۲۲)$$

طرف راست فرمول ۲-۲۲ را بسط داده، با فرض $\sum x_t^2 = A$ داریم

$$\hat{\beta} = \frac{x_1}{A} Y_1 + \frac{x_2}{A} Y_2 + \dots + \frac{x_n}{A} Y_n .$$

از معادله ۲-۳ استفاده می‌کنیم،

$$\hat{\beta}_{OLS} = w_1 Y_1 + w_2 Y_2 + \dots + w_n Y_n = \sum w_t Y_t . \quad (۲-۲۳)$$

با توجه به اینکه ضرایب w_t ثابت است، معادله ۲-۲۳ نشان می‌دهد که $\hat{\beta}_{OLS}$ یک ترکیب خطی از مشاهدات Y_t است.

قضیه ۲. ثابت کنید $\hat{\alpha}_{OLS}$ یک ترکیب خطی از Y_t است.
معادله ۱-۲۶ را دوباره می‌نویسیم،

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} .$$

طبق معادله ۲-۲۳ می‌دانیم، $\hat{\beta} = \sum w_t Y_t$ که در آن $w_t = \frac{x_t}{\sum x_t^2}$. بنابراین خواهیم داشت

$$\begin{aligned}\hat{\alpha} &= \bar{Y} - \bar{X} \sum w_i Y_i, \\ &= \frac{\sum Y_i}{n} - \bar{X} \sum w_i Y_i,\end{aligned}$$

و در نتیجه

$$\hat{\alpha}_{OLS} = \sum \left[\frac{1}{n} - \bar{X} w_i \right] Y_i. \quad (2-24)$$

معادله ۲-۲۴ نشان می‌دهد که $\hat{\alpha}$ یک ترکیب خطی از Y_i است که $(\frac{1}{n} - \bar{X} w_i)$ در واقع ضرایب ثابت این ترکیب خطی است. در زیر به بحث واریانس $\hat{\alpha}$ و $\hat{\beta}$ می‌پردازیم.

روشن دیگری برای واریانس $\hat{\beta}$

طبق معادله ۲-۲۳ می‌دانیم $\hat{\beta}$ یک تابع خطی از Y_i است،

$$\hat{\beta} = \sum w_i Y_i,$$

که در آن w_i ضرایب ثابت بوده و از Y_i مستقل هستند. از دو طرف معادله ۲-۲۳ واریانس می‌گیریم،

$$\text{Var}(\hat{\beta}) = \text{Var}(\hat{\beta} \sum w_i Y_i) = \sum w_i^2 \text{Var}(Y_i).$$

مطابق معادله ۱-۱۴ می‌دانیم، $\text{Var}(Y_i) = \sigma^2$. در نتیجه

$$\text{Var}(\hat{\beta}) = \sum w_i^2 \sigma^2,$$

$$= \sigma^2 \sum w_i^2.$$

با استفاده از معادله ۲-۱۵ خواهیم داشت

$$\text{Var}(\hat{\beta}_{OLS}) = \frac{\sigma^2}{\sum x_i^2},$$

که دقیقاً همان معادله ۲-۱۶ است.

روش دیگری برای واریانس $\hat{\alpha}$

طبق معادله ۲-۲۴ می‌دانیم $\hat{\alpha}$ یک تابع خطی از Y_i است،

$$\hat{\alpha} = \sum \left[\frac{1}{n} - \bar{X} w_i \right] Y_i .$$

از دو طرف این معادله واریانس می‌گیریم،

$$\begin{aligned} \text{Var}(\hat{\alpha}) &= \text{Var} \left[\sum \left(\frac{1}{n} - \bar{X} w_i \right) Y_i \right] , \\ &= \sum \left(\frac{1}{n} - \bar{X} w_i \right)^2 \text{Var} Y_i . \end{aligned}$$

از معادله ۱-۱۴ می‌دانیم $\text{Var}(Y_i) = \sigma^2$ در نتیجه

$$\text{Var}(\hat{\alpha}) = \sigma^2 \sum \left(\frac{1}{n^2} - \frac{2 \bar{X} w_i}{n} + \bar{X}^2 w_i^2 \right) .$$

با توجه به معادله ۲-۱۸ ($\sum w_i = 0$) و معادله ۲-۱۵ ($\sum w_i^2 = \frac{1}{\sum x_i^2}$) و با در نظر گرفتن $\sum \frac{1}{n^2} = \frac{1}{n}$ خواهیم داشت

$$\begin{aligned} \text{Var}(\hat{\alpha}) &= \sigma^2 \left[\frac{1}{n} - \frac{2 \bar{X}}{n} \sum w_i + \bar{X}^2 \sum w_i^2 \right] , \\ &= \sigma^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] , \end{aligned}$$

که دقیقاً همان معادله ۲-۲۱ است. مانند گذشته می‌توان به فرمول دیگری برای $\text{Var}(\hat{\alpha})$ نیز رسید. می‌دانیم

$$\text{Var}(\hat{\alpha}) = \sigma^2 \left[\frac{\sum x_i^2 + n \bar{X}^2}{n \sum x_i^2} \right] .$$

با توجه به

$$\sum x_i^2 = \sum (x_i - \bar{x})^2 + \sum x_i^2 - n \bar{x}^2 ,$$

خواهیم داشت

$$\text{Var}(\hat{\alpha}_{\text{OLS}}) = \sigma^2 \left(\frac{\sum X_i^2}{n \sum X_i^2} \right) .$$

با مفاهیم و فرمولهای میانگین و واریانس $\hat{\alpha}$ و $\hat{\beta}$ آشنا شدیم. قبل از پایان مبحث خصوصیات آماری تخمین‌زنده‌های حداقل مربعات معمولی، کوواریانس $\hat{\alpha}$ و $\hat{\beta}$ را مطرح نموده و به یک نکته در مورد نرمال بودن تابع توزیع احتمال $\hat{\alpha}$ و $\hat{\beta}$ اشاره می‌کنیم.

۶. کواریانس $\hat{\alpha}$ و $\hat{\beta}$

از تعریف کواریانس $\hat{\alpha}$ و $\hat{\beta}$ آغاز می‌کنیم،

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = E[\hat{\alpha} - E(\hat{\alpha})][\hat{\beta} - E(\hat{\beta})] .$$

با استفاده از معادله‌های ۲-۵ و ۲-۸ داریم

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = E(\hat{\alpha} - \alpha)(\hat{\beta} - \beta) .$$

از معادله ۲-۷ داریم

$$(\hat{\alpha} - \alpha) = (\hat{\beta} - \beta) \bar{X} + \bar{U} ,$$

در نتیجه

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = E[(\hat{\beta} - \beta) \bar{X} + \bar{U}](\hat{\beta} - \beta) ,$$

$$= E[(\hat{\beta} - \beta)(\hat{\beta} - \beta) \bar{X} + \bar{U}(\hat{\beta} - \beta)] ,$$

$$= \bar{X} E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)] + E[\bar{U}(\hat{\beta} - \beta)] ,$$

$$= -\bar{X} E(\hat{\beta} - \beta)^2 + E[\bar{U}(\hat{\beta} - \beta)] .$$

از معادله ۲-۱۹ می‌دانیم

$$E[\bar{U}(\hat{\beta} - \beta)] = 0 ,$$

بنابراین

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = -\bar{X} E(\hat{\beta} - \beta)' + 0,$$

و با توجه به $\text{Var}(\hat{\beta}) = E(\hat{\beta} - \beta)'(\hat{\beta} - \beta)$ و با استفاده از معادله ۲-۱۶ خواهیم داشت

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = -\bar{X} \text{Var}(\hat{\beta}),$$

یا

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = \frac{-\bar{X} \sigma^2}{\sum x_i^2}. \quad (2-25)$$

نتیجه می‌گیریم که $\hat{\alpha}$ و $\hat{\beta}$ دارای توزیع نرمال با خصوصیات زیر هستند،

$$\hat{\alpha}_{OLS} \sim N\left(\alpha, \frac{\sigma^2 \sum X_i^2}{n \sum x_i^2}\right), \quad (2-26)$$

$$\hat{\beta}_{OLS} \sim N\left(\beta, \frac{\sigma^2}{\sum x_i^2}\right), \quad (2-27)$$

یک بار دیگر، به بحث نرمال بودن تابع توزیع احتمال $\hat{\beta}$ و $\hat{\alpha}$ برمی‌گردیم. چون $\hat{\beta}$ در معادله ۲-۲۲ یک ترکیب خطی از Y_i است و با توجه به اینکه Y_i در معادله ۱-۸ تابع خطی از U_i است؛ اگر U_i دارای توزیع نرمال باشد، Y_i نیز نرمال بوده و در نتیجه توزیع $\hat{\beta}$ نرمال خواهد بود. $\hat{\alpha}$ نیز طبق معادله ۱-۲۶ تابع خطی از $\hat{\beta}$ است و بنابراین توزیع نرمال دارد. حال اگر توزیع U_i نرمال نباشد، Y_i توزیع نرمال نخواهد داشت؛ با وجود این، با استفاده از قضیه حد مرکزی^۱ می‌توان گفت که $\hat{\beta}$ برای نمونه‌های بزرگ توزیع نرمال دارد. قضیه حد مرکزی به طور بسیار ساده، یعنی اگر حجم یک نمونه به طور مرتب زیاد شود، توزیع میانگین این نمونه به سمت نرمال میل می‌کند. $\hat{\alpha}$ و $\hat{\beta}$ نیز طبق معادله‌های ۲-۲۳ و ۲-۲۴ یک نوع میانگین وزنی از Y_i هستند؛ زیرا در واقع یک ترکیب خطی از Y_i هستند. در مورد $\hat{\alpha}$ و $\hat{\beta}$ این وزنه‌ها به ترتیب w_i و $(\frac{1}{n} - \bar{X} w_i)$ هستند. بنابراین ملاحظه

می شود که روش حداقل مربعات معمولی خصوصیت بسیار برجسته‌ای دارد، بدین معنی که حتی اگر جمله‌های اختلال دارای توزیع نرمال نباشد، تخمین زنده‌های $\hat{\alpha}$ و $\hat{\beta}$ برای نمونه‌های بزرگ توزیع نرمال دارند.

۷. تخمین واریانس جمله اختلال

واریانس $\hat{\beta}$ (معادله ۲-۱۶)، واریانس $\hat{\alpha}$ (معادله ۲-۲۱) و کوواریانس $\hat{\alpha}$ و $\hat{\beta}$ (معادله ۲-۲۵)، همگی تابعی از واریانس جمله اختلال (U_i) هستند. به طور قطع ما مقدار واقعی واریانس U_i را نمی دانیم، چون مقادیر U_i قابل مشاهده نیستند؛ بنابراین باید واریانس U_i ($\sigma_{U_i}^2$) را تخمین بزنیم. یادآوری می کنیم که $\sigma_{U_i}^2$ اگر هم بدون اندیس نوشته شود به معنای واریانس U_i است. بدیهی است اگر σ^2 اندیس دیگری به غیر از U داشته باشد بر واریانس همان اندیس دلالت می کند.

برای محاسبه واریانس U_i ، ضرورتاً باید U_i را داشته باشیم. اگر مقادیر U_i موجود نباشد، باید متغیری که بیشترین نزدیکی را با U_i دارد، در نظر گرفته و واریانس آن را به عنوان تقریبی از واریانس U_i به حساب آوریم. در اقتصادسنجی معمول است که e_i را به عنوان چنین متغیری معرفی می کنند؛ زیرا از بین متغیرهای موجود در یک سیستم رگرسیون خطی، هیچ متغیر دیگری نزدیکتر به U_i یافت نمی شود.

بدین ترتیب، کافی است بعد از تخمین پارامترهای مدل رگرسیون، جمله‌های پسماند را به دست آورده، واریانس آن را محاسبه کنیم تا تخمینی از واریانس U_i به دست آید. بنا بر تعریف، واریانس e_i عبارت است از

$$\text{Var}(e_i) = s^2 = \frac{\sum (e_i - \bar{e})^2}{n}$$

با توجه به معادله ۱-۳۸ می دانیم $\bar{e} = 0$ ، پس

$$\text{Var}(e_i) = s^2 = \frac{\sum e_i^2}{n}$$

بنابراین، باید واریانس e_i را که به ترتیب فوق به دست آمده است به عنوان واریانس U_i

بپذیریم. اما در قسمت ۱-۵ و بدون اثبات، دیدیم که برای رسیدن به تخمین نااریب از واریانس U_i ، باید مجموع مربعات پسماند را بر درجات آزادی آن تقسیم کنیم،

$$\hat{\text{Var}}(U) = \hat{\sigma}^2 = \frac{\sum e_i^2}{n-2} .$$

یادآوری می‌کنیم که برای محاسبه $\sum e_i^2$ باید ابتدا $\hat{\alpha}$ و $\hat{\beta}$ به دست آید و این امر به آن معنی است که در محاسبه $\sum e_i^2$ دو درجه آزادی از دست داده‌ایم. به همین دلیل است که برای رسیدن به تخمین نااریب از σ^2 ، مجموع مربعات پسماند را به جای n ، بر $(n-2)$ تقسیم کرده‌ایم. در اینجا باید ثابت کنیم که $\frac{\sum e_i^2}{n-2}$ تخمین نااریبی از σ^2 است.

برای اثبات، به ترتیب زیر عمل می‌کنیم. ابتدا در تعریف e_i ، به جای Y_i و \hat{Y}_i مقادیر آن را قرار می‌دهیم،

$$\begin{aligned} e_i &= Y_i - \hat{Y}_i , \\ &= \alpha + \beta X_i + U_i - \hat{\alpha} + \hat{\beta} X_i , \\ &= U_i - (\alpha - \hat{\alpha}) - (\hat{\beta} - \beta) X_i . \end{aligned}$$

مقدار $(\alpha - \hat{\alpha})$ را به ترتیب زیر به دست می‌آوریم. می‌دانیم

$$\begin{aligned} \hat{\alpha} &= \bar{Y} - \hat{\beta} \bar{X} , \\ &= (\alpha + \beta \bar{X} + \bar{U}) - \hat{\beta} \bar{X} , \\ &= \alpha - (\hat{\beta} - \beta) \bar{X} + \bar{U} , \end{aligned}$$

در نتیجه

$$(\hat{\alpha} - \alpha) = \bar{U} - (\hat{\beta} - \beta) \bar{X} .$$

بدین ترتیب e_i را می‌توان به صورت زیر نوشت،

$$e_i = (U_i - \bar{U}) - (\hat{\beta} - \beta) x_i .$$

با اینکه $E(U_i) = 0$ ، اما توجه داریم که \bar{U} صفر نیست؛ زیرا میانگین نمونه است.

دو طرف رابطه فوق را مجذور کرده و جمع می‌کنیم،

$$\begin{aligned} \sum e_i^2 &= \sum (U_i - \bar{U})^2 + (\hat{\beta} - \beta) \sum x_i^2 - 2(\hat{\beta} - \beta) \sum (U_i - \bar{U}) x_i, \\ &= \sum U_i^2 - n \bar{U}^2 + (\hat{\beta} - \beta)^2 \sum x_i^2 - 2(\hat{\beta} - \beta) \sum U_i x_i. \end{aligned}$$

از دو طرف رابطه فوق امید ریاضی می‌گیریم،

$$\begin{aligned} E(e_i^2) &= E(\sum U_i^2) - n E(\bar{U}^2) + E(\hat{\beta} - \beta)^2 \sum x_i^2 \\ &\quad - 2 E[(\hat{\beta} - \beta) \sum U_i x_i]. \end{aligned}$$

اما می‌دانیم که

$$E(\sum U_i^2) = \sum E(U_i^2) = \sum E[U_i - E(U_i)]^2 = n \sigma_u^2,$$

$$E(\bar{U}^2) = \text{Var}(\bar{U}) = \frac{\sigma_u^2}{n},$$

$$E(\hat{\beta} - \beta)^2 = E[\hat{\beta} - E(\hat{\beta})]^2 = \text{Var}(\hat{\beta}) = \frac{\sigma_u^2}{\sum x_i^2}.$$

برای محاسبه جمله آخر، از معادله ۲-۱۰ استفاده می‌کنیم،

$$\begin{aligned} E[(\hat{\beta} - \beta) \sum U_i x_i] &= E[(\sum w_i U_i) (\sum U_i x_i)], \\ &= \sigma_u^2 \sum w_i x_i. \end{aligned}$$

با توجه به معادله‌های ۲-۳ و ۲-۱۹ می‌دانیم $\sum w_i x_i = 1$ ؛ بنابراین

$$E[(\hat{\beta} - \beta) \sum U_i x_i] = \sigma_u^2.$$

بدین ترتیب $E(e_i^2)$ به صورت زیر نوشته می‌شود،

$$E(\sum e_i^2) = n \sigma_u^2 - \sigma_u^2 + \sigma_u^2 - 2 \sigma_u^2 = (n - 2) \sigma_u^2,$$

$$E\left(\frac{\sum e_i^2}{n-2}\right) = \sigma_u^2 .$$

نتیجه می‌گیریم که برای تخمین واریانس جمله اختلال (σ_u^2) باید واریانس پسماندها، یعنی

$$\sigma_e^2 = \text{Var}(e_i) = \frac{\sum e_i^2}{n}$$

را به دست آورد. اما برای اینکه واریانس جمله‌های پسماند بتواند تخمین ناریبی از واریانس جمله اختلال را نتیجه دهد، باید مجموع مربعات پسماند ($\sum e_i^2$) را بر $n-2$ که در واقع همان درجات آزادی است تقسیم کنیم. بنابراین

$$s^2 = \text{Var}(e_i) = \sigma_e^2 = \hat{\sigma}_u^2 = \text{Var}(U_i) = \frac{\sum e_i^2}{n-2} . \quad (2-28)$$

در معادله ۱-۴۵ نیز جذر رابطه فوق را به دست آورده، آن را خطای معیار تخمین یا SEE نامیدیم. یا توجه به معادله ۲-۲۸ جذر واریانس U_i را می‌توان خطای معیار رگرسیون^۱ نیز نامید و آن را با SER (خطای معیار رگرسیون) نشان داد

$$\text{SER} = \sqrt{\text{Var}(U_i)} = \hat{\sigma}_u .$$

مثال ۲-۱ مدل رابطه تولید با ساعتهای نیروی کار. موضوع مثال ۱-۲ را در نظر می‌گیریم،

$$Q_i = \alpha + \beta L_i + U_i .$$

(الف) واریانس U_i را به دست آورید.

(ب) واریانس $\hat{\alpha}$ و واریانس $\hat{\beta}$ را تخمین بزنید.

(ج) خطای معیار $\hat{\alpha}$ و خطای معیار $\hat{\beta}$ را حساب کنید.

(د) کوواریانس $\hat{\alpha}$ و $\hat{\beta}$ را محاسبه کنید.

در مثال ۱-۲، تخمین مدل فوق را به صورت زیر به دست آوردیم،

$$\hat{Q}_i = 3/6 + 0/70 L_i ,$$

و با فرض $X_i = L_i$ و $Y_i = Q_i$ داریم

$$\hat{Y}_i = 3/6 + 0/70 X_i .$$

همچنین نتایج زیر را نیز به دست آوردیم،

$$\sum x_i^2 = 28 , \quad \bar{X} = 8 ,$$

$$\sum x_i y_i = 21 , \quad \bar{Y} = 9/6 .$$

در مثال ۱-۴ نیز نتایج زیر را داشتیم،

$$\sum y_i^2 = 30/4 , \quad \sum e_i^2 = 14/60 .$$

الف) برای تخمین واریانس جمله اختلال، از معادله ۲-۲۸ استفاده می‌کنیم،

$$\widehat{\text{Var}}(U_i) = \hat{\sigma}_U^2 = \frac{\sum e_i^2}{n-2} ,$$

$$\hat{\sigma}_U^2 = \frac{14/60}{6-2} = 1/84 .$$

ب) فرمول ۲-۲۱ را می‌نویسیم،

$$\widehat{\text{Var}}(\hat{\alpha}) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] .$$

چون مقدار σ^2 را نمی‌دانیم، باید از مقدار تخمین آن، یعنی $\hat{\sigma}^2$ استفاده کنیم. در آن صورت تخمین واریانس $\hat{\alpha}$ عبارت خواهد بود از

$$\widehat{\text{Var}}(\hat{\alpha}) = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] ,$$

یا

$$\widehat{\text{Var}}(\hat{\alpha}) = 1/84 \left[\frac{1}{6} + \frac{64}{28} \right] = 4/360 .$$

برای تخمین واریانس $\hat{\beta}$ نیز به همین ترتیب عمل می‌کنیم،

$$\begin{aligned}\widehat{\text{Var}}(\hat{\beta}) &= \frac{\hat{\sigma}^2}{\sum x_i^2}, \\ &= \frac{1/13}{28} = 0.065.\end{aligned}$$

ج) برای محاسبه خطای معیار $\hat{\alpha}$ و $\hat{\beta}$ کافی است از واریانس آنها جذر بگیریم،

$$\begin{aligned}\text{SE}(\hat{\alpha}) &= \sqrt{\widehat{\text{Var}}(\hat{\alpha})}, \\ &= \sqrt{4/360} = 2/0.9.\end{aligned}$$

$$\begin{aligned}\text{SE}(\hat{\beta}) &= \sqrt{\widehat{\text{Var}}(\hat{\beta})}, \\ &= \sqrt{0.065} = 0.256.\end{aligned}$$

معمولاً مقادیر خطای معیار را در پرانتز و زیر تخمین پارامتر مربوط به هر یک می‌نویسند؛ بنابراین

$$\hat{Q}_i = 3/6 + 0.75 L_i.$$

$$(2/0.9)(0.256)$$

د) فرمول کواریانس $\hat{\alpha}$ و $\hat{\beta}$ را از معادله ۲-۲۵ می‌نویسیم،

$$\begin{aligned}\text{Cov}(\hat{\alpha}, \hat{\beta}) &= \frac{-\bar{X} \hat{\sigma}^2}{\sum x_i^2}, \\ &= \frac{-8(1/13)}{28} = -0.023.\end{aligned}$$

در پایان، دو نکته را مطرح می‌کنیم: اول اینکه، خطای معیار تخمین پارامترها $\text{SE}(\hat{\alpha})$ یا $\text{SE}(\hat{\beta})$ را نباید با خطای معیار تخمین رگرسیون (SER)، که به آن خطای معیار تخمین (SEE) نیز گفته می‌شود، اشتباه کرد. مفهوم اول، پراکندگی $\hat{\alpha}$ و $\hat{\beta}$ را در حول α و β اندازه گیری می‌کند و مفهوم دوم، معیاری برای اندازه گیری پراکندگی جمله‌های پسماند

است. دوم اینکه، وقتی $Cov(\hat{\alpha} و \hat{\beta})$ منفی شود، می توان نتیجه گرفت که اگر $\hat{\alpha}$ بیشتر از اندازه تخمین زده شود، $\hat{\beta}$ کمتر از میزان تخمین می خورد و برعکس. این نتیجه در تمام مواردی که \bar{X} مثبت باشد صادق است؛ زیرا $\hat{\sigma}^2$ و $\sum x_i^2$ همواره مثبت هستند.

۲-۳. آزمون فرضیه برای هر یک از پارامترها

می دانیم $\hat{\alpha}$ و $\hat{\beta}$ دقیقاً برابر α و β واقعی نیست. سؤال این است که با داشتن $\hat{\alpha}$ و $\hat{\beta}$ چگونه می توان نسبت به مقادیر α و β واقعی استنتاج آماری کرد. باید پاسخ را در مباحث آزمون فرضیه^۱ یافت. با استفاده از توابع توزیع احتمال $\hat{\alpha}$ و $\hat{\beta}$ و با دانستن مقادیر میانگین و واریانس آنها به راحتی می توان فرضیه های مختلف درباره α و β واقعی را آزمون کرد. برای سهولت، بحث را در مورد β متمرکز می کنیم. به همین صورت نیز می توان برای α استنتاج کرد. در این قسمت به ترتیب آزمونهای Z و t برای α و β و آزمون F برای σ^2 مطرح می شود. فرض بر این است که خواننده با مفاهیم مقدماتی آزمونهای آماری به طور خلاصه آشنایی دارد.

۱. آزمون Z برای پارامترها

هرگاه تابع توزیع احتمال $\hat{\beta}$ نرمال باشد، می توان از آزمون Z استفاده کرد. برای استفاده از آزمون Z لازم است که واریانس جمله اختلال، یعنی $\sigma_{U_i}^2 = \text{Var}(U_i)$ را بدانیم. معمولاً مقدار واقعی واریانس U_i ، یعنی $\sigma_{U_i}^2$ ، مشخص نیست؛ بلکه تخمینی از آن در دست است. همان گونه که در بحث از آزمون t نشان خواهیم داد، موقعی می توان از تخمین واریانس U_i ، یعنی $\hat{\sigma}_{U_i}^2$ در آزمون Z استفاده کرد که حجم نمونه بزرگ باشد و معمولاً $n > 25$ را به عنوان نمونه های بزرگ می شناسیم. در ادامه بحث آزمون Z ، فرض بر این است که مقدار واقعی واریانس U_i را می دانیم.

در آزمون فرضیه های مختلف در مورد پارامترهای α یا β معمولاً با استفاده از

یک نمونه n تایی، β را تخمین زده و برای مثال، داریم: $\hat{\beta} = 2/5$. می دانیم مقدار واقعی β اصولاً باید با این مقدار متفاوت باشد. سؤال می کنیم که با $\hat{\beta} = 2/5$ ، آیا می توان این فرضیه را قبول کرد که β واقعی، برای مثال، برابر $2/25$ باشد؟ معمولاً فرضیه ای را که می خواهیم آزمون کنیم، فرضیه «صفر، پوچ، تهی یا عدم»^۱ نامیده و آن را با H_0 نشان می دهیم. در مقابل، خلاف آن فرضیه را فرضیه رقیب^۲ می خوانیم و آن را با H_1 مشخص می سازیم. در مورد این مثال داریم

$$H_0: \beta = 2/25, \quad H_1: \beta \neq 2/25.$$

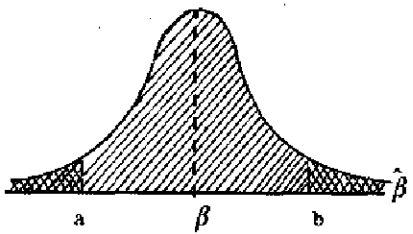
در حقیقت اگر این فرضیه را قبول کنیم یا به طور دقیقتر، نتوانیم رد کنیم که مقدار واقعی β بیشتر یا کمتر از $2/25$ است، فرضیه H_0 قابل قبول نخواهد بود. با توجه به اینکه در این مثال مقادیر بیشتر یا کمتر از $2/25$ ، مبنای استنتاج در آزمون H_0 است، می گویم که آزمون دو طرفه است. در موارد دیگر که فقط یک دامنه مورد توجه است، آزمون یکطرفه خواهد بود.

خصوصیت تصادفی بودن متغیرها اینجا می کند که در این گونه مسائل باید آزمون را در سطح احتمال معینی انجام داد؛ به عبارت دیگر، باید گفت که اگر مثلاً $\hat{\beta}$ برابر $2/5$ است با ۹۵ درصد احتمال این فرضیه را آزمون کنید که مقدار واقعی β برابر $2/25$ باشد.

ماهیت اساسی این گونه آزمونها چنین است. این سؤال که اگر $\hat{\beta}$ برابر $2/5$ باشد، آیا مقدار واقعی β برابر $2/25$ است یا خیر؛ به این بر می گردد که بگویم یک جامعه با پارامتری به نام β داریم که مقدار آن $2/25$ است. یک نمونه با آماره ای^۳ به نام $\hat{\beta}$ داریم که مقدار آن $2/5$ است. با احتمال ۹۵ درصد این فرضیه را آزمون کنید، که این نمونه متعلق به آن جامعه باشد.

۱. Null Hypothesis البته این اصطلاح، به جز در مورد $\beta = 0$ چندان روشن کننده نیست. متأسفانه فرضیه مثلاً $H_0: \beta = 20$ نیز فرضیه صفر نامیده می شود. این اصطلاح را نیمین و پیرسون - که از آماردانان معروفند - در دهه ۱۹۲۰ معرفی کرده اند.

می دانیم چون $\hat{\beta}$ از یک نمونه تصادفی به دست آمده، پس متغیری تصادفی با توزیع احتمال نرمال است که میانگین آن دقیقاً همان میانگین جامعه ($\beta = 2/25$) خواهد بود. این توزیع، در نمودار ۲-۴ رسم شده است. با توجه به اینکه این آزمون دو طرفه است، احتمال ۹۵ درصد به این



نمودار ۲-۴ مفهوم آزمون فرضیه

معنی است که باید ۹۵ درصد سطح زیر منحنی توزیع احتمال را در حول محور مرکزی ملاحظه کرد. سطح هاشور خورده که منعکس کننده این ۹۵ درصد است، نقاط a و b را مشخص می کند. اگر $\hat{\beta}$ به دست آمده از نمونه ($\beta = 2/5$) در فاصله b و a قرار گیرد، نتیجه می گیریم که این نمونه در واقع متعلق به جامعه ای است که میانگین آن $2/25$ است؛ بنابراین، $H_0: \beta = 2/25$ نمی تواند رد شود. در غیر این صورت اگر $\hat{\beta}$ در خارج از این فاصله افتاد، H_0 رد خواهد شد. به همین دلیل است که فاصله b و a را فاصله اعتماد یا فاصله اطمینان^۱ می گویند. نواحی خارج از این فاصله که با هاشور به صورت ضربدر مشخص شده است ناحیه بحرانی^۲ نامیده می شود؛ زیرا اگر $\hat{\beta}$ در آنجا واقع شود، H_0 قبول نخواهد شد.

بنابراین، مهمترین مسأله در آزمون فرضیه، پیدا کردن نقاط a و b است. برای یافتن این نقاط، کافی است از تابع توزیع احتمال $\hat{\beta}$ انتگرال گرفته و دو نقطه در روی محور $\hat{\beta}$ چنان تعیین کنیم که سطح زیر منحنی در حول محور مرکزی برابر $0/95$ شود. با تعیین این نقاط، نواحی بحرانی نیز مشخص می شود. یک روش ساده آزمون فرضیه، این است که اگر $\hat{\beta}$ در ناحیه بحرانی قرار گرفت، نتیجه می گیریم که نمی توانیم فرضیه H_0 را قبول کنیم.

بدیهی است این روش، مستلزم انتگرال گیری از تابع توزیع احتمال نرمال در هر

مسأله‌ای است که می‌خواهیم آن را حل کنیم. برای رفع این مشکل، کافی است تابع توزیع احتمال نرمال را استاندارد کنیم و برای یک بار و همیشه سطوح زیر منحنی آن را به ازای مقادیر مختلف، با روش انتگرال‌گیری محاسبه کنیم. این محاسبات که در حال حاضر با دقت بسیار انجام شده است در جدول Z و به صورت طبقه‌بندی شده، آماده است.^۱ در ذیل این نکته را بیشتر توضیح می‌دهیم.

می‌دانیم اگر یک متغیر تصادفی دارای توزیع نرمال باشد، می‌توان آن را به راحتی استاندارد کرد. برای مثال، در حالت کلی، اگر متغیر تصادفی X_i دارای توزیع احتمال نرمال و میانگین آن $E(X_i)$ و واریانس آن σ_x^2 باشد؛

$$X_i \sim N [E(X_i), \sigma_x^2] ,$$

در آن صورت می‌توان X_i را به متغیر استاندارد Z_i تبدیل کرد، به گونه‌ای که Z_i دارای توزیع احتمال نرمال با میانگین صفر و واریانس یک شود:

$$Z_i = \sim N (0, 1) .$$

برای این منظور کافی است متغیر تصادفی را از میانگین آن کم کرده، بر جذر واریانس (خطای معیار)^۱ تقسیم کنیم،

$$Z_i = \frac{X_i - E(X_i)}{\sigma_x} \sim N (0, 1) .$$

یادآوری می‌شود که X_i در این مثال تنها برای تبیین مفهوم استاندارد کردن یک متغیر تصادفی در حالت کلی به کار رفته است و هیچ رابطه‌ای با X_i به عنوان متغیر برون‌زا یا توضیحی مدل رگرسیون ندارد.

می‌دانیم $\hat{\beta}_{OLS}$ نیز یک متغیر تصادفی است که تابع توزیع احتمال دارد؛ بنابراین، اگر $\hat{\beta}$ را از میانگین آن کم و بر جذر واریانس (خطای معیار) تقسیم کنیم، کسر به دست

۱. به پیوست «آب» مراجعه شود.

آمده، دارای توزیع Z خواهد بود،

$$Z = \frac{\hat{\beta} - E(\hat{\beta})}{\sqrt{\text{Var}(\hat{\beta})}} = \frac{\hat{\beta} - E(\hat{\beta})}{SE(\hat{\beta})}, \quad (2.29)$$

که در آن $SE(\hat{\beta})$ همان خطای خیار $\hat{\beta}$ است. از معادله‌های ۲.۵ و ۲.۱۶ می‌دانیم

$$E(\hat{\beta}) = \beta, \quad \text{Var}(\hat{\beta}) = \frac{\sigma^2}{\sum x_i^2}.$$

با جایگزینی این مقادیر در معادله ۲.۲۹ خواهیم داشت

$$Z = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} = \frac{\hat{\beta} - \beta}{\sqrt{\sigma^2 / \sum x_i^2}}. \quad (2.30)$$

ملاحظه می‌شود Z در معادله ۲.۳۰ دارای توزیعی نرمال با میانگین صفر و واریانس یک است. برای نشان دادن این خصوصیت، باید $E(Z)$ و $\text{Var}(Z)$ را حساب کنیم،

$$E(Z) = \frac{1}{\sqrt{\sigma^2 / \sum x_i^2}} E(\hat{\beta} - \beta),$$

و چون

$$E(\hat{\beta} - \beta) = E(\hat{\beta}) - E(\beta) = \beta - \beta = 0,$$

بنابراین $E(Z) = 0$ است.

با گرفتن واریانس و توجه به این نکته که $\sum x_i^2$ و σ^2 هر یک کمیتهای ثابتی هستند خواهیم داشت

$$\text{Var}(Z) = \frac{1}{\sigma^2 / \sum x_i^2} \text{Var}(\hat{\beta} - \beta),$$

و چون

$$\text{Var}(\hat{\beta} - \beta) = \text{Var}(\hat{\beta}) + \text{Var}(\beta),$$

و با توجه به ثابت بودن مقدار β ، خواهیم داشت؛ $\text{Var}(\beta) = 0$ ، در نتیجه

$$\text{Var}(\hat{\beta} - \beta) = \text{Var}(\hat{\beta}) = \frac{\sigma^2}{\sum x_i^2}$$

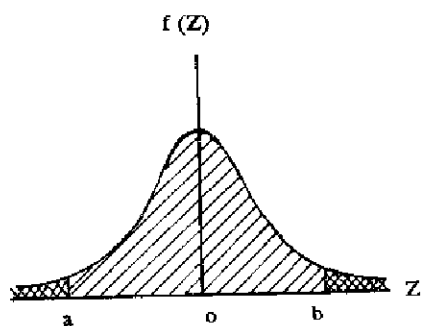
واریانس Z به صورت زیر خواهد بود

$$\text{Var}(Z) = \frac{1}{\sigma^2 / \sum x_i^2} \sigma^2 / \sum x_i^2 = 1$$

نتیجه کلی این است که $\hat{\beta}$ استاندارد دارای توزیع Z با میانگین صفر و واریانس یک است،

$$\frac{\hat{\beta} - \beta}{\text{SE}(\hat{\beta})} = Z \sim N(0, 1) \quad (2-21)$$

نمودار ۲-۵ توزیع احتمال Z را نشان می‌دهد. با توجه به اینکه این آزمون دو طرفه است، احتمال ۹۵ درصد بدین معنی خواهد بود که باید ۹۵ درصد سطح زیر



نمودار ۲-۵ توزیع متغیر استاندارد Z

متحنی توزیع احتمال را در حول محور مرکزی ملاحظه کرد. سطح هاشور خورده که منعکس کننده این ۹۵ درصد است، نقاط a و b را مشخص می‌کند. در اینجا باید دید، آیا $\hat{\beta}$ به دست آمده از نمونه در فاصله a و b قرار می‌گیرد یا خیر؟ اما نکته بسیار مهم این است که ابتدا باید $\hat{\beta}$ را استاندارد کرد تا بتواند با نقاط موجود دیگر روی محور Z

قابل مقایسه شود. اگر $\hat{\beta}$ استاندارد شده در فاصله a و b قرار گیرد، نتیجه می‌گیریم که این نمونه در واقع متعلق به جامعه‌ای است که میانگین آن صفر است؛ بنابراین فرضیه $H_0: \beta = 2/25$ نمی‌تواند رد شود. البته صفر بودن میانگین جامعه در یک مقیاس استاندارد به معنی $2/25$ بودن این میانگین در مقیاس عادی است. اگر $\hat{\beta}$ استاندارد شده در خارج از فاصله a و b قرار بگیرد، فرضیه H_0 قبول نخواهد شد. به همین دلیل است که فاصله a و b را فاصله اطمینان می‌گویند؛ مانند بحث قبل، نواحی خارج از این فاصله - که

با هاشور ضربدری مشخص شده است - نواحی بحرانی نامیده می‌شود؛ زیرا اگر $\hat{\beta}$ استاندارد شده در آنجا واقع شود، H_0 قبول نخواهد شد.

به این ترتیب مهمترین مسأله در آزمون فرضیه، یافتن نقاط b و a است و این امر با انتگرال‌گیری از تابع توزیع نرمال استاندارد، به سهولت به دست می‌آید. برخلاف گذشته لازم نیست برای هر مسأله، جداگانه انتگرال‌گیری شود. چون فقط و فقط یک تابع برای منحنی نرمال استاندارد وجود دارد و با توجه به اینکه میانگین و واریانس این تابع همواره صفر و یک است بنابراین کافی است به ازای مقادیر مختلف موجود در محور Z ، انتگرال‌گیری شده و سطوح مختلف زیر منحنی نرمال استاندارد محاسبه شود. این محاسبات در جدول Z طبقه‌بندی شده است.

بدین ترتیب در سطح احتمال ۹۵ درصد، نقاط a و b که قرینه نیز هستند از جدول Z به دست می‌آید. این مقادیر را « Z جدول» می‌نامیم. اکنون باید $\hat{\beta}$ را استاندارد کرد؛ یعنی آن را به متغیر Z تبدیل کرد. کافی است $\hat{\beta}$ را از میانگین آن کم کرده و بر خطای معیارش تقسیم کنیم؛ این مقدار را « Z محاسباتی» یا «آماره آزمون»^۱ می‌گوییم،

$$Z = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} \quad (2-32)$$

اگر آماره آزمون در فاصله مقادیر Z جدول قرار گرفت، نتیجه می‌گیریم که $\hat{\beta}$ به دست آمده از نمونه در واقع متعلق به جامعه‌ای است که میانگین آن در H_0 تحت آزمون قرار گرفته است. در چنین حالتی، فرضیه H_0 رد نمی‌شود. اگر آماره آزمون در ناحیه بحرانی قرار گرفت، H_0 رد خواهد شد، به عبارت دقیقتر قبول نمی‌شود. قبل از پایان مباحث این قسمت، ضروری است که به یک اصطلاح و یک نکته اشاره کنیم.

معمولاً برای هر آزمون ابتدا یک «اندازه یا سطح آزمون»^۲ تعیین می‌شود که آن را با α نشان می‌دهند. اندازه یا سطح آزمون ۱ درصد یا ۵ درصد بسیار معمول است.^۳ البته

1. Test Statistic

2. Size or Level of Test

۳. آزمون ۱ درصد و ۵ درصد بیشتر توسط فیشر (R. A. Fisher ۱۸۹۰ - ۱۹۶۲) از بین‌گذاران آمار جدید رایج شده است.

به جای لفظ اندازه یا سطح آزمون به طور خلاصه لفظ آزمون یا سطح نیز به کار می‌رود؛ یعنی گفته می‌شود آزمون ۱ درصد و آزمون ۵ درصد یا سطح ۱ درصد و سطح ۵ درصد. α در حقیقت سطح معنی‌داری آماری آزمون است و نشان می‌دهد چقدر احتمال دارد فرضیه H_0 را که در واقع صحیح است، رد کنیم؛ یعنی اشتباه نوع اول را مرتکب شویم. اصطلاح سطح معنی‌داری آماری را به طور خلاصه «سطح معنی‌دار» می‌گوییم. بنابراین اصطلاحات سطح معنی‌داری ۵ درصد، سطح ۵ درصد و آزمون ۵ درصد نه تنها همگی یک حقیقت را بیان می‌کنند بلکه با اصطلاحات سطح اطمینان ۹۵ درصد، فاصله اطمینان ۹۵ درصد و احتمال ۹۵ درصد - در واقع همان $(1 - \alpha)$ - نیز معادل است. ناگفته نماند که بعضی از متخصصان اقتصادسنجی به α ، اندازه یا سطح آزمون، ولی به $(1 - \alpha)$ سطح معنی‌داری آزمون می‌گویند؛ بنابراین در متون مختلف باید به تفاوت این اصطلاحات توجه کرد. البته همه متون، اصطلاح «فاصله اطمینان» را به کار می‌برند؛ مثلاً فاصله اطمینان ۹۵ درصد منعکس‌کننده این واقعیت است که مجموع سطوح نواحی بحرانی برابر ۵ درصد است.

با توجه به قرینگی منحنی تغییرات z و t ، محاسبات انتگرال، فقط برای نیمی از منحنی تغییرات انجام شده است؛ بنابراین در آزمونهای دو طرفه باید به مقدار $\frac{\alpha}{2}$ توجه کرد؛ مثلاً در مثال فوق و با توجه به نمودار ۲-۵ می‌توان گفت که چون مساحت سطح زیر منحنی برابر یک است؛ بنابراین a نقطه‌ای است که سطح زیر منحنی از α - تا آن نقطه برابر $0/025$ است. به همین ترتیب b نقطه‌ای است که سطح زیر منحنی از آن نقطه تا $+ \alpha$ ، برابر $0/025$ خواهد بود. مجموع این دو سطح، برابر $\alpha = 0/05$ می‌شود؛ به عبارت دیگر نقاط a و b متناظر با کمیت $Z_{\frac{\alpha}{2}}$ در جدول Z است؛ یعنی نقاطی است که سطح زیر منحنی از آن به بعد (یا قبل از آن) برابر $\frac{0/05}{2}$ است.

مثال ۲-۲ فرض کنید که در تخمین یک مدل رگرسیون خطی، $\hat{\beta}$ برابر $29/48$ شده است. همچنین فرض می‌کنیم که واریانس $\hat{\beta}$ را می‌دانیم، به گونه‌ای که جذر آن - که در واقع

همان خطای معیار است - برابر ۳۶ شده است. می خواهیم در سطح معنی دار ۵ درصد این فرضیه را آزمون کنیم که β واقعی برابر با ۲۵ است.

حل داریم:

$$\hat{\beta} = ۲۹/۴۸ \quad H_0: \beta = ۲۵$$

$$SE(\hat{\beta}) = ۳۶ \quad H_1: \beta \neq ۲۵$$

ملاحظه می شود که آزمون ما دو طرفه است. ابتدا آماره Z را می سازیم؛ به عبارت دیگر $\hat{\beta}$ را استاندارد می کنیم که بتوان آن را روی محور Z مشخص کرد. کافی است $\hat{\beta}$ را از میانگین آن کم کرده و بر خطای معیارش تقسیم کنیم،

$$Z = \frac{\hat{\beta} - E(\hat{\beta})}{SE(\hat{\beta})}$$

می دانیم $E(\hat{\beta}) = \beta$ است و چون بنا بر فرضیه H_0 داریم $\beta = ۲۵$ ، بنابراین

$$Z = \frac{۲۹/۴۸ - ۲۵}{۳۶} = ۰/۱۲$$

در اینجا باید این آماره آزمون را با مقدار به دست آمده از جدول Z مقایسه کنیم اگر در ناحیه بحرانی قرار گرفت، فرضیه H_0 قبول می شود. چون می خواهیم در سطح معنی دار

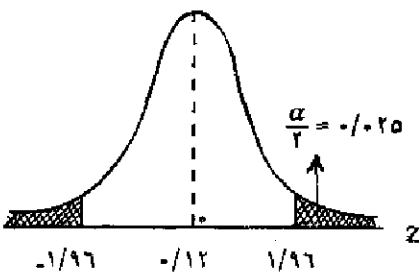
۵ درصد، یعنی به ازای $\alpha = ۵\%$ این فرضیه

را آزمون کنیم و با توجه به اینکه آزمون ما

دو طرفه است، باید در جدول Z ، مقدار Z

متناظر با $\frac{\alpha}{۲} = \frac{۰/۰۵}{۲} = ۰/۰۲۵$ را پیدا کنیم.

ملاحظه می شود که



$$Z_{0.025} = \pm 1/96$$

مفید است مقدار آماره آزمون را در

آزمون Z : مقایسه آماره آزمون و Z جدول

نمودار مشخص کنیم. آماره آزمون برابر

۰/۱۲ است؛ بنابراین در ناحیه بحرانی قرار نمی‌گیرد. نتیجه می‌گیریم که آزمون H_0 یعنی $\beta = 25$ نمی‌تواند رد شود.

فرضیه $\beta = 0$ در آزمون Z
مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \alpha + \beta X_i + U_i$$

اگر بخواهیم به این سؤال پاسخ دهیم که آیا X_i تأثیری بر Y_i دارد یا خیر، باید این فرضیه را آزمون کنیم که آیا $\beta = 0$ قابل قبول است؟ بنابراین فرضیه H_0 و فرضیه H_1 در این مورد چنین خواهد بود،

$$H_0: \beta = 0 \quad \text{و} \quad H_1: \beta \neq 0$$

در این حالت، مقدار $\hat{\beta}$ استاندارد شده یعنی آماره آزمون عبارت خواهد بود از

$$Z = \frac{\hat{\beta} - E(\hat{\beta})}{SE(\hat{\beta})} = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})}$$

و چون بنا بر فرضیه H_0 داریم، $\beta = 0$ ، آنگاه

$$Z = \frac{\hat{\beta}}{SE(\hat{\beta})} \quad (2.23)$$

نتیجه می‌گیریم که نحوه عمل در آزمون فرضیه $\beta = 0$ بسیار ساده است. باید مقدار $\hat{\beta}$ را بر خطای معیار آن تقسیم کنیم تا آماره آزمون به دست آید. اگر آماره آزمون در ناحیه بحرانی قرار گرفت، معنی دار است و فرضیه H_0 رد می‌شود، در غیر این صورت نمی‌توان آن را رد کرد.

معمولاً، فرضیه $\beta = 0$ را می‌توان به روش ساده‌تر - که البته تقریبی است - نیز آزمود. می‌دانیم در سطح احتمال ۹۵ درصد، مقدار Z به دست آمده از جدول Z برابر با $\pm 1/96$ است. اگر این مقدار را برابر ± 2 بگیریم (که قاعدتاً متناظر با احتمال بیشتر از

۹۵ درصد نیز هست)، آنگاه هنگامی H_0 رد می‌شود که قدر مطلق آماره آزمون بیشتر از $|\pm 2|$ بشود،

$$\left| \frac{\hat{\beta}}{SE(\hat{\beta})} \right| > |\pm 2|$$

به عبارت دیگر، بدون توجه به علامت، مقدار $\hat{\beta}$ باید از دو برابر خطای معیار خود بزرگتر باشد،

$$|\hat{\beta}| > 2 SE(\hat{\beta}) \quad (2-34)$$

در صورت درست بودن نامساوی (۲-۳۴)، فرضیه H_0 رد می‌شود، یعنی نتیجه می‌گیریم که متغیر توضیحی X_1 در واقع تأثیر قابل ملاحظه‌ای بر تغییرات Y_1 دارد.

مثال ۲-۳ فرض کنید با استفاده از یک نمونه ۵۰۰ تایی مشاهدات درآمد و مصرف، تابع مصرف زیر را تخمین زده‌ایم،

$$\hat{C}_1 = 250/7 + 0/84 Y_1$$

می‌دانیم انحراف معیار میل نهایی به مصرف برابر ۰/۱۴ است $(SE(\hat{\beta}) = 0/14)$

که β میل نهایی به مصرف فرض شده است. می‌خواهیم در سطح معنی دار ۵ درصد این فرضیه را آزمون کنیم که آیا درآمد، بر تغییرات مصرف تأثیری دارد یا خیر؟

به دو دلیل می‌توان از آزمون Z استفاد کرد: اولاً حجم مشاهدات بسیار است؛

ثانیاً، فرض شده است که مقدار واقعی واریانس تخمین پارامتر را می‌دانیم. البته وجود

یکی از این موارد، کافی است که به ما اجازه دهد از آزمون Z استفاده کنیم. فرضیه‌ای که

می‌خواهیم آزمون کنیم چنین است.

$$H_0: \beta = 0$$

$$H_1: \beta \neq 0$$

ابتدا باید آماره آزمون را تشکیل دهیم. با استفاده از معادله ۲-۳۳ داریم

$$Z = \frac{\hat{\beta}}{SE(\hat{\beta})} = \frac{0/84}{0/14} = 6$$

چون آزمون دو طرفه است، در سطح معنی دار ۵ درصد، یعنی به ازای

$$\frac{\alpha}{2} = \frac{0/05}{2} = 0/025$$

داریم

$$Z_{0/025} = \pm 1/96$$

بدیهی است که آماره آزمون به مراتب از این مقدار بیشتر می شود و معنی دار است.

در نتیجه فرضیه H_0 رد می شود؛ یعنی فرض تأثیر درآمد بر مصرف را نمی توان رد کرد.

راه حل بسیار ساده تر این است که چون مقدار تخمین $\hat{\beta}$ یعنی ۰/۸۴ از دو برابر

انحراف معیار آن بزرگتر است،

$$0/84 > 2 (0/14)$$

در نتیجه بنا بر رابطه ۲-۳۴ فرضیه H_0 رد خواهد شد.

فاصله اطمینان برای پارامترها در آزمون Z

مدل رگرسیون زیر را در نظر می گیریم،

$$Y_i = \alpha + \beta X_i + U_i$$

می دانیم $\hat{\beta}_{OLS}$ دقیقاً برابر مقدار واقعی β نخواهد بود. سؤال این است که آیا می توان با دانستن $\hat{\beta}$ ، درباره مقدار واقعی β اظهار نظر کرد؟ بدیهی است چون مقدار واقعی β را نمی دانیم، با آگاهی از $\hat{\beta}$ ، فقط می توان فاصله ای را تعیین کرد که با احتمال معینی انتظار می رود β واقعی در آن محدوده قرار گیرد. این محدوده را فاصله اعتماد یا فاصله اطمینان می گویند.

برای به دست آوردن این فاصله، کافی است $\hat{\beta}$ را به روش سابق استاندارد کنیم. $\hat{\beta}$

را از میانگین آن کم کرده و حاصل را بر خطای معیار $\hat{\beta}$ تقسیم کنیم. مقدار $\hat{\beta}$ استاندارد

شده عبارت خواهد بود از

$$Z = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})}$$

فاصله اطمینان را باید با یک احتمال معینی تخمین زد که به آن سطح معنی داری می‌گویند. فرض کنیم آزمون ما مثلاً α درصد است؛ پس سطح اطمینان آزمون $(1 - \alpha)$ خواهد بود. مطابق نمودار ۲-۵، هدف ما در واقع این است که حدودی مانند a و b را تعیین کنیم که آماره آزمون در آن فاصله قرار گیرد. نقطه a متناظر با $Z_{\frac{\alpha}{2}}$ و نقطه b متناظر با $-Z_{\frac{\alpha}{2}}$ است که از جدول Z به دست می‌آید؛ بنابراین آماره آزمون باید در فاصله بین این دو قرار گیرد، یعنی

$$-Z_{\frac{\alpha}{2}} \text{ جدول} < \text{آماره آزمون} < Z_{\frac{\alpha}{2}} \text{ جدول}$$

یا

$$-Z_{\frac{\alpha}{2}} < \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} < Z_{\frac{\alpha}{2}}$$

دو طرف نامساوی را در $SE(\hat{\beta})$ ضرب می‌کنیم؛

$$-Z_{\frac{\alpha}{2}} SE(\hat{\beta}) < \hat{\beta} - \beta < Z_{\frac{\alpha}{2}} SE(\hat{\beta})$$

از دو طرف $\hat{\beta}$ را کم کرده، بعد از ضرب در (-1) ، خواهیم داشت

$$\hat{\beta} - Z_{\frac{\alpha}{2}} SE(\hat{\beta}) < \beta < \hat{\beta} + Z_{\frac{\alpha}{2}} SE(\hat{\beta}) \quad (۲-۳۵)$$

نامساوی ۲-۳۵ فاصله اطمینان β واقعی را تعیین می‌کند. با داشتن فاصله اطمینان، می‌توان روش دیگری نیز برای آزمون فرضیه‌ها پیشنهاد کرد. کافی است بعد از ساختن فاصله اطمینان، ببینیم آیا فرضیه $\beta = k$: H_0 در این فاصله قرار می‌گیرد یا خیر. اگر $\beta = k$ در فاصله اطمینان واقع شد، در آن صورت می‌گوییم فرضیه H_0 رد نمی‌شود، در غیر این صورت نمی‌توان H_0 را قبول کرد.

مثال ۲-۴ در مثال ۲-۲، اولاً، برای β یک فاصله اطمینان ۹۵ درصد بسازید. ثانیاً، این

فرضیه را آزمون کنید که آیا β برابر ۲۵ است؟
اولاً، با استفاده از معادله ۲-۳۵ داریم:

$$29/48 - 1/96(36) < \beta < 29/48 + 1/96(36) ,$$

$$-41/08 < \beta < 100/04 .$$

ثانیاً، چون مقدار $\beta = 25$ در این فاصله اطمینان قرار می‌گیرد، بنابراین نمی‌تواند رد شود.

مثال ۲-۵ در مورد مثال ۲-۳، اولاً، برای میل نهایی به مصرف، یک فاصله اطمینان ۹۵ درصد بسازید، ثانیاً، فرضیه $\beta = 0$ را آزمون کنید.
اولاً، با استفاده از رابطه ۲-۳۵ داریم

$$0/84 - 1/96(0/14) < \beta < 0/84 + 1/96(0/14) ,$$

$$0/5656 < \beta < 1/1144 .$$

ثانیاً، چون مقدار $\beta = 0$ در این فاصله اطمینان قرار نمی‌گیرد، بنابراین فرضیه H_0 نمی‌تواند قبول شود.

مثال ۲-۶ در تخمین یک مدل رگرسیون، $\hat{\beta} = N/4$ و $SE(\hat{\beta}) = 2/2$ به دست آمده است. برای مقدار واقعی β یک فاصله اطمینان ۹۵ درصد بسازید.
می‌دانیم

$$\hat{\beta} - Z_{\alpha/2} SE(\hat{\beta}) < \beta < \hat{\beta} + Z_{\alpha/2} SE(\hat{\beta}) ,$$

با توجه به $Z_{\alpha/2} = \pm 1/96$ ، خواهیم داشت

$$-N/4 + 1/96(2/2) < \beta < N/4 + 1/96(2/2) ,$$

یا

$$4/1 < \beta < 12/7 .$$

یعنی $\hat{\beta} = 1/4$ که حاصل یک نمونه آماری است از جامعه‌ای گرفته شده است که مقدار واقعی β در آن با ۹۵ درصد احتمال بین $1/4$ و $12/7$ قرار دارد.

۲. آزمون t برای پارامترها

می‌دانیم آزمون هر فرضیه‌ای در مورد $\hat{\beta}$ مستلزم استاندارد کردن $\hat{\beta}$ است. $\hat{\beta}$ استاندارد شده (Z) را که از فرمول ۲-۳۰ به دست می‌آید، یک بار دیگر می‌نویسیم،

$$Z = \frac{\hat{\beta} - \beta}{\sqrt{\sigma^2 / \sum x_i^2}}$$

محاسبه Z مستلزم داشتن σ^2 است و معمولاً σ^2 را نمی‌دانیم؛ زیرا مقدار واقعی واریانس U_i یا واریانس Y_i مجهول است؛ بنابراین باید σ^2 را تخمین زد. در معادله ۲-۲۸ و در بحث مربوط به U_i و e_i ، تخمین واریانس U_i به صورت زیر به دست آمد،

$$\hat{\sigma}_U^2 = \frac{\sum e_i^2}{n - 2}$$

اگر در فرمول ۲-۳۰ به جای σ^2 ، مقدار تخمین آن ($\hat{\sigma}_U^2$) را قرار دهیم، در آن صورت آماره به دست آمده توزیع t با $(n - 2)$ درجه آزادی دارد،

$$\frac{\hat{\beta} - \beta}{\sqrt{\hat{\sigma}_U^2 / \sum x_i^2}} \sim t(n - 2)$$

با فرض $s^2 = \hat{\sigma}_U^2$ ، خواهیم داشت

$$\frac{\hat{\beta} - \beta}{\sqrt{s^2 / \sum x_i^2}} \sim t(n - 2) \quad (2-36)$$

برای تبیین فرمول ۲-۳۶ به طور خلاصه به بعضی مفاهیم آماری اشاره می‌کنیم. برای مباحث مقدماتی اقتصادسنجی همین تعاریف کافی است، ولی برای تحلیل دقیق قاعدتاً باید به کتابهای آمار مراجعه کرد. می‌دانیم اگر n متغیر تصادفی داشته باشیم که هر یک دارای توزیع نرمال مستقل با میانگین صفر و واریانس یک باشد آنگاه مجموع

مربعات آنها، توزیع «مربع کای»^۱ با n درجه آزادی خواهد داشت. معمولاً توزیع مربع کای با علامت χ^2 نشان می‌دهند؛^۲ بنابراین اگر $Z \sim N(0, 1)$ یک متغیر نرمال استاندارد باشد و n مقدار تصادفی z_1, z_2, \dots, z_n را از این توزیع به دست آوریم و هر یک را مجذور و نتایج را با هم جمع کنیم، آنگاه آماره^۳ به دست آمده توزیع χ^2 با n درجه آزادی خواهد داشت،^۲

$$(z_1^2 + z_2^2 + \dots + z_n^2) \sim \chi^2(n) . \quad (2.37)$$

وقتی درجات آزادی به سمت بی‌نهایت میل کند، توزیع χ^2 به سمت توزیع نرمال میل می‌کند. می‌توان نشان داد که $\frac{\sum_{i=1}^n e_i^2}{\sigma_u^2}$ دارای توزیع χ^2 با $(n - 2)$ درجه آزادی است،

$$\frac{\sum_{i=1}^n e_i^2}{\sigma_u^2} \sim \chi^2(n - 2) . \quad (2.38)$$

در اینجا فرض می‌کنیم که دو متغیر تصادفی با توزیعهای نرمال استاندارد و χ^2 به شرح زیر داریم

$$Z \sim N(0, 1) , \quad V \sim \chi^2(\nu) ,$$

که Z و V از یکدیگر مستقل و ν درجات آزادی توزیع V است. اگر متغیر نرمال استاندارد را بر جذر متغیری که توزیع χ^2 داشته باشد و بر درجات آزادی نیز تقسیم شده باشد، تقسیم کنیم آماره^۳ به دست آمده، توزیع t با ν درجه آزادی خواهد داشت؛ بنابراین

$$\frac{Z}{\sqrt{V/\nu}} \sim t(\nu) .$$

با توجه به معادله^۳ ۲-۳۰ می‌دانیم که مقدار $\hat{\beta}$ استاندارد، توزیع نرمال استاندارد دارد،

۱. Chi - square Distribution

۲. χ^2 بیست و دومین حرف از حروف الفبای یونانی است که صورت بزرگ آن χ است. به پیوست «ج» مراجعه شود. ۳. به پیوستهای «۲-ب» و «۳-ب» مراجعه شود.

$$Z = \frac{\hat{\beta} - \beta}{\sqrt{\sigma_u^2 / \sum x_i^2}} \sim N(0, 1) .$$

از طرف دیگر با توجه به معادله ۲-۳۸ نیز می‌دانیم که $\frac{\sum e_i^2}{\sigma_u^2}$ دارای توزیع χ^2 با $(n-2)$ درجه آزادی است؛ اگر $\frac{\sum e_i^2}{\sigma_u^2}$ را بر درجات آزادی آن تقسیم کنیم، خواهیم داشت

$$\frac{\sum e_i^2}{(n-2) \sigma_u^2} ,$$

و اگر Z به دست آمده از معادله ۲-۳۰ را بر آن تقسیم کنیم، آماره به دست آمده توزیع t به $(n-2)$ درجه آزادی خواهد داشت،

$$\left[\frac{(\hat{\beta} - \beta) \sqrt{\sum x_i^2}}{\sqrt{\sigma_u^2}} + \frac{\sqrt{\sum e_i^2}}{\sqrt{n-2} \sqrt{\sigma_u^2}} \right] \sim t(n-2) ,$$

یا

$$\frac{(\hat{\beta} - \beta)}{\sqrt{\frac{\sum e_i^2}{n-2} / \sum x_i^2}} \sim t(n-2) .$$

با جایگزینی معادله ۲-۲۸ خواهیم داشت

$$\frac{(\hat{\beta} - \beta)}{\sqrt{\hat{\sigma}_u^2 / \sum x_i^2}} \sim t(n-2) ,$$

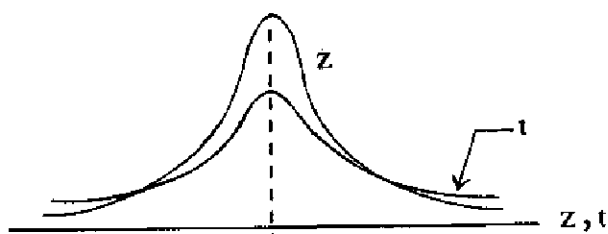
و با فرض $s^2 = \hat{\sigma}_u^2$ داریم

$$\frac{(\hat{\beta} - \beta)}{\sqrt{s^2 / \sum x_i^2}} \sim t(n-2) ,$$

که دقیقاً همان معادله ۲-۳۶ است. بنابراین نتیجه می‌گیریم که اگر در فرمول Z به جای σ_u^2 مقدار تخمین آن یعنی $\hat{\sigma}_u^2$ را قرار دهیم، توزیع t به دست می‌آید.^۱

۱. توزیع t را ویلیام گاست در سال ۱۹۰۸ در مجله بیومتریکا (*Biometrika*) منتشر کرد. گاست مقالات خود را با نام مستعار استیودنت (Student) به چاپ می‌رساند. به پیوست «د» مراجعه شود.

مانند توزیع نرمال، توزیع t نیز در حول میانگین صفر، قرینگی دارد. در مقایسه منحنی تغییرات t با Z می توان گفت که دنباله های منحنی t پهن تر است و این امر بویژه در حجم کم مشاهدات بسیار قابل ملاحظه است. به موازات افزایش حجم نمونه، منحنی تغییرات t به سمت Z میل می کند. معمولاً به ازای $n > 30$ تفاوت مقادیر t و Z بسیار کم و قابل اغماض است. به همین دلیل در مواردی که حجم مشاهدات بیشتر از ۳۰ است، علی رغم نداشتن مقدار واقعی واریانس U_1 ، می توان از جدول Z استفاده کرد. منحنی Z و منحنی t در نمودار ۲-۶ نشان داده شده است.



نمودار ۲-۶. نمایش منحنیهای Z ، t

آزمون فرضیه های مختلف درباره $\hat{\beta}$ در تابع توزیع احتمال t دقیقاً مشابه مباحثی است که برای Z مطرح کردیم. ابتدا باید فرضیه های H_0 و H_1 را تشکیل دهیم، آنگاه سطح معنی دار بودن آزمون را تعیین کنیم. $\hat{\beta}$ را استاندارد کرده، آماره آزمون را به دست می آوریم. برای به دست آوردن t از جدول، علاوه بر سطح معنی دار بودن آزمون، باید درجات آزادی را نیز در نظر بگیریم. در واقع عناصر موجود در جدول t بر حسب سطوح معنی داری و درجات آزادی، مدون شده است. همچنین با توجه به قرینگی منحنی تغییرات t ، فقط نیمی از محاسبات انتگرال انجام شده است. اگر قدر مطلق آماره آزمون از $|t|$ در جدول بزرگتر باشد در آن صورت آماره آزمون معنی دار است و فرضیه H_0 رد می شود؛ به عبارت دیگر هرگاه

$$\left| \frac{\hat{\beta} - \beta}{\sqrt{s^2 / \sum x_i^2}} \right| > |t| \quad (2-39)$$

۱. به پیوست (الف - ب) مراجعه شود.

باشد آنگاه H_0 رد خواهد شد، در غیر این صورت نمی توان H_0 را رد کرد.

مثال ۲-۷ مدل رابطه تولید با ساعتهای نیروی کار، موضوع مثال ۲-۱ را دوباره در نظر می گیریم؛

$$Q_i = \alpha + \beta L_i + U_i .$$

اولاً، فرضیه $\beta = 1$ را در مقابل فرضیه $\beta \neq 1$ در سطح معنی دار ۵ درصد آزمون کنید.

ثانیاً، در همین سطح معنی داری، فرضیه $\beta = 1$ را در مقابل فرضیه $\beta > 1$ آزمون کنید.

از مثال ۲-۱ این اطلاعات را داریم

$$\hat{\beta} = 0.75 \quad , \quad SE(\hat{\beta}) = 0.256 \quad ,$$

$$\hat{\sigma}^2 = 1/83 \quad .$$

با توجه به اینکه به جای مقدار واقعی واریانس U_i ، فقط تخمین آن را داریم؛ بنابراین $\hat{\beta}$ استاندارد توزیع t خواهد داشت،

$$t = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} ,$$

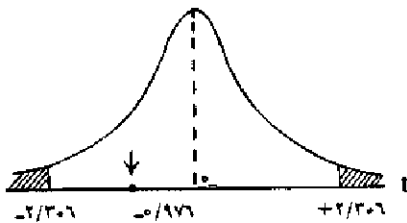
$$= \frac{0.75 - 1}{0.256} = -0.976 \quad .$$

اولاً، با حجم مشاهدات $n = 10$ و درجات آزادی $\nu = 10 - 2 = 8$ ، مقدار t را از جدول t و در سطح معنی دار ۵ درصد به دست می آوریم. چون آزمون دو طرفه است، باید t متناظر با سطح بحرانی $\frac{0.05}{2} = 0.025$ را ملاحظه کنیم، داریم

$$t_{0.025} = \pm 2.306 \quad .$$

باید قدر مطلق آماره آزمون را با قدر مطلق مقدار جدول t مقایسه کنیم. اگر آماره آزمون از مقدار جدول t بیشتر شد، معنی دار است و فرضیه H_0 رد می شود. ملاحظه می شود که

$2/306 < 0/976$ - است؛ بنابراین فرضیه H_0 رد نخواهد شد.



نمودار ۲-۷ نواحی بحرانی و آزمون t

با استفاده از نمودار نیز می توان آزمون را انجام داد. اگر آماره آزمون در ناحیه بحرانی قرار بگیرد، فرضیه H_0 رد می شود. نواحی بحرانی با استفاده از جدول t به دست می آید و به ترتیب با سطوح از $-\infty$ تا $-2/306$ و از $2/306$ تا $+\infty$ مشخص

می شود. چون آماره آزمون، یعنی $-0/976$ در هیچکدام از نواحی بحرانی واقع نشده است؛ پس فرضیه H_0 رد نمی شود.

ثانیاً، اگر بخواهیم فرضیه $\beta = 1$: H_0 را در مقابل فرضیه $\beta > 1$: H_1 آزمون کنیم، در آن صورت آزمون ما یکطرفه است. دیگر لازم نیست سطح معنی داری را نصف کنیم، بلکه با درجات آزادی $\nu = 8$ و $\alpha = 0/05$ ، مقدار t را از جدول به دست می آوریم؛ یعنی نقطه ای را در روی محور t مشخص می کنیم که سطح زیر منحنی از آن نقطه به بعد برابر ۵ درصد باشد و این دقیقاً همان ناحیه بحرانی است، خواهیم داشت $t = 1/860$. بنابراین چون آماره آزمون در ناحیه بحرانی قرار نمی گیرد، فرضیه H_0 رد نمی شود.

فرضیه $\beta = 0$ در آزمون t

در اینجا دقیقاً مانند مباحثی که در آزمون Z مطرح شد، عمل می کنیم. فرضیه H_0 و H_1 به صورت زیر خواهد بود؛

$$H_0: \beta = 0 \quad , \quad H_1: \beta \neq 0$$

$\hat{\beta}$ استاندارد شده (آماره آزمون) را به دست می آوریم؛

$$t = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})}$$

چون بنا بر فرضیه H_0 داریم، $\beta = 0$ پس

$$t = \frac{\hat{\beta}}{SE(\hat{\beta})} \quad (2-40)$$

بنابراین باید $\hat{\beta}$ را بر انحراف معیار آن تقسیم کنیم تا آماره آزمون به دست آید. سپس مقدار به دست آمده را با جدول t مقایسه می‌کنیم، اگر در ناحیه بحرانی قرار گرفت، معنی‌دار است و فرضیه H_0 رد می‌شود؛ در غیر این صورت نمی‌توان آن را رد کرد.

مثال ۲-۸ در مثال ۲-۷ فرضیه $H_0: \beta = 0$ را در مقابل فرضیه $H_1: \beta \neq 0$ در سطح معنی‌دار ۵ درصد آزمون کنید.

حل: با اینکه از لحاظ نظریه اقتصادی، این فرضیه مفهومی ندارد، برای بررسی چگونگی استفاده از فرمول ۲-۴۰ آماره آزمون را محاسبه می‌کنیم،

$$t = \frac{0/75}{0/256} = 2/929$$

چون آزمون دو طرفه است، بنابراین $t = \pm 2/306$. ملاحظه می‌شود که آماره آزمون در ناحیه بحرانی قرار می‌گیرد و فرضیه H_0 رد می‌شود؛ یعنی متغیر توضیحی از نظر آماری بر متغیر درون‌زا تأثیر می‌گذارد.

فاصله اطمینان برای پارامترها در آزمون t

همچنان که در مباحث آزمون Z گفتیم، با داشتن $\hat{\beta}$ می‌توان برای مقدار واقعی β یک فاصله اطمینان ساخت. اگر بخواهیم در سطح معنی‌دار α این فاصله را بسازیم کافی است مقدار t را از جدول برای سطح معنی‌دار $\frac{\alpha}{2}$ به دست آورده و آماره آزمون را در فاصله $\pm t$ جدول قرار دهیم،

$$t_{\frac{\alpha}{2}} \text{ جدول} < \text{آماره آزمون} < -t_{\frac{\alpha}{2}} \text{ جدول}$$

یعنی

$$-t_{\frac{\alpha}{2}} < \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} < t_{\frac{\alpha}{2}}$$

دو طرف نامساوی را در $SE(\hat{\beta})$ ضرب کرده و از دو طرف $\hat{\beta}$ را کم می‌کنیم. اگر نتیجه را در (-1) ضرب کنیم خواهیم داشت

$$\hat{\beta} - t_{\alpha} SE(\hat{\beta}) < \beta < \hat{\beta} + t_{\alpha} SE(\hat{\beta}) \quad (2-41)$$

نامساوی ۲-۴۱ فاصله اطمینان α درصد برای β واقعی است. با داشتن فاصله اطمینان، می‌توان روش دیگری نیز برای آزمون فرضیه‌ها پیشنهاد کرد. اگر $\beta = k$ در فاصله اطمینان واقع شد، می‌گوییم فرضیه H_0 رد نمی‌شود، در غیر این صورت H_0 را رد می‌کنیم.

مثال ۲-۹ در مثال ۲-۷ اولاً، برای α و β فاصله اطمینان ۹۵ درصد بسازید. ثانیاً، برای β فاصله اطمینان یکطرفه ۹۵ درصدی بسازید، ثالثاً، برای حالت یک، این فرضیه را آزمون کنید که آیا β برابر $1/12$ است.

اولاً، با استفاده از رابطه ۲-۴۱ داریم

$$0.75 - 2/306(0/256) < \beta < 0.75 + 2/306(0/256) ,$$

$$0/159 < \beta < 1/340 .$$

می‌توان راه حل فوق را به بیان دیگری نیز مطرح کرد. فاصله اطمینان ۹۵ درصد، یعنی فاصله‌ای که با احتمال ۹۵ درصد، مقدار آماره آزمون در نواحی بحرانی قرار نگیرد؛ به عبارت دیگر، آماره آزمون در فاصله مقادیر $\pm t$ به دست آمده از جدول t واقع شود؛ بنابراین

$$\Pr \left[-2/306 < \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} < 2/306 \right] = 0/95 ,$$

که در آن \Pr به معنای احتمال است. به ازای $SE(\hat{\beta}) = 0/256$ ، داریم

$$0/159 < \beta < 1/340 .$$

به همین ترتیب برای α عمل می‌کنیم.

$$\Pr[-2/306 < \frac{\hat{\alpha} - \alpha}{SE(\hat{\alpha})} < 2/306] = 0/90 ,$$

که به ازای $SE(\hat{\alpha}) = 2/09$ داریم:

$$-1/22 < \alpha < 8/42 .$$

ثانیاً، فاصله‌های اطمینان یکطرفه، به حد بالا یا حد پایین پارامترها مربوط است. باید نقاطی را در روی محور t پیدا کنیم که سطح زیر منحنی از آن به بعد یا قبل از آن ۹۵ درصد باشد. یا مراجعه به جدول t داریم

$$\Pr(t < 1/86) = 0/90 , \Pr(t > -1/86) = 0/90 .$$

بنابراین برای یک فاصله اطمینان یکطرفه، حد بالای β عبارت است از

$$\hat{\beta} + 1/86 SE(\hat{\alpha}) = 0/70 + 1/86(0/206) = 0/70 + 0/47 = 1/22$$

و در نتیجه فاصله اطمینان ۹۵ درصدی یکطرفه برابر است با $(-\infty, 1/22)$. به همین ترتیب حد پایین β را به دست می‌آوریم،

$$\hat{\beta} - 1/86 SE(\hat{\alpha}) = 0/70 - 1/86(0/206) = 0/70 - 0/47 = -0/27 ,$$

ملاحظه می‌شود که فاصله اطمینان ۹۵ درصد یکطرفه، برابر است با $(-0/27, +\infty)$.

ثالثاً، چون $\beta = 1/12$ در فاصله اطمینان β برای حالت اول یعنی

$$1/330 < \beta < 0/164 \text{ قرار می‌گیرد؛ بنابراین } H_0 \text{ را نمی‌توان رد کرد.}$$

مثال ۲.۱۰ تخمین تابع مصرف به شرح زیر مفروض است،

$$C_t = 13 + 0/89 Y_t^d ,$$

(0/6) (0/01)

که در آن C_t مصرف و Y_t^d درآمد قابل تصرف است. اعداد داخل پرانتز مقادیر انحراف

معیار تخمین پارامترها هستند. می‌دانیم $n = 10$.

اولاً، اگر ضریب ثابت را α بگیریم، فرضیه $\alpha = 0$: H_0 را در مقابل فرضیه $\alpha \neq 0$: H_1 در سطوح معنی دار ۱ درصد و ۵ درصد آزمون کنید.

ثانیاً، اگر میل نهایی به مصرف را β بگیریم، برای یک بار فاصله اطمینان ۹۵ درصد و بار دیگر فاصله اطمینان ۹۹ درصد بسازید.

اولاً، چون آزمون، دو طرفه است باید مقدار t را از جدول به ازای $\alpha = \frac{0.01}{2} = 0.005$ و درجه آزادی $8 = 10 - 2 = 8$ به دست آوریم. خواهیم داشت $t = \pm 2.350$ ، حال باید آماره آزمون را به دست آوریم،

$$t = \frac{\hat{\alpha} - 0}{SE(\hat{\alpha})} = \frac{13}{5.6} = 2.32$$

با توجه به اینکه این آماره در ناحیه بحرانی قرار نمی‌گیرد، فرضیه H_0 رد نمی‌شود. اما در سطح معنی دار ۵ درصد مقدار به دست آمده از جدول t برابر است با $t_{0.025} = 2.306$ و ملاحظه می‌شود که آماره آزمون در ناحیه بحرانی قرار می‌گیرد؛ در نتیجه فرضیه H_0 رد می‌شود. این مثال نشان می‌دهد که در یک سطح معینی از معنی داری، می‌توان فرضیه‌ای را قبول کرد، در حالی که در سطح دیگر رد خواهد شد.

ثانیاً، برای به دست آوردن فاصله اطمینان ۹۵ درصد برای β ، باید مقدار t را از جدول با ۸ درجه آزادی و به ازای سطح زیر منحنی ۰/۰۲۵ به دست آورد. خواهیم داشت $t = \pm 2.306$ بنابراین

$$0.189 - 2.306(0.01) < \beta < 0.189 + 2.306(0.01),$$

یا

$$0.166 < \beta < 0.193.$$

اما برای فاصله اطمینان ۹۹ درصد کافی است مقدار t را به ازای سطح زیر منحنی ۰/۰۰۵ حساب کنیم؛ خواهیم داشت: $t = \pm 2.350$. در نتیجه فاصله اطمینان برای β عبارت است از:

$$0.189 - 2.350(0.01) < \beta < 0.189 + 2.350(0.01),$$

یا

$$0.1856 < \beta < 0.9223 .$$

ملاحظه می شود که به ازای سطوح بالاتری از احتمال، فاصله اطمینان بیشتر می شود.

فاصله اطمینان و دقت تخمین پارامترها

در مثال ۲-۱۰ دیدیم که هر چه سطح احتمال بالاتر باشد، فاصله اطمینان بیشتر می شود. بدیهی است این نتیجه، چندان رضایت بخش نیست. حالت مطلوب این است که در سطوح بالای احتمال، فواصل کوچک اطمینان داشته باشیم. سؤال این است که چگونه ممکن است به این هدف نزدیک شویم. یک بار دیگر رابطه ۲-۴۱ را می نویسیم.

$$\hat{\beta} - t \text{ SE}(\hat{\beta}) < \beta < \hat{\beta} + t \text{ SE}(\hat{\beta}) .$$

چون مقادیر $\hat{\beta}$ و t خارج از کنترل ما هستند بنابراین برای اینکه فاصله اطمینان β کمتر شود باید $\text{SE}(\hat{\beta})$ کوچکتر شود. اما می دانیم $\text{SE}(\hat{\beta})$ در واقع جذر واریانس $\hat{\beta}$ است؛ بنابراین باید واریانس $\hat{\beta}$ به کمترین مقدار خود برسد تا فاصله اطمینان برای β حداقل شود. حال ببینیم چگونه می توان واریانس $\hat{\beta}$ را کوچک کرد.

می دانیم

$$\text{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2} ,$$

که در آن $\hat{\sigma}^2 = s^2 = \frac{\sum e_i^2}{n-2}$. برای اینکه واریانس $\hat{\beta}$ کوچک شود، باید صورت کسر؛ یعنی $\hat{\sigma}^2$ کم شود و مخرج کسر، یعنی $\sum x_i^2$ افزایش یابد. به نظر می رسد که ممکن است افزایش حجم مشاهدات صورت کسر $\text{Var}(\hat{\beta})$ را از طریق افزایش مخرج کسر $\frac{\sum e_i^2}{n-2}$ کوچکتر کند. افزایش حجم مشاهدات می تواند مخرج کسر $\text{Var}(\hat{\beta})$ را نیز از طریق افزایش $\sum_{i=1}^n x_i^2$ بزرگتر کند و در نتیجه واریانس $\hat{\beta}$ را کاهش دهد. چون افزایش حجم مشاهدات باعث بزرگتر شدن $\sum e_i^2$ خواهد شد - که خود $\hat{\sigma}^2$ را بیشتر می کند - بنابراین

تأثیر مستقیمی در افزایش واریانس $\hat{\sigma}^2$ نیز خواهد داشت. البته انتظار این است که خالص تأثیر افزایش حجم مشاهدات باعث کمتر شدن واریانس $\hat{\sigma}^2$ شود؛ می‌دانیم هرچه واریانس یک تخمین زنده کمتر باشد، دقت آن تخمین زنده بیشتر خواهد بود. اما مهمترین عاملی که می‌تواند واریانس $\hat{\sigma}^2$ را کاهش دهد، این است که از طریق پراکندگی بیشتر مشاهدات، مقدار $\sum x_i^2$ بیشتر شود. می‌دانیم

$$\sum x_i^2 = \sum (x_i - \bar{X})^2 + n\bar{X}^2 \quad (2-42)$$

در نتیجه هر چقدر مقادیر x_i پراکندگی بیشتری داشته باشد، تفاوت هر یک از آنها با میانگین \bar{X} کمیت بزرگتری خواهد بود. وقتی این تفاوتها، یعنی $(x_i - \bar{X})$ ، مجذور شود، به مراتب بزرگتر شده و در نتیجه حاصل جمع بزرگتری خواهد داد. بدین ترتیب، می‌توان گفت که بهترین روش در کاهش واریانس $\hat{\sigma}^2$ ، که خود موجب افزایش دقت تخمین شده و در نهایت فاصله اطمینان $\hat{\sigma}^2$ را محدودتر می‌کند، این است که پراکندگی مشاهدات متغیر برونزا یا توضیحی بیشتر باشد.

این نکته در مباحث اقتصادسنجی کاربردی، اهمیت فراوان دارد؛ برای مثال، وقتی می‌خواهیم تابع مصرف را تخمین بزنیم که در آن درآمد به منزله متغیر توضیحی در نظر گرفته شده است، نباید نمونه‌ای که می‌گیریم به گروه خاصی از افراد در طبقه معینی از درآمد منحصر باشد؛ زیرا در این صورت درآمد افراد موجود در نمونه، اختلاف چندانی با هم ندارد، بنابراین تفاوت بسیاری با میانگین نخواهد داشت. در نتیجه $\sum (x_i - \bar{X})^2$ بسیار کوچک خواهد بود. برعکس وقتی مشاهدات ما به تمام گروههای درآمدی مربوط باشد، اختلاف عناصر نمونه از میانگین کل بزرگ شده و $\sum x_i^2$ افزایش خواهد یافت. به همین دلیل نمونه‌گیریها نباید به مناطق جغرافیایی خاص یا به دوره‌هایی محدود شود که سیاستهای شدید کنترل درآمد بر روند تغییرات اقتصادی حاکم است. برای مثال، اگر می‌خواهیم تابع سرمایه‌گذاری در بخش خصوصی را تخمین بزنیم که قیمت به عنوان یک متغیر توضیحی فرض شده است در آن صورت نباید دامنه تغییرات قیمت را به دوره‌هایی محدود کرد که سیاست شدید کنترل قیمت‌ها اعمال می‌شود.

آزمونهایی که تا اینجا مطرح کردیم، به فرضیه‌هایی مربوط می‌شود که می‌تواند در مورد هر یک از پارامترها مطرح شود. از مرحله بعد باید از آزمونهایی صحبت کنیم که درباره تخمین کل مدل نظر می‌دهد؛ یعنی به اصطلاح معنی‌دار بودن مدل را آزمون می‌کند. اما قبل از طرح آن، به آزمون واریانس جمله اختلال می‌پردازیم.

۳. آزمون واریانس جمله اختلال

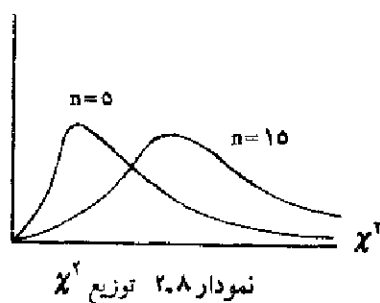
در مباحث آزمون t و در معادله ۲-۳۸ دیدیم که کمیت $\frac{\sum_{i=1}^n e_i^2}{\sigma^2}$ ، دارای توزیع χ^2 با $(n - 2)$ درجه آزادی است،

$$\frac{\sum_{i=1}^n e_i^2}{\sigma^2} \sim \chi^2 (n - 2)$$

بنابراین با استفاده از جدول χ^2 می‌توان در سطح معنی‌دار α ، برای $\frac{\sum e_i^2}{\sigma^2}$ یک فاصله اطمینان به دست آورد.

توزیع χ^2 همان‌گونه که در نمودار ۲-۸ ملاحظه می‌شود، معمولاً از مرکز مختصات

احتمال χ^2

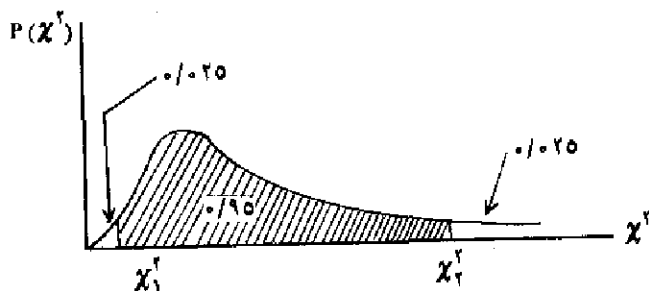


شروع شده و به سمت راست تا بی‌نهایت ادامه دارد. شکل دقیق منحنی تغییرات χ^2 تابعی از درجات آزادی است. به ازای افزایش درجات آزادی، قرینگی منحنی توزیع χ^2 به طور مرتب بیشتر می‌شود. وقتی مقدار درجات آزادی خیلی بزرگ شود، توزیع χ^2 به سمت توزیع نرمال میل

خواهد کرد. مقادیر مختلف χ^2 را می‌توان به ازای سطوح مختلف معنی‌داری و یا درجات آزادی از جدول χ^2 به دست آورد. چون توزیع χ^2 قرینگی ندارد، هر مقدار χ^2 که از جدول به دست آید، در واقع نقطه‌ای است که مساحت سطح زیر منحنی از آن نقطه به سمت راست برابر سطح معنی‌داری مورد نظر است؛ برای مثال، با درجه آزادی ۱۰، نقطه‌ای روی محور χ^2 که مساحت سطح زیر منحنی از آن نقطه به سمت راست ۵ درصد است

برابر $18/307$ خواهد بود. ^۱ در جدولهایی که برحسب توزیع تراکمی X^1 تنظیم شده‌اند معمولاً مساحت سطح زیر منحنی از مبدأ مختصات تا نقطه X^1 ، در ردیف اول از بالای جدول قرار دارد. مقدار X^1 که از متن جدول به دست می‌آید در واقع نقطه‌ای است که مساحت سطح زیر منحنی از مبدأ مختصات تا آن نقطه برابر سطح مندرج در ردیف اول است، برای مثال، به ازای درجه آزادی ۱۰ مقدار $18/3 = X^1$ نقطه‌ای روی محور X^1 است که مساحت سطح زیر منحنی از مبدأ مختصات تا آن نقطه برابر ۹۵ درصد است.

در اینجا فرض کنید می‌خواهیم یک فاصله اطمینان ۹۵ درصدی برای $\frac{\sum_{i=1}^n e_i^2}{\sigma^2}$ به دست آوریم؛ یعنی سطحی از منحنی تغییرات X^1 را مشخص کنیم که ۹۵ درصد در وسط و $0/025$ در هر طرف باشد. مقدار هاشور خورده نمودار ۹-۲ جواب مسأله است؛ بنابراین باید تقاطع X_1^1 و X_2^1 را با استفاده از جدول به دست آوریم.



نمودار ۲-۹ فاصله اطمینان ۹۵ درصد برای X^1

می‌توان از جدول توزیع تراکمی X^1 استفاده کرد. X_1^1 نقطه‌ای است که سطح زیر منحنی از مبدأ مختصات تا آن برابر $0/025$ است؛ بنابراین آن را با $X^1/0/025$ نشان می‌دهیم. به همین ترتیب X_2^1 نقطه‌ای است که انتگرال $P(X^1)$ از 0 تا X_2^1 برابر $0/975$ است. بنابراین می‌توان X_2^1 را با $X^1/975$ نشان داد. در نتیجه، فاصله اطمینان ۹۵ درصد برای $\frac{\sum_{i=1}^n e_i^2}{\sigma^2}$ برابر است با

۲. به پیوست «۳-ب» رجوع شود.

۱. به پیوست «۲-ب» رجوع شود.

$$\chi^2_{\alpha/2, n-2} < \frac{\sum e_i^2}{\sigma^2} < \chi^2_{1-\alpha/2, n-2}$$

به عبارت دقیقتر، احتمال اینکه $\frac{\sum e_i^2}{\sigma^2}$ بین $\chi^2_{\alpha/2, n-2}$ و $\chi^2_{1-\alpha/2, n-2}$ قرار گیرد، ۹۵ درصد است، یعنی

$$\Pr(\chi^2_{\alpha/2, n-2} < \frac{\sum e_i^2}{\sigma^2} < \chi^2_{1-\alpha/2, n-2}) = 0.95$$

اما هدف اصلی ما پیدا کردن فاصله اطمینان برای σ^2 است. به عبارت دیگر، وقتی در یک تخمین $\hat{\sigma}_U^2$ را به دست می آوریم، به طور طبیعی این سؤال مطرح می شود که آیا می توان با داشتن $\hat{\sigma}^2$ به حدود یا فاصله اطمینانی برای مقدار واقعی σ^2 رسید؟ برای پاسخ به این سؤال کافی است از معادله ۲-۲۸ استفاده کنیم، می دانیم

$$\hat{\sigma}_U^2 = \frac{\sum e_i^2}{n-2}$$

در نتیجه

$$\sum e_i^2 = (n-2) \hat{\sigma}^2 \quad (2-43)$$

با جایگزینی ۲-۴۳ در فاصله اطمینان ۹۵ درصد برای $\frac{\sum e_i^2}{\sigma^2}$ ، خواهیم داشت

$$\chi^2_{\alpha/2, n-2} < \frac{(n-2) \hat{\sigma}^2}{\sigma^2} < \chi^2_{1-\alpha/2, n-2}$$

دو طرف نامساوی را معکوس کرده، در $\hat{\sigma}^2 (n-2)$ ضرب می کنیم. خواهیم داشت

$$\frac{(n-2) \hat{\sigma}^2}{\chi^2_{1-\alpha/2, n-2}} < \sigma^2 < \frac{(n-2) \hat{\sigma}^2}{\chi^2_{\alpha/2, n-2}} \quad (2-44)$$

نامساوی ۲-۴۴ فاصله اطمینان ۹۵ درصد را برای σ^2 مشخص می کند.

در پایان به این نکته اشاره می کنیم که فاصله های اطمینان $\hat{\alpha}$ و $\hat{\beta}$ که قبلاً به دست

آوردیم - حول مقادیر α و β قرینگی دارند؛ در حالی که فاصله اطمینان σ^2 حول مقدار $\hat{\sigma}^2$

قرینه نیست. دلیل این امر عدم وجود قرینگی در توزیع χ^2 است؛ در حالی که توزیع Z و t هر کدام قرینگی کامل دارند.

مثال ۲-۱۱ با توجه به مدل تولید و نیروی کار، موضوع مثال ۲-۷، برای σ^2 فاصله اطمینان ۹۵ درصد بسازید. نامساوی ۲-۴۴ را می نویسیم،

$$\frac{(n-2)\hat{\sigma}^2}{\chi^2_{.975}} < \sigma^2 < \frac{(n-2)\hat{\sigma}^2}{\chi^2_{.025}}$$

با استفاده از مثال ۲-۷ می دانیم $n = 10$ و $\hat{\sigma}^2 = 1/83$ است، در اینجا باید $\chi^2_{.025}$ و $\chi^2_{.975}$ را با استفاده از جدول χ^2 تراکمی و با $10 - 2 = 8$ درجه آزادی به دست آوریم. خواهیم داشت

$$\chi^2_{.025}(8) = 2/18 \quad , \quad \chi^2_{.975}(8) = 17/5 \quad ,$$

بنابراین:

$$\Pr\left(\frac{8\hat{\sigma}^2}{17/5} < \sigma^2 < \frac{8\hat{\sigma}^2}{2/18}\right) = 0/95 \quad .$$

به عبارت دیگر در سطح احتمال ۹۵ درصد، فاصله اطمینان برای σ^2 عبارت است از

$$\frac{8(1/83)}{17/5} < \sigma^2 < \frac{8(1/83)}{2/18} \quad ,$$

یا

$$0/83 < \sigma^2 < 6/71 \quad .$$

۲-۴ آنالیز واریانس و آزمون معنی دار بودن مدل رگرسیون

می دانیم آنالیز واریانس یکی از مفاهیم بسیار مهم آمار است. در این قسمت به کاربرد آن در رگرسیون خطی ساده اشاره می کنیم. در فصل اول، قسمت ۱-۴ وقتی که بحث ما

در باره ضریب تعیین بود، قضیه بسیار مهم ۱-۴۰ را ثابت کردیم، که

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2 ,$$

یا

$$TSS = ESS + RSS .$$

این قضیه نشان می دهد که تغییرات متغیر درون زا (TSS) را می توان به دو قسمت تقسیم کرد: تغییرات توضیح داده شده (ESS)، و تغییرات توضیح داده نشده (RSS).

یکی از هدفهای آنالیز واریانس این است که بتواند معنی دار بودن تغییرات توضیح داده شده را آزمون کند. البته در مورد مدل رگرسیون ساده، مفهوم این آزمون چیزی نیست جز اینکه $\beta = 0$ را آزمون کنیم.

معمولاً تقسیم بندی تغییرات متغیر درون زا به تغییرات توضیح داده شده و توضیح داده نشده را همراه با درجات آزادی و میانگین هر کدام در یک جدول به شرح ذیل تنظیم می کنند و آن را جدول تجزیه واریانس می نامند.

جدول ۲-۱ آنالیز واریانس برای مدل رگرسیون ساده

منبع تغییرات	مجموع مربعات	درجات آزادی	میانگین مربعات یا واریانس
تخمین رگرسیون	$ESS = \sum \hat{y}_i^2 = \hat{\beta} \sum x_i y_i$	۱	$ESS / 1$
پسماند	$RSS = \sum e_i^2 = \sum y_i^2 - \hat{\beta} \sum x_i y_i$	$(n-2)$	$RSS / (n-2)$
کل تغییرات	$TSS = \sum y_i^2$	$(n-1)$	$F = \frac{ESS / 1}{RSS / (n-2)}$

برای تبیین بیشتر این جدول، به نکات زیر اشاره می کنیم. اثبات رابطه

$$ESS = \hat{\beta} \sum x_i y_i ,$$

بسیار ساده است. با استفاده از معادله ۱-۳۱، داریم

$$ESS = \sum \hat{y}_i^2 = \sum \hat{\beta}^2 x_i^2 ,$$

$$= \hat{\beta}' \sum x_i^2 = \hat{\beta} \frac{\sum x_i y_i}{\sum x_i^2} \cdot \sum x_i^2 ,$$

یا

$$ESS = \hat{\beta} \sum x_i y_i . \quad (۲-۴۵)$$

برای اثبات

$$RSS = \sum y_i^2 - \hat{\beta} \sum x_i y_i ,$$

می‌گوییم بنا بر تعریف

$$\begin{aligned} RSS &= \sum (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2 , \\ &= \sum [(Y_i - \bar{Y}) - \hat{\beta} (X_i - \bar{X})]^2 , \\ &= \sum y_i^2 + \hat{\beta}^2 \sum x_i^2 - 2 \hat{\beta} \sum x_i y_i , \\ &= \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} , \end{aligned}$$

در نتیجه

$$\sum e_i^2 = RSS = \sum y_i^2 - \hat{\beta} \sum x_i y_i , \quad (۲-۴۶)$$

و با جایگزینی $\hat{\beta}$ خواهیم داشت^۱

$$\sum e_i^2 = RSS = \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} . \quad (۲-۴۷)$$

۱. فرمول ۲-۴۷ را می‌توان به صورت زیر نیز نوشت. از $\sum y_i$ فاکتور می‌گیریم.

$$\sum e_i^2 = \sum y_i \left[1 - \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2} \right] ,$$

در نتیجه

$$\sum e_i^2 = (1 - r^2) \sum y_i^2 .$$

از این فرمول در فصل چهارم و در قسمت ۴-۳ استفاده خواهیم کرد.

ملاحظه می شود که اگر تغییرات توضیح داده شده (ESS) را با تغییرات توضیح داده شده (RSS) جمع کنیم حاصل، $\sum \hat{y}_i$ کل تغییرات خواهد بود. بحث فوق در واقع روش دیگری برای اثبات قضیه بسیار مهم ۱-۴۰ است.

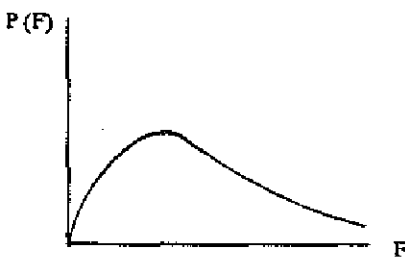
بحث بیشتر در مورد جدول ۲-۱ مستلزم بررسی توزیع احتمال F است؛ بنابراین در همینجا به مفاهیم اساسی این توزیع اشاره کرده، سپس به آزمون معنی دار بودن مدل رگرسیون بر می گردیم.

۱. توزیع و آزمون F

توزیع F برحسب دو متغیر مستقل از هم که هر یک توزیع χ^2 دارند تعریف می شود. متغیرهای تصادفی U و V با توزیهای χ^2 و مستقل از یکدیگر مفروض است. همچنین فرض بر این است که درجات آزادی U و V به ترتیب ν_1 و ν_2 است. در این صورت آماره زیر، توزیع F با درجات آزادی ν_1 و ν_2 خواهد داشت،

$$F = \frac{\frac{U}{\nu_1}}{\frac{V}{\nu_2}}$$

جدول F برحسب درجات آزادی ν_1 و ν_2 مدون شده است. ν_1 و ν_2 به ترتیب



نمودار ۲-۱۰ تابع توزیع احتمال F

درجات آزادی صورت و مخرج کسر F است. توزیع F مانند توزیع χ^2 از مرکز مختصات شروع می شود و به طرف راست تا بی نهایت ادامه دارد. نمودار ۲-۱۰ منحنی تغییرات F را نشان می دهد. توزیع F معمولاً برای مواردی به کار می رود که می خواهیم چند فرضیه را به طور همزمان در مورد

پارامترها آزمون کنیم؛ بنابراین آزمون F بیشتر برای آن دسته از مدل‌های رگرسیون مفید است که بیش از یک متغیر توضیحی دارند (این نکته را در فصل ششم و در قسمت ۶-۴، در مباحث آنالیز واریانس در رگرسیون چند متغیره توضیح خواهیم داد). همچنین - همان گونه که در آمار دیده‌ایم - می‌توان از توزیع F برای آزمون تساوی دو واریانس نیز استفاده کرد.

در اینجا پس از یادآوری تعریف توزیع F به بحث آنالیز واریانس در مدل‌های رگرسیون ساده می‌پردازیم. با توجه به معادله ۲-۳۰ داریم

$$\frac{\hat{\beta} - \beta}{\sqrt{\sigma^2 / \sum x_i^2}} \sim N(0, 1) .$$

می‌دانیم اگر توزیع نرمال استاندارد را مجذور کنیم، توزیع χ^2 خواهیم داشت؛ بنابراین با توجه به رابطه ۲-۳۷ اگر کسر فوق را مجذور کنیم، توزیع χ^2 با یک درجه آزادی به دست می‌آید،

$$\frac{(\hat{\beta} - \beta)^2}{\sigma^2 / \sum x_i^2} \sim \chi^2(1) . \quad (2-48)$$

از معادله ۲-۳۸ نیز می‌دانیم که

$$\frac{\sum e_i^2}{\sigma^2} \sim \chi^2(n-2) .$$

با توجه به تعریف توزیع F، اگر معادله ۲-۴۸ و ۲-۳۸ را به ترتیب بر درجات آزادی آنها تقسیم کرده و نتایج را نیز بر هم تقسیم کنیم و کسر حاصل را F بنامیم، آنگاه F دارای توزیع F با درجات آزادی ۱ و (n - ۲) خواهد بود؛ به گونه‌ای که درجه آزادی ۱ مربوط به صورت و درجه آزادی (n - ۲) متعلق به مخرج است؛

$$F = \frac{(\hat{\beta} - \beta)^2 \sum x_i^2 / 1}{\sum e_i^2 / (n - 2)} \sim F(1, n - 2) . \quad (2-49)$$

گفتیم که هدف آنالیز واریانس آزمون معنی دار بودن تغییرات توضیح داده شده است. به عبارت دیگر، اگر بخواهیم آنچه از تغییرات متغیر درون را - که توضیح داده شده است - آزمون کنیم، در واقع مانند این است که فرضیه $H_0: \beta = 0$ ، یعنی کل مدل را آزمون کرده باشیم. حال اگر فرض کنیم که فرضیه H_0 صحیح است؛ یعنی $\beta = 0$ ، آنگاه آماره F در معادله ۲-۴۹ را می توان چنین نوشت،

$$F = \frac{\hat{\beta}' \sum x_i^2 / 1}{\sum e_i^2 / (n - 2)} \sim F(1, n - 2) . \quad (2-50)$$

مقدار F در معادله ۲-۵۰ را آماره آزمون می نامیم و می گوئیم اگر فرضیه $\beta = 0$ صحیح باشد، آنگاه مقدار آماره آزمون در معادله ۲-۵۰ توزیع F خواهد داشت. بدیهی است برای هر سطح معنی داری که برای آزمون در نظر می گیریم، می توان F متناظر با نقطه بحرانی را از جدول F به دست آورد. اگر آماره آزمون در ناحیه بحرانی قرار بگیرد، فرضیه H_0 رد می شود، در غیر این صورت H_0 رد نمی شود؛ پس نتیجه می گیریم که مدل، معنی دار نیست؛ یعنی X_i نمی تواند تغییرات Y_i را توضیح دهد.

رابطه ۲-۵۰ را به صورت ساده تری می نویسیم. با استفاده از معادله ۱-۳۱ داریم

$$\hat{\beta}' \sum x_i^2 = \sum \hat{y}_i^2 = ESS ,$$

در نتیجه خواهیم داشت

$$F = \frac{ESS / 1}{RSS / (n - 2)} \sim F(1, n - 2) . \quad (2-51)$$

به فرمول ۲-۵۱ که بر اساس فرضیه $H_0: \beta = 0$ استوار است، آزمون مدل رگرسیون^۱ یا آزمون معنی دار بودن مدل رگرسیون^۲ می گویند. ملاحظه می شود که رابطه ۲-۵۱ چیزی نیست جز نسبت دو مقداری که در ستون آخر جدول تجزیه واریانس (جدول ۲-۱)

1. Testing the Regression Equation
2. Testing the Significance of Regression Equation

به دست آوردیم.

مثال ۲-۱۲ مدل رابطه تولید با ساعتهای نیروی کار، موضوع مثالهای ۱-۲ و ۲-۱ را در نظر بگیرید. برای این مدل، جدول آنالیز واریانس را تشکیل دهید و آزمون F را برای $\beta = 0$ در سطح معنی دار ۵ درصد انجام دهید. با استفاده از محاسبه‌های قبلی می‌دانیم

$$\sum x_i y_i = 21, \quad \hat{\beta} = 0.75,$$

$$\sum x_i^2 = 28, \quad \sum y_i^2 = 30/4.$$

و با استفاده از فرمول ۲-۴۵ داریم

$$ESS = 0.75(21) = 15/75,$$

و با استفاده از فرمول ۲-۴۶ خواهیم داشت

$$RSS = 30/4 - 0.75(21) = 14/65.$$

به این ترتیب جدول آنالیز واریانس برای این مثال به صورت جدول ۲-۴ به دست می‌آید.

جدول ۲-۴ مقادیر عددی جدول آنالیز واریانس برای اعداد جدول ۱-۴

منبع تغییرات	مجموع مربعات	درجات آزادی	میانگین مربعات یا واریانس
تغییرات X_j	۱۵/۷۵	۱	۱۵/۷۵
تغییرات e_i	۱۴/۶۵	۸	۱/۸۳
کل تغییرات	۳۰/۴۰	۹	$F = \frac{15/75}{1/83} = 8/606$

آماره آزمون از فرمول ۲-۵۱ به سهولت قابل محاسبه است.

$$F = \frac{15/75}{1/83} = 8/606.$$

مقدار به دست آمده از جدول F یا درجات آزادی صورت و مخرج به ترتیب برابر ۱ و ۸ و در سطح معنی دار ۵ درصد، برابر با ۵/۳۲ است. بدین ترتیب در سطح معنی دار ۵ درصد، فرضیه H_0 رد می شود.

۲. رابطه بین r^2 ، t ، F

در مثال فوق دیدیم که آماره آزمون F برابر ۸/۶۰۶ به دست آمد. قبلاً در مثال ۲-۸ ملاحظه شد که برای همین مورد، مقدار t برابر با ۲/۹۳ است؛ بنابراین مقدار F برابر مجذور مقدار t است. می توان این نکته را در حالت کلی ثابت کرد. از معادله ۲-۳۶ می دانیم که

$$t = \frac{\hat{\beta} - \beta}{\sqrt{s^2 / \sum x_i^2}}$$

و بنابر فرضیه $H_0: \beta = 0$ ، خواهیم داشت

$$t = \frac{\hat{\beta}}{\sqrt{s^2 / \sum x_i^2}}$$

مقدار s^2 را از معادله ۲-۲۸ در رابطه فوق جایگزین کرده، t^2 را محاسبه می کنیم،

$$t^2 = \frac{\hat{\beta}^2 \sum x_i^2}{\sum e_i^2 / (n - 2)} = \frac{\hat{\beta}^2 \sum x_i^2 / 1}{\sum e_i^2 / (n - 2)}$$

که دقیقاً همان فرمول ۲-۵۰ است. نتیجه می گیریم که برای آزمون $H_0: \beta = 0$ ، خواهیم داشت

$$F = t^2 \quad (2-51)$$

در اینجا برای یافتن رابطه بین F و r^2 ، معادله ۱-۲۹ را دوباره می نویسیم،

$$r^2 = \frac{ESS}{TSS}$$

در نتیجه

$$ESS = r^2 TSS \quad (2-52)$$

همچنین از معادله ۱-۴۳ داریم

$$RSS = (1 - r^2) TSS \quad (۲-۵۳)$$

با جایگزینی معادله‌های ۲-۵۲ و ۲-۵۳ در ۲-۵۱ خواهیم داشت

$$F = \frac{r^2/1}{(1 - r^2)/(n - 2)} \quad ,$$

یا

$$F = \frac{(n - 2) r^2}{1 - r^2} \quad (۲-۵۴)$$

برای یافتن رابطه بین r و t کافی است معادله ۲-۵۱ را در معادله ۲-۵۴ قرار دهیم،

$$t^2 = \frac{(n - 2) r^2}{1 - r^2} \quad , \quad (۲-۵۵)$$

به همین ترتیب، خواهیم داشت

$$r^2 = \frac{t^2}{t^2 + (n - 2)} \quad (۲-۵۶)$$

فرمول ۲-۵۶، رابطه بین مقدار آماره t برای فرضیه $\beta = 0$ و ضریب تعیین r^2 را مشخص می‌کند.

در پایان مباحث آزمونها، به این نکته اشاره می‌کنیم که سه آزمون فوق در حقیقت صورتهای مختلف یک واقعیت است؛ به عبارت دیگر، اولاً، آزمونی که شامل r^2 (فرمول ۲-۵۴)، آزمون معنی دار بودن ضریب تعیین است؛ ثانیاً، آزمونی که بر اساس $\hat{\beta}$ (فرمول ۲-۳۶) استوار می‌شود، در واقع آزمون معنی دار بودن شیب مدل رگرسیون خواهد بود؛ ثالثاً، از طریق آنالیز واریانس کوشش می‌شود که معنی دار بودن تغییرات توضیح داده شده آزمون شود.

مثال ۲-۱۳ با استفاده از t یا F در مثال ۲-۱۲ مقدار r^2 را به دست آورید.

می‌دانیم t^2 یا F در این مثال برابر ۸/۶ است. با استفاده از فرمول ۲-۵۶، r^2

به سهولت به دست می‌آید،

$$r^2 = \frac{8/6}{8/6 + 8} = \frac{8/6}{16/6} = 0.52 .$$

۲-۵ خصوصیات مطلوب تخمین زنده‌ها

در این قسمت، ابتدا مفهوم خصوصیات مطلوب یک تخمین زنده را در حالت کلی بررسی کرده، سپس انطباق این خصوصیات مطلوب را با تخمین زنده‌های حداقل مربعات معمولی (OLS) در یک مدل رگرسیون ساده مطالعه خواهیم کرد. بنابراین ضرورتاً از مباحث اقتصادسنجی خارج شده، مطالب صرفاً آماری را ارائه می‌دهیم.

متغیر تصادفی X_i مفروض است. فرض کنید θ یکی از پارامترهایی است که تابع توزیع احتمال این متغیر تصادفی را مشخص می‌کند. هدف ما تخمین θ است جامعه مشاهدات X_i ، شامل تمام مقادیر ممکن این متغیر بوده و θ نیز یکی از مشخصه‌های پارامتریکی^۱ این جامعه است. X_i می‌تواند یک متغیر پیوسته یا ناپیوسته فرض شود؛ برای مثال، می‌توان X_i را درآمد خانوار فرض کرد؛ بنابراین میانگین درآمد، یعنی θ ، یکی از مشخصه‌های پارامتریکی این جامعه محسوب خواهد شد. می‌خواهیم میانگین این جامعه را تخمین بزنیم. بدیهی است این مثال ساده‌ترین حالت در مسأله تخمین است. در موارد پیشرفته، هدف این است که بتوان چند پارامتر در مورد چند متغیر را همزمان تخمین زد؛ با وجود این، برای بررسی مشخصه‌های مطلوب تخمین زنده‌ها، بهتر است از ساده‌ترین حالت شروع کنیم؛ زیرا نتایج حاصل برای مدل‌های پیچیده‌تر نیز دقیقاً صادق خواهد بود.

برای تخمین یک پارامتر، معمولاً دو سری اطلاعات در دسترس است؛ اطلاعات اولیه^۱ و اطلاعات حاصل از نمونه. معلوماتی که در مورد جامعه مشاهدات X_i داریم، اطلاعات اولیه نامیده می‌شود؛ مانند فرضهای صورت تابع توزیع احتمال، یا مقادیر بعضی دیگر از پارامترهای جامعه غیر از θ یا حتی فرضهایی که در مورد قلمرو تغییرات

θ ارائه شده است. اطلاعات حاصل از نمونه، شامل کلیه اطلاعاتی است که از نمونه n تایی X_1, X_2, \dots, X_n به دست می آید. سؤال این است که چگونه می توان از این دوسری اطلاعات برای تخمین پارامتر θ استفاده کرد. برای پاسخ به این سؤال باید به فرمولی مراجعه کرد که برای تخمین پارامتر به دست می آوریم، به این فرمول «تخمین زننده» می گویند. معمولاً همواره می توان به بیش از یک فرمول برای تخمین θ رسید؛ به عبارت دیگر، برای تخمین θ ، همواره تخمین زننده های متعددی وجود دارد. سؤال این است که چگونه می توان از بین تخمین زننده های مختلف، مناسبترین را انتخاب کرد. بدیهی است هر تخمین زننده، خصوصیات خاص خود را دارد. بنابراین برای اینکه بتوان مناسبترین تخمین زننده را برگزید؛ ابتدا باید خصوصیات مطلوب تخمین زننده ها را بررسی کرد. بهترین تخمین زننده آن است که بیشترین خصوصیات مطلوب را دارا باشد.

در قسمت ۱-۲ دیدیم که به یک تخمین زننده از θ اصطلاحاً $\hat{\theta}$ می گوئیم. می دانیم $\hat{\theta}$ نیز یک متغیر تصادفی است، بنابراین دارای میانگین و واریانس است. اگر میانگین $\hat{\theta}$ را با $E(\hat{\theta})$ نشان دهیم، آنگاه بنا بر تعریف واریانس $\hat{\theta}$ برابر است با

$$\text{Var}(\hat{\theta}) = E[\hat{\theta} - E(\hat{\theta})]^2 .$$

می دانیم انحراف معیار $\hat{\theta}$ برابر است با جذر واریانس $\hat{\theta}$ ،

$$\text{SE}(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})} .$$

حال مفاهیم زیر را تعریف می کنیم.

$$\theta - \hat{\theta} = \text{خطای نمونه گیری} ,$$

$$\text{اریب} = E(\hat{\theta}) - \theta ,$$

$$\text{میانگین مربع خطا} = E(\hat{\theta} - \theta)^2 .$$

«خطای نمونه گیری»^۱، در واقع تفاوت مقدار تخمین و مقدار واقعی پارامتری است که باید تخمین زده شود. «اریب»^۲ به تفاوت بین میانگین توزیع تخمین زنده $E(\hat{\theta})$ و مقدار واقعی پارامتر θ گفته می‌شود. مقدار اریب برای هر تخمین زنده معمولاً مقدار معینی است که می‌تواند صفر یا غیر صفر باشد. سرانجام میانگین مربع خطا (MSE)^۳ مفهومی است که با پراکندگی توزیع یک تخمین زنده مرتبط بوده بنابراین بسیار نزدیک به مفهوم واریانس است. تفاوت بین واریانس و میانگین مربع خطای یک تخمین زنده این است که واریانس در واقع پراکندگی توزیع تخمین زنده را در حول مقدار میانگین آن اندازه گیری می‌کند؛ در حالی که میانگین مربع خطا، پراکندگی توزیع تخمین زنده را در اطراف مقدار واقعی پارامتر مشخص می‌سازد. اگر میانگین توزیع $\hat{\theta}$ ، یعنی $E(\hat{\theta})$ بر مقدار واقعی پارامتر یعنی θ منطبق شود، $Var(\hat{\theta})$ و میانگین مربع خطای $\hat{\theta}$ ، یعنی $MSE(\hat{\theta})$ ، با هم برابر می‌شود، در غیر این صورت قطعاً با یکدیگر متفاوتند. رابطه بین میانگین مربع خطای $\hat{\theta}$ و واریانس $\hat{\theta}$ را می‌توان به صورت زیر به دست آورد. می‌دانیم

$$MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2$$

به طرف راست و در داخل پرانتز $E(\hat{\theta}) \pm$ را اضافه می‌کنیم

$$\begin{aligned} MSE(\hat{\theta}) &= E[\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta]^2, \\ &= E\{[\hat{\theta} - E(\hat{\theta})] + [E(\hat{\theta}) - \theta]\}^2, \\ &= E[\hat{\theta} - E(\hat{\theta})]^2 + [E(\hat{\theta}) - \theta]^2 + \\ &\quad 2E\{[\hat{\theta} - E(\hat{\theta})][E(\hat{\theta}) - \theta]\}. \end{aligned}$$

جمله آخر برابر صفر است؛ زیرا بنا بر قواعد امید ریاضی، می‌دانیم $E(ax) = aE(x)$ و

$$E(a) = a, \text{ بنابراین}$$

$$\begin{aligned} E \{ [\hat{\theta} - E(\hat{\theta})][E(\hat{\theta}) - \theta] \} &= [E(\hat{\theta}) - \theta] E[\hat{\theta} - E(\hat{\theta})], \\ &= [E(\hat{\theta}) - \theta][E(\hat{\theta}) - E(\hat{\theta})], \\ &= 0. \end{aligned}$$

بدین ترتیب $MSE(\hat{\theta})$ عبارت است از

$$MSE(\hat{\theta}) = E[\hat{\theta} - E(\hat{\theta})]^2 + [E(\hat{\theta}) - \theta]^2,$$

یا

$$MSE(\hat{\theta}) = \text{Var}(\hat{\theta}) + (\text{اریب})^2. \quad (2.57)$$

یعنی میانگین مربع خطای $\hat{\theta}$ برابر است با واریانس $\hat{\theta}$ به علاوه مربع مقدار اریب. در نتیجه مقدار میانگین مربع خطای $\hat{\theta}$ ، هیچگاه نمی تواند از مقدار واریانس $\hat{\theta}$ کمتر باشد:

$$MSE(\hat{\theta}) \geq \text{Var}(\hat{\theta}). \quad (2.58)$$

تفاوت $MSE(\hat{\theta})$ و $\text{Var}(\hat{\theta})$ دقیقاً برابر مربع مقدار اریب است.

بعد از ذکر این مقدمه، به بررسی خصوصیات از $\hat{\theta}$ می پردازیم که معمولاً به عنوان خصوصیات مطلوب تخمین زننده ها شناخته می شود. بحث را در دو قسمت خصوصیات مطلوب در نمونه های محدود یا نمونه های کوچک^۱ و خصوصیات مطلوب حدی در نمونه های بزرگ^۲ ارائه می دهیم. در حالت اول خصوصیات از $\hat{\theta}$ را بررسی می کنیم که حاصل نمونه هایی با حجم محدود و معین است و در حالت دوم خصوصیات $\hat{\theta}$ متعلق به نمونه هایی را به دست می آوریم که حجم آنها به سمت بی نهایت میل می کند.

1. Finite or Small Sample Properties
2. Asymptotic or Large Sample Properties

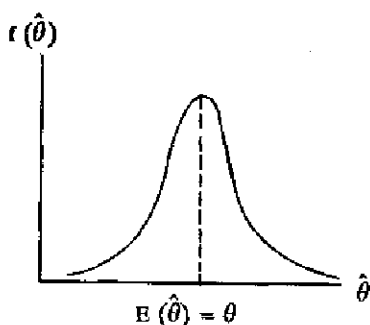
۱. خصوصیات مطلوب تخمین زنده‌ها در نمونه‌های کوچک
 در این قسمت به بررسی خصوصیات مطلوب $\hat{\theta}$ ، هنگامی که حاصل نمونه‌ای با حجم
 معین و محدود است، می‌پردازیم:

ناریبی

اولین و شاید مهمترین خصوصیات $\hat{\theta}$ ناریبی آن است. بنا بر تعریف $\hat{\theta}$ یک تخمین زنده
 ناریب از θ است.

هرگاه

$$E(\hat{\theta}) = \theta \quad (2-59)$$



نمودار ۲-۱۱ ناریبی $\hat{\theta}$

نمودار ۲-۱۱ خصوصیت ناریبی $\hat{\theta}$ را نشان
 می‌دهد. ملاحظه می‌شود که چون توزیع $\hat{\theta}$
 قرینگی دارد؛ بنابراین میانگین این توزیع
 یعنی $E(\hat{\theta})$ دقیقاً در محور مرکزی توزیع،

قرار می‌گیرد. اگر مقدار واقعی پارامتر، یعنی θ ، نیز دقیقاً در همین نقطه واقع شود،
 می‌گوییم $\hat{\theta}$ یک تخمین زنده ناریب از θ است.

مثالی که در مورد خصوصیت ناریبی یک تخمین زنده می‌توان ذکر کرد، این
 است که بگوییم میانگین نمونه، یعنی \bar{X} ، یک تخمین ناریب از میانگین جامعه، یعنی
 μ ، است؛ زیرا

$$E(\bar{X}) = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum E(X_i) = \mu$$

باید توجه کرد که ناریب بودن - به منزله تنها معیار مطلوب - نمی‌تواند چندان
 رضایت‌بخش باشد، زیرا هیچ نکته‌ای درباره پراکندگی توزیع تخمین زنده بیان
 نمی‌کند. یک تخمین زنده ناریب که واریانس بسیار بزرگی داشته باشد می‌تواند در
 بسیاری موارد به نتایجی برسد که ممکن است قابل قبول نباشد؛ زیرا می‌دانیم واریانس
 بزرگ به معنای دقت کم در تخمین پارامترهاست؛ از طرف دیگر، تخمین زنده‌ای که

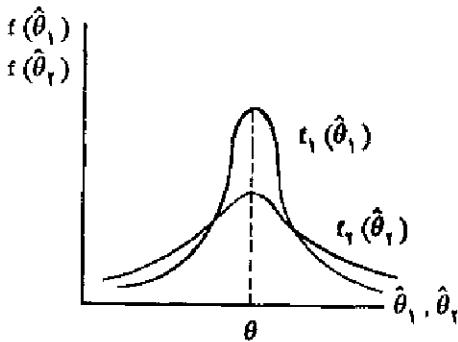
واریانس بسیار کم دارد، اما مقدار اریب آن مشخص نیست با همین استدلال، نمی تواند چندان مفید باشد.

از بحث فوق می توان نتیجه گرفت که تخمین زنده ای مطلوب است که میانگین مربع خطای $\hat{\theta}$ را حداقل کند؛ زیرا میانگین مربع خطا ($\hat{\theta}$) عبارت است از مجموع واریانس و مربع اریب. اما متأسفانه هنوز نمی توان تخمین زنده ای یافت که بتواند حداقل بودن میانگین مربع خطا را تضمین کند؛ زیرا هر تابعی از میانگین مربع خطا شامل مقدار واقعی θ است، در حالی که ما می خواهیم خود θ را تخمین بزنیم. بدین ترتیب معیار حداقل نمودن میانگین مربع خطا یک بحث نظری بیش نیست و در عمل ممکن است هیچگاه تحقق نیابد.

کارایی^۱

مباحثی که تا به حال مطرح کردیم، زمینه مناسبی برای طرح خصوصیت کارایی تخمین زنده هاست که در این قسمت مطرح خواهد شد. متأسفانه در آمار هنوز تعریف واحدی برای کارایی یک تخمین زنده ارائه نشده است. بعضی از نویسندگان، اصطلاح کارایی را در مواردی به کار می برند که یک تخمین زنده بتواند میانگین مربع خطا را حداقل کند. همان گونه که دیدیم حداقل نمودن میانگین مربع خطا تنها یک مفهوم نظری است و در عمل نمی توان $\hat{\theta}$ را چنان تعیین کرد که میانگین مربع خطا حداقل شود؛ بنابراین منظور، یک کارایی نسبی است؛ یعنی تخمین زنده ای از کارایی نسبی بیشتری برخوردار است که در مقایسه با تخمین زنده های دیگر، واریانس و اریب کمتری داشته باشد. عده ای دیگر نیز اصطلاح کارایی را فقط برای خصوصیات مطلوب حدی تخمین زنده ها به کار می برند و در نمونه های محدود و کوچک از این مفهوم استفاده نمی کنند. گروهی نیز معتقدند، تخمین زنده ای کاراست که اولاً نااریب باشد، ثانیاً حداقل واریانس را داشته باشد. به نظر می رسد که تعریف آخر بیشتر معمول است

و ما نیز در این کتاب از این تعریف استفاده خواهیم کرد.

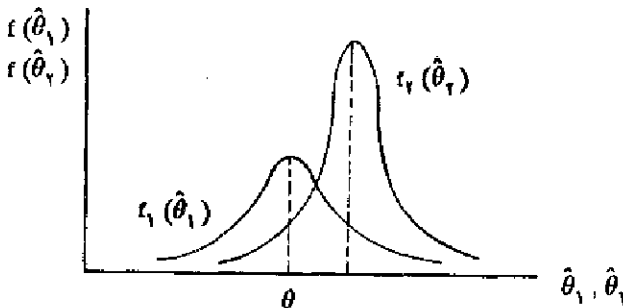


نمودار ۲.۱۲ ارب و واریانس

برای تبیین بیشتر مفهوم کارایی به دو نمودار ۲-۱۲ و ۲-۱۳ توجه می‌کنیم. در نمودار ۲-۱۲ دو تخمین‌زننده $\hat{\theta}_1, \hat{\theta}_2$ با توابع توزیع احتمال $f_1(\hat{\theta}_1)$ و $f_2(\hat{\theta}_2)$ مفروض هستند. هر دو می‌خواهند θ را تخمین بزنند. با اینکه هر دو تخمین‌زننده نااریب هستند، $\hat{\theta}_1$ پراکندگی بیشتری نسبت به $\hat{\theta}_2$ دارد. در

این حالت می‌توان گفت که $\hat{\theta}_2$ از کارایی کمتری نسبت به $\hat{\theta}_1$ برخوردار است. اگر نتوانیم هیچ تخمین‌زننده نااریب دیگری پیدا کنیم که نسبت به $\hat{\theta}_1$ واریانس کمتری داشته باشد، می‌گوییم $\hat{\theta}_1$ در میان تمام تخمین‌زننده‌های نااریب، حداقل واریانس را دارد؛ بنابراین یک تخمین‌زننده کارآست.

در نمودار ۲-۱۳ ملاحظه می‌کنیم که $\hat{\theta}_1$ نااریب است، اما واریانس بیشتری نسبت به $\hat{\theta}_2$ دارد - که خود دارای ارب است. نمی‌توان گفت کدام تخمین‌زننده بر دیگری برتری دارد، مگر اینکه وزنه‌ای - که بیان‌کننده درجه اهمیت است - به معیارهای نااریبی و واریانس کمتر منسوب کنیم. بدیهی است تعیین این ضرایب یا وزنه‌ها تابعی از



نمودار ۲.۱۳ ارب و واریانس

مشخصات مسأله‌ای است که در عمل با آن مواجهیم؛ به عبارت دیگر، در اقتصادسنجی،

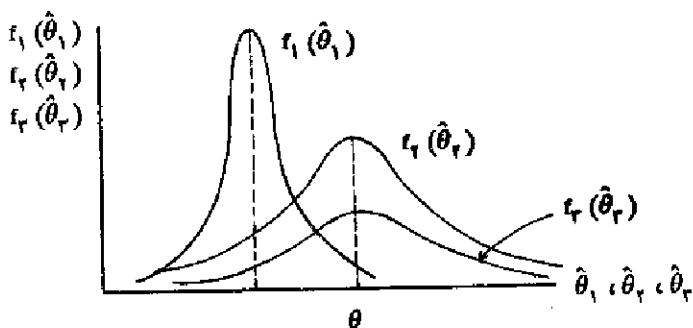
محقق با توجه به خصوصیات مسأله‌ای که در مقابل اوست، تصمیم می‌گیرد که آیا ناریبی نسبت به واریانس کمتر مرجح است یا برعکس. بعد از ذکر این مقدمه، می‌توان تعریف کارآیی را به صورت دقیقتری بیان کرد.

$\hat{\theta}$ یک تخمین‌زننده کارآ از θ است، اگر شرایط زیر برقرار باشد. اولاً، $\hat{\theta}$ ناریب باشد؛ $E(\hat{\theta}) = \theta$ ، ثانیاً، واریانس $\hat{\theta}$ از واریانس هر تخمین‌زننده ناریب دیگری از θ کمتر باشد،

$$\text{Var}(\hat{\theta}) \leq \text{Var}(\bar{\theta}), \quad (2.60)$$

که در آن $\bar{\theta}$ هر تخمین‌زننده ناریب دیگری از θ است.

در بعضی موارد به تخمین‌زننده‌های کارآ تخمین‌زننده ناریب حداقل واریانس (MVUE)^۱ نیز می‌گویند. اصطلاح دیگری که برای تخمین‌زننده‌های کارآ به کار می‌رود، بهترین تخمین‌زننده ناریب^۲ است. باید به این نکته توجه داشت که بنا بر تعریف فوق، یک تخمین‌زننده با حداقل واریانس - که کوچکترین مقدار اریب را داشته باشد - دیگریک تخمین‌زننده کارآ نخواهد بود. این نکته در نمودار ۲-۱۴ نشان داده شده است. در این نمودار سه تخمین‌زننده $\hat{\theta}_1$ ، $\hat{\theta}_2$ و $\hat{\theta}_3$ می‌خواهند θ را تخمین بزنند.



نمودار ۲-۱۴ اریب و واریانس

ملاحظه می‌شود که $\hat{\theta}_1$ حداقل واریانس را دارد، اما کارآ نیست، زیرا اریب دارد. از طرف

1. Minimum Variance Unbiased Estimator

2. Best Unbiased Estimator

دیگر، $\hat{\theta}_p$ و $\hat{\theta}_p$ هر دو ناریب هستند اما $\hat{\theta}_p$ واریانس بیشتری نسبت به $\hat{\theta}_p$ دارد، بنابراین $\hat{\theta}_p$ کارآ نیست. در مجموع می توان گفت هر گاه هیچ تخمین زنده ناریب دیگری نتوان یافت که واریانس کمتری از $\hat{\theta}_p$ داشته باشد، $\hat{\theta}_p$ یک تخمین زنده کارآست.

تعیین کارآیی در عمل بسیار مشکلتر از این است که مشخص کنیم آیا یک تخمین زنده اریب دارد یا خیر؛ زیرا برای یافتن اریب کافی است از $\hat{\theta}$ امید ریاضی بگیریم و اگر برابر θ شد، در آن صورت ناریبی $\hat{\theta}$ را نتیجه می گیریم. اما برای کارآیی باید معیار کمترین واریانس در بین همه تخمین زنده ها را دقیقاً بررسی کنیم. ممکن است دهها تخمین زنده وجود داشته باشد؛ بنابراین محاسبه واریانس همه آنها برای یافتن کمترین واریانس، قاعدتاً بسیار مشکل است. به همین دلیل مفید است که از مفهوم کارآیی نسبی^۱ نیز استفاده شود. در این صورت می گوئیم مثلاً $\hat{\theta}_p$ از کارآیی نسبی بیشتری نسبت به $\hat{\theta}_p$ در نمودار ۱۴-۲ برخوردار است^۲. راه دیگر این است که قلمرو تخمین زنده ها را محدود کنیم تا محاسبه واریانسها در آن قلمرو محدود، راحت تر انجام شود. برای این منظور ابتدا باید خطی بودن را خصوصیت مطلوب تخمین زنده ها معرفی کنیم.

تخمین زنده های خطی

یکی دیگر از خصوصیات مطلوب تخمین زنده ها این است که $\hat{\theta}$ یک تخمین زنده خطی باشد. خطی بودن در اینجا بدین معنی است که $\hat{\theta}$ یک تابع خطی از مشاهدات موجود در نمونه است؛ به عبارت دیگر، برای نمونه ای شامل X_1, X_2, \dots, X_n ، تخمین زنده $\hat{\theta}$ موقعی خطی است که داشته باشیم

$$\hat{\theta} = a_1 X_1 + a_2 X_2 + \dots + a_n X_n \quad (2-61)$$

1. Relative Efficiency

۲. البته برای تعیین درجه کارآیی یک تخمین زنده، وقتی تابع توزیع احتمال آن کاملاً مشخص باشد، می توان از نامساوی کرامر - راتو (Cramer - Rao Inequality) استفاده کرد که در کتابهای آمار به تفصیل بحث شده است.

با اضافه کردن خصوصیت خطی بودن تخمین‌زنده‌ها به خصوصیت کارآیی، می‌توان مفهوم دیگری را در این زمینه مطرح کرد که در ذیل به آن اشاره می‌کنیم.

خصوصیت بهترین تخمین‌زنده ناریب خطی

اصطلاح BLUE به معنای «بهترین تخمین‌زنده ناریب خطی»^۱ است. در اقتصادسنجی اصطلاح «بهترین» معمولاً مترادف «حداقل واریانس» به کار می‌رود؛ یعنی وقتی می‌گوییم $\hat{\theta}$ بهترین است، بدین معنی نیست که $\hat{\theta}$ از هر نظر بهترین تخمین‌زنده است، بلکه منظور این است که $\hat{\theta}$ حداقل واریانس را دارد.

$\hat{\theta}$ بهترین تخمین‌زنده خطی (BLUE) از θ است اگر سه شرط زیر را داشته باشد:

۱. $\hat{\theta}$ یک تابع خطی از مشاهدات موجود در نمونه، یعنی $\hat{\theta} = \sum a_i X_i$ باشد؛

۲. $\hat{\theta}$ ناریب، یعنی $E(\hat{\theta}) = \theta$ باشد؛

۳. واریانس $\hat{\theta}$ کمتر از واریانس هر تخمین‌زنده ناریب دیگر برای θ ، یعنی

$$\text{Var}(\hat{\theta}) \leq \text{Var}(\tilde{\theta})$$

که هر تخمین‌زنده ناریب دیگری از θ است.

برای مقایسه خصوصیت بهترین تخمین‌زنده خطی و کارآیی، مفید است سه

حالت زیر را ملاحظه کنیم.

حالت اول: $\hat{\theta}$ یک تخمین‌زنده کارآ و نیز یک تابع خطی از مشاهدات X_i است.

در این حالت بهترین تخمین‌زنده‌های خطی و کارآ در واقع یکی هستند. برای مثال، چون می‌توان نشان داد که \bar{X} یک تخمین‌زنده کارآ از میانگین یک توزیع احتمال نرمال و \bar{X} یک تابع خطی از X_i است؛ بنابراین \bar{X} بهترین تخمین‌زنده خطی از μ نیز خواهد بود.

حالت دوم: $\hat{\theta}$ یک تخمین‌زنده کارآ، اما تقریباً یک تابع خطی از مشاهدات X_i

است. در این حالت بهترین تخمین‌زنده‌های خطی دیگر تخمین‌زنده‌های کارآ نیستند. البته در این حالت می‌توان نشان داد که واریانس بهترین تخمین‌زنده خطی بسیار نزدیک به واریانس تخمین‌زنده کارآ است.

حالت سوم: $\hat{\theta}$ یک تخمین زنده کارآ، اما تابع غیرخطی از مشاهدات X است. در این حالت بهترین تخمین زنده خطی دیگر کارآ نیست، اما می توان نشان داد که واریانس آن به مراتب بیشتر از واریانس تخمین زنده کارآ خواهد بود.

آخرین نکته این است که خصوصیت خطی بودن در واقع یک خصوصیت بسیار مطلوب در تخمین پارامترهایی چون میانگین است. که در واقع گشتاورهای مرتبه اول است. برای تخمین پارامترهایی که صورت گشتاورهای مرتبه دوم را دارد؛ مانند واریانس، خطی بودن تخمین زنده، یک خصوصیت بسیار نامطلوب است. در تخمین این گونه پارامترها، خصوصیت مطلوب در واقع خصوصیت بهترین تخمین زنده ناریب درجه دو^۱ خواهد بود، که آن را با BQUE نشان می دهیم.

جامعیت

آخرین خصوصیت مطلوب تخمین زنده ها در نمونه های محدود، شرط «جامعیت»^۲ است. بنا بر تعریف، تخمین زنده ای جامع است که بتواند در تخمین پارامتر جامعه، از کلیه اطلاعات موجود در نمونه استفاده کند. یا توجه به اینکه هر مشاهده ای در مجموعه مشاهدات موجود در نمونه باید مشخص کننده ساختار جامعه باشد، در محاسبه $\hat{\theta}$ انتظار این است که تمام اطلاعات موجود در نمونه مشارکت داشته باشند. برای مثال، میانه نمی تواند تخمینی جامع از میانگین جامعه را باشد؛ زیرا برای محاسبه میانه به طور کلی از اولویت و طبقه بندی مشاهدات موجود در نمونه استفاده می شود و مقادیر تمام مشاهدات نقش چندانی ایفا نمی کند.

البته ممکن است استدلال شود که معیار استفاده از تمام اطلاعات موجود در نمونه نباید امتیاز برجسته ای در خصوصیات مطلوب یک تخمین زنده به حساب آید. در اقتصادسنجی مهم این است که از پارامترهای موجود در جامعه تخمینهای خوب و رضایت بخشی به دست آوریم، بنابراین مهم نیست که آیا از تمام اطلاعات موجود در

نمونه استفاده کرده‌ایم یا خیر. با وجود این، به نظر می‌رسد، این استدلال چندان معتبر نباشد؛ زیرا می‌توان ثابت کرد که معیار جامعیت شرط لازم برای کارآیی است.^۱ به عبارت دیگر، یک تخمین‌زننده کارآ نیست، مگر اینکه از تمام اطلاعات موجود در نمونه استفاده کرده باشد. همچنین ثابت شده است که مجموع خصوصیات نااریبی و جامعیت، دقیقاً موجب کارآیی می‌شود - که در آمار به نام قضیهٔ بلک ول - راثو^۲ معروف است.

در نتیجه سه شرط نااریبی، کارآیی و بهترین تخمین‌زننده نااریب خطی، در واقع، مهمترین خصوصیات مطلوب تخمین‌زننده‌ها هستند. این خصوصیات مطلوب - همان گونه که دیدیم - برحسب میانگین و واریانس تخمین‌زننده‌ها تعریف می‌شوند، بنابراین برای استنتاج در مورد هر یک از آنها باید میانگین و واریانس توزیع احتمال تخمین‌زننده‌ها را دقیقاً شناخت.

۲. خصوصیات حدی تخمین‌زننده‌ها

قبلاً گفتیم که خصوصیات حدی^۳ به نمونه‌هایی مربوط می‌شود که حجم آنها بسیار بزرگ است و به بی‌نهایت میل می‌کند. ابتدا مفهوم توزیع حدی^۴ را تعریف می‌کنیم. در بسیاری موارد، شکل توزیع یک متغیر تصادفی، تابعی از حجم نمونه است؛ برای مثال، بنا بر قضیهٔ حد مرکزی، می‌دانیم که به ازای افزایش حجم نمونه، توزیع میانگین نمونه، یعنی \bar{X} ، به سمت یک توزیع نرمال میل می‌کند. بنابراین، می‌گوییم توزیع نرمال، یک توزیع حدی برای میانگین نمونه است. به طور کلی، می‌توان گفت هرگاه توزیع یک تخمین‌زننده به ازای افزایش حجم نمونه به سمت خاصی میل کند، آن شکل خاص، توزیع حدی تخمین‌زننده نامیده می‌شود.

واژهٔ حدی در اینجا مطلقاً بدین معنی نیست که توزیع حدی برای یک

۱. برای توضیحات بیشتر در این زمینه به کتاب (Lindgrén, 1976) صفحات ۲۶۴ به بعد مراجعه شود.

2. Blackwell - Rao Theorem

3. Asymptotic Properties

4. Asymptotic Distribution

تخمین زنده، در واقع صورت نهایی توزیع $\hat{\theta}$ است. وقتی n به سمت بی نهایت میل کند، آنچه در عمل اتفاق می افتد، این است که به ازای افزایش حجم نمونه، توزیع احتمال تخمین زنده در یک نقطه معینی به مرحله فروپاشی^۱ می رسد، و چه بسا ممکن است این نقطه دقیقاً منطبق با مقدار واقعی پارامتر شود. برای مثال، دیدیم که میانگین توزیع \bar{X} برابر μ و واریانس آن $\frac{\sigma^2}{n}$ است که σ^2 واریانس جامعه است. بدیهی است وقتی n به سمت بی نهایت میل کند، $\frac{\sigma^2}{n}$ به سمت صفر میل می کند و توزیع \bar{X} در نقطه میانگین جامعه (μ) فرو می پاشد. از نظر هندسی، توزیع در این نقطه با یک خط عمودی مشخص می شود که قاعدتاً منحنی نرمال نخواهد بود؛ بنابراین ملاحظه می شود که منظور از توزیع حدی یک تخمین زنده، شکل نهایی توزیع آن نیست، بلکه در واقع صورتی از توزیع است که تخمین زنده قبل از مرحله فروپاشی به آن رسیده است. در مثال توزیع میانگین نمونه \bar{X} ، دیدیم که به ازای افزایش حجم نمونه، واریانس این توزیع به طور مرتب کمتر می شود و در عین حال شکل توزیع نیز بیشتر به نرمال میل می کند. قبل از مرحله فروپاشی، توزیع به صورت نرمال و با واریانس بسیار کوچکی خواهد بود.

حال که با مفهوم توزیع حدی آشنا شدیم، این سؤال مطرح می شود که چگونه می توان توزیع حدی تخمین زنده ها را به دست آورد. برای بسیاری از تخمین زنده ها پاسخ به این سؤال بسیار ساده است؛ زیرا شکل توزیع آنها تابعی از حجم نمونه نیست. توزیع میانگین نمونه برای یک جامعه نرمال، مثال خوبی در این مورد است؛ زیرا این توزیع، به رغم حجم نمونه، همواره نرمال و میانگین آن μ و واریانس آن $\frac{\sigma^2}{n}$ است. در بسیاری موارد، شکل توزیع تخمین زنده در نمونه های کوچک چندان روشن نیست؛ زیرا معمولاً شکل توزیع به موازات افزایش حجم نمونه و میل آن به سمت بی نهایت مشخص می شود. برای مثال، می توان گفت توزیع میانگین نمونه، به عنوان تخمینی از میانگین جامعه، دارای چنین خصوصیتی برای جامعه های غیر نرمال است؛ زیرا در حقیقت بر اساس قضیه حد مرکزی می دانیم که به ازای افزایش حجم نمونه و میل آن

1. Degenerate

به سمت بی‌نهایت، \bar{X} توزیع نرمال خواهد داشت.

توزیع حدی نیز با میانگین و واریانس آن مشخص می‌شود. میانگین این توزیعها را در اصطلاح «میانگین حدی»^۱ و واریانس آن را «واریانس حدی»^۲ می‌گویند. برای یافتن میانگین حدی، کافی است از امید ریاضی آن حد بگیریم. برای مثال، می‌دانیم میانگین $\hat{\theta}$ برابر با $E(\hat{\theta})$ است، اما میانگین حدی آن برابر با

$$\hat{\theta} \text{ میانگین حدی} = \lim_{n \rightarrow \infty} E(\hat{\theta}) \quad (۲-۶۲)$$

خواهد بود. به همین ترتیب واریانس حدی $\hat{\theta}$ را تعریف می‌کنیم،

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}) = \lim_{n \rightarrow \infty} E[\hat{\theta} - E(\hat{\theta})]^2 \quad (۲-۶۳)$$

در ذیل به سه خصوصیت حدی برای تخمین‌زنده‌ها اشاره می‌شود.

نااریب حدی

یک تخمین‌زنده، موقعی نااریب حدی^۳ است که به ازای افزایش حجم نمونه، به طور مرتب از مقدار اریب آن کم شود و وقتی حجم نمونه به سمت بی‌نهایت میل می‌کند، مقدار اریب صفر شود،

$$\lim_{n \rightarrow \infty} E(\hat{\theta}) = \theta \quad (۲-۶۴)$$

به عبارت دیگر، یک تخمین‌زنده، موقعی نااریب حدی است که در حد نااریب شود. باید در نظر داشت که اگر یک تخمین‌زنده نااریب باشد آنگاه نااریب حدی نیز هست، اما نه برعکس. دلیل این امر روشن است: وقتی تخمین‌زنده‌ای نااریب باشد، امید ریاضی آن به ازای هر حجمی از نمونه برابر مقدار واقعی پارامتر است؛ به بیان دیگر، اگر برای تخمین پارامتر θ ، نمونه‌ای با حجم n داشته باشیم و به کمک این نمونه مقدار $\hat{\theta}$ را به دست آوریم، آنگاه با اینکه ممکن است برای مقادیر کوچک n مقدار $E(\hat{\theta})$ برابر θ

1. Asymptotic Mean

2. Asymptotic Variance

3. Asymptotically Unbiased

نشود، اما اگر برای مقادیر بزرگ و بزرگتر n ، مقدار $E(\hat{\theta})$ به سمت θ میل کند، به گونه‌ای که در حد مقدار $E(\hat{\theta})$ دقیقاً برابر θ شود، در آن صورت می‌گوییم $\hat{\theta}$ یک تخمین نااریب حدی از θ است.

مثال ۲-۱۴ نشان دهید که واریانس نمونه s^2 ، یک تخمین نااریب حدی از واریانس جامعه σ^2 است. اگر واریانس نمونه را به صورت زیر تعریف کنیم،

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{n} , \quad (2-65)$$

با اضافه کردن $\pm \mu$ به داخل پرانتز صورت کسر خواهیم داشت

$$\begin{aligned} s^2 &= \frac{1}{n} \left\{ \sum [(X_i - \mu) - (\bar{X} - \mu)]^2 \right\} , \\ &= \frac{1}{n} [\sum (X_i - \mu)^2 + n(\bar{X} - \mu)^2 - 2(\bar{X} - \mu) \sum (X_i - \mu)] . \end{aligned}$$

اما می‌دانیم $\sum (X_i - \mu) = n(\bar{X} - \mu)$ بنابراین

$$\begin{aligned} s^2 &= \frac{1}{n} [\sum (X_i - \mu)^2 - n(\bar{X} - \mu)^2] , \\ &= \frac{1}{n} \sum (X_i - \mu)^2 - (\bar{X} - \mu)^2 . \end{aligned}$$

حال اگر امید ریاضی بگیریم، با توجه به $E(\bar{X}) = \mu$ خواهیم داشت

$$\begin{aligned} E(s^2) &= \frac{1}{n} \sum E(X_i - \mu)^2 - E(\bar{X} - \mu)^2 , \\ &= \frac{1}{n} \sum E[X_i - E(X_i)]^2 - E[\bar{X} - E(\bar{X})]^2 , \\ &= \frac{1}{n} \sum \text{Var}(X_i) - \text{Var}(\bar{X}) . \end{aligned}$$

با توجه به اینکه $\text{Var}(X_i) = \sigma^2$ و $\text{Var}(\bar{X}) = \frac{1}{n} \sigma^2$ و $\sum \text{Var}(X_i) = n\sigma^2$ ، در نتیجه

$$E(s^2) = \sigma^2 - \frac{1}{n} \sigma^2 ,$$

یا

$$E(s^2) = \left(1 - \frac{1}{n}\right) \sigma^2 . \quad (2-66)$$

ملاحظه می شود که $E(s^2) \neq \sigma^2$ و s^2 تخمین نااریبی از σ^2 است. اما نکته مهم این است که اگر n به سمت بی نهایت میل کند؛ یعنی حجم نمونه به طور مرتب بزرگ و بزرگتر شود، آنگاه در حد $E(s^2)$ برابر σ^2 خواهد شد؛ زیرا با افزایش n مقدار $\frac{1}{n}$ به سمت صفر میل می کند،

$$\lim_{n \rightarrow \infty} E(s^2) = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right) \sigma^2 = \sigma^2 , \quad (2-67)$$

یعنی s^2 یک تخمین زنده نااریب حدی از σ^2 است. به همین دلیل است که s^2 را در مباحث آماری به صورت زیر تعریف می کنند که تخمین نااریبی از σ^2 باشد،

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{n-1} . \quad (2-68)$$

زیرا از معادله ۲-۶۶ داریم

$$E(s^2) = \frac{n-1}{n} \sigma^2 ,$$

یا

$$\frac{n}{n-1} E(s^2) = \sigma^2 ,$$

$$E\left(\frac{ns^2}{n-1}\right) = \sigma^2 . \quad (2-69)$$

با توجه به معادله ۲-۶۵ می دانیم

$$ns^2 = \sum (X_i - \bar{X})^2 ,$$

که با جایگزینی در ۲-۶۹ خواهیم داشت

$$E\left[\frac{\sum (X_i - \bar{X})^2}{n-1}\right] = \sigma^2 .$$

بنابراین اگر واریانس نمونه را به صورت معادله ۲-۶۸ تعریف کنیم، آنگاه $E(s^2) = \sigma^2$ است، در نتیجه واریانس نمونه تخمین ناریبی از واریانس جامعه خواهد بود.

سازگاری

در مقدمه بحث خصوصیات حدی دیدیم که به ازای افزایش حجم نمونه و میل آن به سمت بی نهایت، معمولاً شکل توزیع تخمین زنده به مرحله فروپاشی می رسد. تخمین زنده $\hat{\theta}$ را در نظر می گیریم. نقطه ای را که در آن توزیع $\hat{\theta}$ فرو می پاشد، نقطه احتمال حدی θ^* می گوئیم و آن را با علامت $\text{Plim } \hat{\theta}$ نشان می دهیم. فرض کنید θ^* نقطه احتمال حدی $\hat{\theta}$ است، البته نمی دانیم که آیا θ^* با مقدار واقعی پارامتر θ برابر است یا خیر. معادله

$$\text{Plim } \hat{\theta} = \theta^*$$

در واقع بدین معنی است که احتمال برابری $\hat{\theta}$ با θ^* در حد برابر یک است؛ یعنی توزیع در نقطه θ^* فرو می پاشد؛ بنابراین می توان نوشت

$$\text{Lim Pr } (\theta^* - \varepsilon \leq \hat{\theta} \leq \theta^* + \varepsilon) = 1$$

که در آن ε عدد مثبت بسیار کوچکی است.

در اینجا می گوئیم تخمین زنده ای سازگار است که توزیع آن در نقطه ای فرو می ریزد که دقیقاً برابر مقدار واقعی پارامتر جامعه است؛ به عبارت دیگر $\hat{\theta}$ یک تخمین زنده سازگار از θ است هر گاه

$$\text{Plim } \hat{\theta} = \theta \quad (2-70)$$

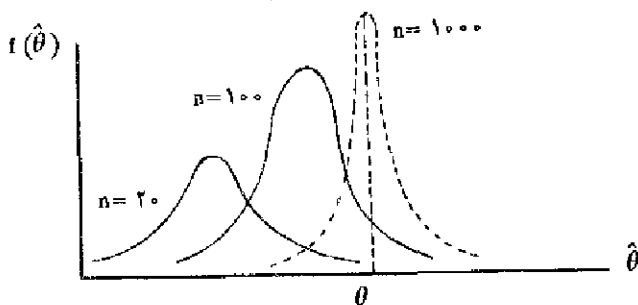
برای اینکه بدانیم یک تخمین زنده سازگار است یا خیر، باید به جهت حرکت مقدار اریب و واریانس آن به ازای افزایش حجم نمونه و میل آن به سمت بی نهایت توجه کنیم. اگر به موازات افزایش حجم نمونه، مقدار اریب و مقدار واریانس هر دو کاهش

پیدا کند و این سیر تا مرحله‌ای ادامه یابد که به ازای $n \rightarrow \infty$ ، مقادیر اریب و واریانس به سمت صفر میل کنند، تخمین‌زننده سازگار خواهد بود. بنابراین احتمال حدی $\hat{\theta} = \theta$ است، هرگاه دو شرط زیر صادق باشد.

$$\lim_{n \rightarrow \infty} E(\hat{\theta}) = \theta, \quad (2-71)$$

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}) = 0.$$

نمودار ۲-۱۵ یک تخمین‌زننده سازگار را نشان می‌دهد.



نمودار ۲-۱۵ سازگاری تخمین‌زننده

بنا بر معادله ۲-۵۷ میانگین مربع خطای $\hat{\theta}$ برابر با واریانس به علاوه مربع اریب است، بنابراین وقتی در سازگاری، اریب و واریانس هر دو به سمت صفر میل کنند، $MSE(\hat{\theta})$ نیز به سمت صفر میل خواهد کرد؛ در نتیجه می‌توان سازگاری تخمین‌زننده را برحسب میانگین مربع خطا به صورت زیر تعریف کرد.

$\hat{\theta}$ یک تخمین‌زننده سازگار از θ است، هرگاه در حد میانگین مربع خطا برابر صفر شود،

$$\text{Plim}_{n \rightarrow \infty} \hat{\theta} = \theta \quad \text{اگر} \quad \lim_{n \rightarrow \infty} MSE(\hat{\theta}) = 0. \quad (2-72)$$

می‌توان ثابت کرد که شرط فوق در واقع شرط لازم است ولی کافی نیست؛ به عبارت دیگر می‌توان یک تخمین‌زننده سازگار یافت که میانگین مربع خطای آن به ازای $n \rightarrow \infty$ به سمت صفر میل نکند (در مورد این موضوع در کتابهای آمار پیشرفته بحث بیشتری شده است).

خصوصیت بسیار مهم تخمین زنده‌های سازگار این است که هر تابع پیوسته‌ای از یک تخمین زنده سازگار، خود سازگار است. این قضیه را در ذیل توضیح می‌دهیم.

قضیه اسلاتسکی

بنابر قضیه اسلاتسکی^۱، اگر احتمال حدی $\hat{\theta} = \theta$ و $g(\hat{\theta})$ یک تابع پیوسته از $\hat{\theta}$ باشد، آنگاه

$$\text{Plim } g(\hat{\theta}) = g(\theta) . \quad (۲-۷۳)$$

در اینجا تنها یادآوری می‌کنیم که این قضیه، کاربردهای بسیاری در روشهای تخمین دارد، که در فصلهای آینده به آن اشاره خواهد شد. بر طبق این قضیه می‌توان گفت که اگر $\hat{\theta}$ یک تخمین زنده سازگار از θ باشد، $\frac{1}{\hat{\theta}}$ نیز تخمین زنده سازگاری از $\frac{1}{\theta}$ و $\log \hat{\theta}$ نیز تخمین زنده سازگاری از $\log \theta$ خواهد بود. این قاعده، در حالت کلی برای خصوصیت ناربیبی صادق نیست؛ یعنی اگر $\hat{\theta}$ یک تخمین زنده ناربیب از θ باشد، بدین معنی نخواهد بود که $\frac{1}{\hat{\theta}}$ یا $\log \hat{\theta}$ نیز تخمینهای ناربیبی از $\frac{1}{\theta}$ یا $\log \theta$ هستند. یکی از امتیازهای احتمال حدی سادگی عملیات آن است. بدون ورود به اثبات قضایا، بعضی از خصوصیات عملیاتی احتمال حدی را مطرح می‌کنیم. احتمال حدی یک عدد ثابت، عدد ثابت است، یعنی $\text{Plim } a = a$. همچنین داریم

$$\text{Plim } g(\hat{\theta})^2 = (\text{Plim } \hat{\theta})^2 ,$$

$$\text{Plim } g(\hat{\theta})^{-1} = (\text{Plim } \hat{\theta})^{-1}$$

اگر $\hat{\theta}_1$ و $\hat{\theta}_2$ به ترتیب تخمینهای سازگاری از θ_1 و θ_2 باشند، $f_1(\hat{\theta}_1)$ و $f_2(\hat{\theta}_2)$ نیز تخمین سازگاری از $f_1(\theta_1)$ و $f_2(\theta_2)$ خواهد بود که در آن f_1 و f_2 هر تابع پیوسته‌ای می‌تواند فرض شود. بنابراین برای مثال داریم

$$\text{Plim} (\hat{\theta}_1^2 \log \hat{\theta}_2) = \theta_1^2 \log \theta_2 ,$$

$$\text{Plim} (\hat{\theta}_1 \hat{\theta}_2) = \text{Plim} \hat{\theta}_1 \text{Plim} \hat{\theta}_2 = \theta_1 \theta_2 ,$$

$$\text{Plim} \frac{\hat{\theta}_1}{\hat{\theta}_2} = \frac{\text{Plim} \hat{\theta}_1}{\text{Plim} \hat{\theta}_2} = \frac{\theta_1}{\theta_2} .$$

همچنین داریم

$$\text{Plim} (\hat{\theta}_1 \pm \hat{\theta}_2) = \text{Plim} \hat{\theta}_1 \pm \text{Plim} \hat{\theta}_2 = \theta_1 \pm \theta_2 .$$

امتیاز عملیاتی احتمال حدی بر امید ریاضی به وضوح روشن است. برای مثال، می‌دانیم $E(\hat{\theta}_1) = E(\hat{\theta}_2)$ یا $E(\frac{\hat{\theta}_1}{\hat{\theta}_2}) \neq \frac{E(\hat{\theta}_1)}{E(\hat{\theta}_2)}$ فقط وقتی برابر $E(\hat{\theta}_1) = E(\hat{\theta}_2)$ است که $\hat{\theta}_1$ و $\hat{\theta}_2$ از یکدیگر مستقل باشد، در حالی که این شرط برای احتمال حدی ضروری نیست.

کارایی حدی

آخرین خصوصیت حدی که در اینجا مطرح می‌کنیم، خصوصیت کارایی حدی^۱ است که به پراکندگی توزیع حدی یک تخمین‌زننده مربوط می‌شود. این خصوصیت معمولاً برای تخمین‌زننده‌هایی تعریف می‌شود که میانگین حدی و واریانس حدی محدود و معینی دارند. کارایی حدی در واقع معیار خوبی است که می‌توان با آن از بین گروهی از تخمین‌زننده‌های نااریب حدی، مناسبترین را انتخاب کرد. می‌دانیم توزیع تخمین‌زننده‌های سازگار به ازای $n \rightarrow \infty$ در نقطه‌ای فرو می‌پاشد که برابر مقدار واقعی پارامتر جامعه است؛ بنابراین اگر تعداد بسیاری از تخمین‌زننده‌ها این خصوصیت را داشته باشند آنگاه تخمین‌زننده‌ای بر بقیه برتری دارد که سریعتر به این نقطه فروپاشی می‌رسد. واضح است که توزیع حدی این تخمین‌زننده کمترین واریانس را خواهد داشت؛ زیرا می‌دانیم، توزیع حدی در واقع صورتی از توزیع است که $\hat{\theta}$ قبل از نیل به

1. Asymptotic Efficiency

مرحله فروپاشی به آن رسیده است؛ و تخمین زنده‌ای که حداقل واریانس را دارد به طور طبیعی زودتر از تخمین زنده‌های سازگار دیگر به این مرحله خواهد رسید. با توجه به تعاریف فوق، می‌توان تعریف دقیقتری از کارایی حدی ارائه کرد.

$\hat{\theta}$ یک تخمین زنده کارای حدی از θ است اگر تمام شرایط زیر صادق باشد.

۱. $\hat{\theta}$ یک توزیع حدی با میانگین و واریانس حدی محدود و معین باشد؛

۲. $\hat{\theta}$ سازگار باشد؛

۳. هیچ تخمین زنده سازگار دیگری از θ را نتوان یافت که واریانس حدی

کوچکتری از $\hat{\theta}$ داشته باشد.

مثال ۲-۱۵ فرض کنید متغیر تصادفی X_1 دارای توزیع نرمال با میانگین μ و واریانس σ^2 است. می‌خواهیم μ را از یک نمونه تصادفی یا مشاهدات X_1, X_2, \dots, X_n تخمین بزنیم. سه تخمین زنده به شرح زیر پیشنهاد شده است،

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i ,$$

$$\hat{\mu} = \frac{1}{n+1} \sum_{i=1}^n X_i ,$$

$$\bar{\mu} = \frac{1}{2} X_1 + \frac{1}{2n} \sum_{i=2}^n X_i .$$

خصوصیات مطلوب تخمین زنده‌ها را در مورد هر یک از این سه تخمین زنده بررسی کنید.

ابتدا خصوصیات نااریبی، کارایی و بهترین تخمین زنده نااریب خطی را در حالت نمونه‌های محدود مطالعه می‌کنیم، آنگاه فرض می‌شود که «به سمت بی‌نهایت میل می‌کند. در این حالت خصوصیات نااریبی حدی، سازگاری و سرانجام کارایی حدی در مورد هر یک از سه تخمین زنده فوق به طور خلاصه بررسی خواهد شد.

نااریبی

از \bar{X} ، $\hat{\mu}$ و $\bar{\mu}$ به ترتیب امید ریاضی می‌گیریم،

$$E(\bar{X}) = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum E(X_i),$$

$$= \frac{1}{n} \sum \mu = \frac{1}{n} \cdot n\mu = \mu,$$

یعنی \bar{X} یک تخمین زنده ناریب از μ است.

$$E(\bar{\mu}) = E\left[\frac{1}{n+1} \sum_{i=1}^n X_i\right] = \frac{1}{n+1} \sum E(X_i),$$

$$= \frac{n}{n+1} \mu,$$

یعنی $\bar{\mu}$ تخمینی از μ و دارای اریب است.

$$E(\bar{\mu}) = E\left[\frac{1}{\gamma} X_1 + \frac{1}{\gamma n} \sum_{i=1}^n X_i\right],$$

$$= \frac{1}{\gamma} E(X_1) + \frac{1}{\gamma n} \sum_{i=1}^n E(X_i),$$

$$= \frac{1}{\gamma} \mu + \left(\frac{n-1}{\gamma n}\right) \mu = \frac{\gamma n - 1}{\gamma n} \mu,$$

یعنی $\bar{\mu}$ اریب دارد.

کارایی

می دانیم کارایی فقط برای تخمین زنده های ناریب تعریف می شود. با توجه به اینکه $\bar{\mu}$ و \bar{X} هر دو اریب دارند فقط \bar{X} می تواند به عنوان مورد مناسبی برای کارایی مطرح شود. واریانس \bar{X} برابر با $\frac{\sigma^2}{n}$ است. همچنین می دانیم که نمونه مفروض از یک جامعه نرمال گرفته شده است. بدون اثبات می گوئیم که با استفاده از نامساوی کرامر - راثو می توان نشان داد که \bar{X} یک تخمین زنده کارآ از μ است.

بهترین تخمین زنده ناریب خطی

ملاحظه می شود که فقط \bar{X} می تواند مورد مناسب این خصوصیت باشد؛ چون دو تخمین زنده دیگر اریب دارند. \bar{X} خصوصیت خطی بودن را نیز دارا است؛ زیرا

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} X_1 + \frac{1}{n} X_2 + \dots + \frac{1}{n} X_n .$$

در قسمت قبل دیدیم که \bar{X} کارآست، یعنی حداقل واریانس را در میان تمام تخمین‌زننده‌های ناریب دیگر از جمله خطی دارد. بنابراین \bar{X} بهترین تخمین‌زننده ناریب خطی است. در اینجا به بررسی خصوصیات حدی در مثال مذکور می‌پردازیم.

ناریب حدی

کافی است از امید ریاضی تخمین‌زننده‌ها حد بگیریم.

$$\lim_{n \rightarrow \infty} E(\bar{X}) = \lim_{n \rightarrow \infty} \mu = \mu ,$$

یعنی \bar{X} ناریب حدی است.

$$\lim_{n \rightarrow \infty} E(\hat{\mu}) = \lim_{n \rightarrow \infty} \frac{n}{n+1} \mu = \mu ,$$

یعنی $\hat{\mu}$ نیز ناریب حدی است.

$$\lim_{n \rightarrow \infty} E(\tilde{\mu}) = \lim_{n \rightarrow \infty} \left(\frac{2n-1}{2n} \right) \mu = \mu ,$$

و بدین ترتیب $\tilde{\mu}$ نیز ناریب حدی خواهد بود.

سازگاری

چون \bar{X} ناریب است، پس

$$MSE(\bar{X}) = \text{Var}(\bar{X}) = \frac{\sigma^2}{n} .$$

اگر از میانگین مربع خطای (\bar{X}) حد بگیریم، داریم

$$\lim_{n \rightarrow \infty} MSE(\bar{X}) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0 ,$$

یعنی \bar{X} سازگار است.

برای بررسی سازگاری $\hat{\mu}$ ابتدا $MSE(\hat{\mu})$ را به دست می‌آوریم. با توجه به معادله ۲-۵۷

داریم

$$\text{MSE}(\hat{\mu}) = \text{Var}(\hat{\mu}) + (\hat{\mu} \text{ اریب})^2$$

بنابراین

$$\text{MSE}(\hat{\mu}) = \text{Var}\left[\frac{1}{n+1} \sum_{i=1}^n X_i\right] + \left[\left(\frac{n}{n+1}\right)\mu - \mu\right]^2$$

$$= \left(\frac{1}{n+1}\right)^2 \sum_{i=1}^n \text{Var}(X_i) + \left(\frac{-1}{n+1}\right)^2 \mu^2$$

$$= \frac{n\sigma^2}{(n+1)^2} + \frac{1}{(n+1)^2} \mu^2 = \frac{n\sigma^2 + \mu^2}{(n+1)^2}$$

از $\text{MSE}(\hat{\mu})$ حد گرفته، داریم

$$\lim_{n \rightarrow \infty} \text{MSE}(\hat{\mu}) = \lim_{n \rightarrow \infty} \frac{n\sigma^2 + \mu^2}{(n+1)^2} = 0$$

یعنی $\hat{\mu}$ سازگار است. به همین ترتیب برای $\bar{\mu}$ عمل می‌کنیم،

$$\text{MSE}(\bar{\mu}) = \text{Var}(\bar{\mu}) + (\bar{\mu} \text{ اریب})^2$$

بنابراین

$$\text{MSE}(\bar{\mu}) = \text{Var}\left[\frac{1}{\xi} X_1 + \frac{1}{\xi n} \sum_{i=2}^n X_i\right] + \left[\left(\frac{\xi n - 1}{\xi n}\right)\mu - \mu\right]^2$$

$$= \frac{1}{\xi} \text{Var}(X_1) + \left(\frac{1}{\xi n}\right)^2 \sum_{i=2}^n \text{Var}(X_i) + \left(\frac{-1}{\xi n}\right)^2 \mu^2$$

$$= \frac{n^2 \sigma^2 + (n-1)\sigma^2 + \mu^2}{\xi n^2} = \frac{(n^2 + n - 1)\sigma^2 + \mu^2}{\xi n^2}$$

از $\text{MSE}(\bar{\mu})$ حد می‌گیریم، داریم

$$\lim_{n \rightarrow \infty} \text{MSE}(\bar{\mu}) = \frac{\sigma^2}{\xi} \neq 0$$

چون $\lim MSE(\hat{\mu})$ برابر صفر نیست، پس $\hat{\mu}$ یک تخمین زنده سازگار نخواهد بود.

کارایی حدی

قط \bar{X} و $\hat{\mu}$ موردهای مناسبی برای کارایی حدی هستند؛ زیرا شرط سازگاری را دارند. چون \bar{X} برای هر حجمی از نمونه کارآست؛ برای حالتی که حجم نمونه به سمت بی نهایت میل کند نیز کارآ خواهد بود؛ بنابراین \bar{X} کارایی حدی دارد. در مورد $\hat{\mu}$ ، می دانیم

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \frac{1}{(n+1)^2} \sum_{i=1}^n \text{Var}(X_i) , \\ &= \frac{1}{(n+1)^2} \cdot n \sigma^2 = \frac{n^2}{(n+1)^2} \cdot \frac{\sigma^2}{n} . \end{aligned}$$

از $\text{Var}(\hat{\mu})$ حد می گیریم،

$$\begin{aligned} \lim_{n \rightarrow \infty} \text{Var}(\hat{\mu}) &= \lim_{n \rightarrow \infty} \left(\frac{n}{n+1} \right)^2 \cdot \frac{\sigma^2}{n} , \\ &= \lim_{n \rightarrow \infty} \left(\frac{n}{n+1} \right)^2 \text{Var}(\bar{X}) , \end{aligned}$$

چون به ازای $n \rightarrow \infty$ ، مقدار $\frac{n}{n+1}$ به سمت یک میل می کند

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\mu}) = \text{Var} \bar{X} .$$

با توجه به اینکه \bar{X} کارایی حدی دارد و واریانس $\hat{\mu}$ در حد برابر واریانس \bar{X} شده است نتیجه می گیریم که $\hat{\mu}$ نیز کارایی حدی دارد. نتایج ششگانه فوق را می توان در جدول ۲-۳ خلاصه کرد.

جدول ۲.۳ خصوصیات مطلوب برای سه تخمین زنده

$\bar{\mu}$	$\hat{\mu}$	\bar{X}	خصوصیات مطلوب
-	-	+	نارایی
-	-	+	کارآیی
-	-	+	BLUE
+	+	+	نارایی حدی
-	+	+	سازگاری
-	+	+	کارآیی حدی

نتیجه کلی این که \bar{X} مناسبترین تخمین زنده برای نمونه‌های محدود است، اما برای نمونه‌های بزرگ، \bar{X} و $\hat{\mu}$ به یک اندازه مناسبند. تخمین زنده سوم، یعنی $\bar{\mu}$ ، هیچ خصوصیت مطلوبی بجز نارایی حدی ندارد.

۲-۶ خصوصیات مطلوب تخمین زنده‌های حداقل مربعات معمولی: قضیه گاس-مارکف*

در قسمت ۲-۵ به بررسی خصوصیات یک تخمین زنده در حالت کلی پرداختیم و به این نتیجه رسیدیم که خصوصیات مطلوب یک تخمین زنده را می‌توان در حالت کلی به صورت بهترین تخمین زنده نارایی خطی جمع‌بندی کرد. در واقع یک تخمین زنده، موقعی مطلوب است که اولاً تابع خطی از مشاهدات متغیر درون‌زا باشد، ثانیاً نارایی باشد و بالاخره در مقایسه با تخمین زنده‌های خطی نارایی دیگر، از حداقل واریانس برخوردار باشد. در مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i$$

دیدیم که تخمین زنده‌های $\hat{\alpha}$ و $\hat{\beta}$ ، با توجه به معادله‌های ۱-۲۵ و ۱-۲۶ عبارتند از

۱. علامت * بدین معنی است که سطح موضوع از دوره کارشناسی بالاتر است.

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} \quad , \quad \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

بنا بر قضیه گاس - مارکف^۱ می توان نشان داد که $\hat{\alpha}$ و $\hat{\beta}$ خصوصیات مطلوب بهترین تخمین زنده ناریب خطی را دارند. اثبات این قضیه را در ذیل یا در روش مطرح می کنیم.

۱. روش مستقیم

در این روش باید نشان دهیم که تخمین زنده های $\hat{\alpha}$ و $\hat{\beta}$ در معادله های ۱-۲۵ و ۱-۲۶ خصوصیات بهترین تخمین زنده ناریب خطی را دارا هستند. بحث را روی $\hat{\beta}$ متمرکز می کنیم؛ زیرا برای $\hat{\alpha}$ می توان از استدلالی مشابه، استفاده کرد.

اولاً، برای اثبات اینکه $\hat{\beta}$ یک تابع خطی از مشاهدات Y_i است، کافی است معادله ۲-۲۳ را یک بار دیگر بنویسیم،

$$\hat{\beta}_{OLS} = w_1 Y_1 + w_2 Y_2 + \dots + w_n Y_n = \sum w_i Y_i$$

که در آن w_i با توجه به معادله ۲-۳ عبارت است از

$$w_i = \frac{x_i}{\sum x_i^2}$$

بنابراین، با توجه به ثابت بودن مقادیر w_i ، می توان نتیجه گرفت که $\hat{\beta}_{OLS}$ یک تابع خطی برحسب Y_i است.

ثانیاً، برای اثبات ناریبی $\hat{\beta}_{OLS}$ به معادله ۲-۵ رجوع می کنیم،

$$E(\hat{\beta}_{OLS}) = \beta$$

ثالثاً، باید نشان دهیم که اگر هر تخمین زنده خطی ناریب دیگری را در نظر بگیریم، واریانس آن بیشتر از واریانسی است که در معادله ۲-۱۶ یعنی

$$\text{Var}(\hat{\beta}_{OLS}) = \frac{\sigma^2}{\sum x_i^2}$$

۱. شکل اولیه این قضیه را ابتدا گاس در سال ۱۸۲۱ اثبات کرد. به پیوست «د» مراجعه شود.

برای $\hat{\beta}_{OLS}$ به دست آورده‌ایم. برای اثبات، یک تخمین زنده خطی ناریب دیگر از β مثلاً $\hat{\beta}_*$ را در نظر می‌گیریم. چون $\hat{\beta}_*$ خطی است پس

$$\hat{\beta}_* = \sum d_i Y_i \quad (2.74)$$

و با توجه به اینکه $\hat{\beta}_*$ ناریب است

$$E(\hat{\beta}_*) = \beta \quad (2.75)$$

مدل $Y_i = \alpha + \beta X_i + U_i$ را در معادله ۲.۷۴ قرار داده خواهیم داشت

$$\hat{\beta}_* = \alpha \sum d_i + \beta \sum d_i X_i + \sum d_i U_i .$$

از دو طرف رابطه فوق امید ریاضی می‌گیریم،

$$E(\hat{\beta}_*) = \alpha \sum d_i + \beta \sum d_i X_i .$$

برای اینکه $\hat{\beta}_*$ یک تخمین زنده ناریب از β شود، ضروری است رابطه‌های زیر برقرار باشد،

$$\sum d_i = 0 \quad , \quad \sum d_i X_i = 1 \quad (2.76)$$

بر اساس فرض استقلال مقادیر Y_i از یکدیگر و نیز فرض واریانس همسانی، اگر از دو طرف معادله ۲.۷۴ واریانس بگیریم، با توجه به $\text{Var}(Y) = \sigma^2$ ، خواهیم داشت

$$\text{Var}(\hat{\beta}_*) = \sum d_i^2 \sigma^2 = \sigma^2 \sum d_i^2 \quad (2.77)$$

باید نشان دهیم که واریانس $\hat{\beta}_*$ که از معادله ۲.۷۷ به دست می‌آید از مقدار واریانس $\hat{\beta}_{OLS}$ کمتر است (قبلاً آن را از معادله ۲.۱۶ به دست آوردیم). در معادله ۲.۲۳ دیدیم که $\hat{\beta}_{OLS} = \sum w_i Y_i$ است. حال رابطه زیر را تعریف می‌کنیم،

$$d_i = w_i + (d_i - w_i) .$$

دو طرف رابطه فوق را می‌گذرد کرده و برای تمام مقادیر t جمع می‌کنیم،

$$\sum d_t^2 = \sum w_t^2 + \sum (d_t - w_t)^2 + 2 \sum w_t (d_t - w_t) . \quad (2.78)$$

نشان می‌دهیم که

$$\sum w_t d_t = \frac{1}{\sum x_t^2} . \quad (2.79)$$

برای اثبات، اگر مقدار w_t از معادله ۲-۳، یعنی $w_t = \frac{x_t}{\sum x_t^2}$ را در طرف چپ رابطه ۲-۷۸ قرار دهیم، خواهیم داشت

$$\sum w_t d_t = \sum \frac{d_t x_t}{\sum x_t^2} = \frac{\sum d_t x_t}{\sum x_t^2} . \quad (2.80)$$

با استفاده از رابطه ۲-۷۶ داریم

$$\begin{aligned} \sum d_t x_t &= \sum d_t (x_t + \bar{x}) , \\ &= \sum d_t x_t + \bar{x} \sum d_t = \sum d_t x_t = 1 , \end{aligned}$$

که با جایگزینی نتیجه فوق در معادله ۲-۸۰ خواهیم داشت

$$\sum w_t d_t = \frac{1}{\sum x_t^2} ,$$

که دقیقاً همان رابطه ۲-۷۹ است. از طرف دیگر، طبق معادله ۲-۳ داریم

$$w_t^2 = \frac{x_t^2}{(\sum x_t^2)^2} ,$$

یا

$$\sum w_t^2 = \frac{\sum x_t^2}{(\sum x_t^2)^2} = \frac{1}{\sum x_t^2} . \quad (2.81)$$

بنابراین

$$\sum w_t (d_t - w_t) = \sum w_t d_t - \sum w_t^2 ,$$

که با استفاده از معادله‌های ۲-۷۹ و ۲-۸۱ خواهیم داشت

$$\sum w_i (d_i - w_i) = \frac{1}{\sum x_i^2} - \frac{1}{\sum x_i^2} = 0 .$$

نتیجه فوق را در معادله ۲-۷۸ قرار می‌دهیم،

$$\sum d_i^2 = \sum w_i^2 + \sum (d_i - w_i)^2 .$$

دو طرف رابطه فوق را در σ^2 ضرب می‌کنیم،

$$\sigma^2 \sum d_i^2 = \sigma^2 \sum w_i^2 + \sigma^2 \sum (d_i - w_i)^2 . \quad (2-82)$$

در معادله ۲-۱۴ دیدیم که

$$\text{Var}(\hat{\beta}_{OLS}) = \sigma^2 \sum w_i^2 .$$

با استفاده از رابطه فوق و معادله ۲-۷۷ و معادله ۲-۸۲ را به صورت زیر می‌نویسیم،

$$\text{Var}(\hat{\beta}_*) = \text{Var}(\hat{\beta}_{OLS}) + \sigma^2 \sum (d_i - w_i)^2 .$$

جمله $\sigma^2 \sum (d_i - w_i)^2$ همواره مثبت و بزرگتر از صفر است، بنابراین

$$\text{Var}(\hat{\beta}_{OLS}) < \text{Var}(\hat{\beta}_*) .$$

به ازای $d_i = w_i$ خواهیم داشت

$$\text{Var}(\hat{\beta}_*) = \text{Var}(\hat{\beta}_{OLS}) .$$

اما واضح است که در این حالت تخمین زنده $\hat{\beta}_*$ دقیقاً همان $\hat{\beta}_{OLS}$ خواهد بود.

۲. روش غیرمستقیم

روش دیگر، اثبات خصوصیت بهترین تخمین زنده ناریب خطی برای تخمین زنده‌های حداقل مربعات معمولی این است که از ابتدا $\hat{\alpha}$ و $\hat{\beta}$ را چنان تخمین بزنیم که نه تنها

مجموع مربعات پسماند، یعنی $\sum e_i^2$ ، را حداقل کند؛ بلکه خصوصیت بهترین تخمین‌زننده ناریب خطی را نیز دارا باشد. این روش، بسیار عمومی است و می‌تواند کاربردهای بسیاری در مسائل مختلف داشته باشد. مانند گذشته بحث را روی $\hat{\beta}$ متمرکز می‌کنیم؛ زیرا به راحتی می‌توان نتایج مشابهی را برای $\hat{\alpha}$ استنتاج کرد.

مدل رگرسیون $Y_i = \alpha + \beta X_i + U_i$ را در نظر می‌گیریم. با توجه به معادله ۱-۲۵

می‌دانیم

$$\hat{\beta}_{OLS} = \frac{\sum x_i y_i}{\sum x_i^2}$$

مقدار $\hat{\beta}_{OLS}$ که از رابطه فوق به دست می‌آید یقیناً مقدار $\sum e_i^2$ را حداقل می‌کند. در اینجا یک تخمین‌زننده غیر حداقل مربعات معمولی از β ، مانند $\hat{\beta}$ را چنان تعریف می‌کنیم که بهترین تخمین‌زننده ناریب خطی باشد. می‌خواهیم ثابت کنیم که $\hat{\beta}$ دقیقاً برابر $\hat{\beta}_{OLS}$ خواهد بود.

با توجه به معادله‌های ۲-۷۶ و ۲-۷۷ می‌دانیم لازمه اینکه $\hat{\beta}$ بهترین تخمین‌زننده ناریب خطی باشد این است که:

$$\sum d_i = 0 \quad , \quad \sum d_i X_i = 1 \quad , \quad \text{Var}(\hat{\beta}_*) = \sigma^2 \sum d_i^2$$

در اینجا اول اینکه، باید ضرایب d_i را چنان تعیین کنیم که سه رابطه فوق به طور همزمان برقرار باشد و دوم اینکه، باید نشان دهیم که در چنین حالتی $\hat{\beta}$ دقیقاً همان $\hat{\beta}_{OLS}$ خواهد بود. کافی است d_i را چنان تعیین کنیم که معادله

$$\text{Var}(\hat{\beta}_*) = \sigma^2 \sum d_i^2$$

به شرط اینکه رابطه‌های زیر حداقل شود،

$$\sum d_i = 0 \quad , \quad \sum d_i X_i = 1$$

ملاحظه می‌شود که این مسأله در واقع یک حداقل‌سازی مقید است که به راحتی می‌توان آن را به استفاده از ضرایب لاگرانژ حل کرد. تابع زیر را تعریف می‌کنیم،

$$\varphi = \sigma^2 \sum d_i^2 - \lambda_1 (\sum d_i) - \lambda_2 (\sum d_i X_i - 1) .$$

باید از φ نسبت به d_i و λ_1 و λ_2 مشتق جزئی گرفته و آنها را مساوی صفر قرار دهیم؛

$$\frac{\partial \varphi}{\partial d_1} = 0 , \frac{\partial \varphi}{\partial d_2} = 0 , \dots , \frac{\partial \varphi}{\partial d_n} = 0 \text{ و } \frac{\partial \varphi}{\partial \lambda_1} = 0 , \frac{\partial \varphi}{\partial \lambda_2} = 0 .$$

بنابراین خواهیم داشت،

$$\begin{aligned} 2 d_1 \sigma^2 - \lambda_1 - \lambda_2 X_1 &= 0 , \\ 2 d_2 \sigma^2 - \lambda_1 - \lambda_2 X_2 &= 0 , \\ &\vdots \\ 2 d_n \sigma^2 - \lambda_1 - \lambda_2 X_n &= 0 , \quad (2.13) \\ - \sum d_i &= 0 , \\ - (\sum d_i X_i) + 1 &= 0 . \end{aligned}$$

ملاحظه می شود که $(n + 2)$ معادله داریم و مجهولات ما عبارتند از: n مقدار d_i و نیز λ_1 و λ_2 . می توان n معادله اول را به صورت زیر نوشت،

$$\begin{aligned} d_1 &= \frac{1}{2 \sigma^2} (\lambda_1 + \lambda_2 X_1) , \\ d_2 &= \frac{1}{2 \sigma^2} (\lambda_1 + \lambda_2 X_2) , \\ &\vdots \\ d_n &= \frac{1}{2 \sigma^2} (\lambda_1 + \lambda_2 X_n) . \end{aligned} \quad (2.14)$$

معادله های فوق را با هم جمع می کنیم، خواهیم داشت

$$\sum_{i=1}^n d_i = \frac{1}{2 \sigma^2} (\lambda_1 n + \lambda_2 \sum X_i) . \quad (2.15)$$

همچنین اگر معادله اول از معادله های ۲-۱۴ را در X_1 و معادله دوم را در X_2 و معادله سوم را در X_3 و ... ضرب کنیم و نتایج را با یکدیگر جمع کنیم، خواهیم داشت.

$$\sum d_t X_t = \frac{1}{\gamma \sigma^2} [\lambda_1 \sum X_t + \lambda_2 \sum X_t'] \quad (2.86)$$

معادله‌های ۲.۸۵ و ۲.۸۶ را در دو معادله آخر سیستم معادله‌های ۲.۸۳ جایگزین کرده داریم

$$-\frac{1}{\gamma \sigma^2} (\lambda_1 n + \lambda_2 \sum X_t) = 0 ,$$

$$-\frac{1}{\gamma \sigma^2} [\lambda_1 \sum X_t + \lambda_2 \sum X_t'] = -1 .$$

دو معادله فوق را حل کرده و λ_1 و λ_2 را به صورت زیر به دست می‌آوریم،

$$\lambda_1 = \frac{-\gamma \sigma^2 \sum X_t}{n \sum X_t' - (\sum X_t)^2} , \quad \lambda_2 = \frac{\gamma n \sigma^2}{n \sum X_t' - (\sum X_t)^2} .$$

مقادیر λ_1 و λ_2 از معادله‌های فوق را در سیستم معادله‌های ۲.۸۴ قرار می‌دهیم و معادله‌های حاصل را برای d_1 ، d_2 ، تا d_n حل می‌کنیم. خواهیم داشت

$$d_t = \frac{-(\sum X_t) + n X_t}{n (\sum X_t') - (\sum X_t)^2} = \frac{(X_t - \bar{X})}{\sum (X_t - \bar{X})^2} ,$$

یا

$$d_t = \frac{x_t}{\sum x_t'^2} , \quad , \quad t = 1, 2, \dots, n . \quad (2.87)$$

مقادیر d_t که از رابطه فوق به دست می‌آید، $\hat{\beta}_*$ را تا آریب و خطی و واریانس آن را حداقل می‌کند. مقدار d_t از معادله فوق را در معادله ۲.۷۴ قرار می‌دهیم،

$$\hat{\beta}_* = \sum d_t Y_t = \sum \frac{x_t Y_t}{\sum x_t'^2} = \frac{\sum x_t Y_t}{\sum x_t'^2} ,$$

یا

$$\hat{\beta}_* = \frac{\sum x_t Y_t}{\sum x_t'^2} ,$$

که دقیقاً برابر تخمین زننده $\hat{\beta}$ است و با استفاده از روش حداقل مربعات معمولی به دست

می آید؛ یعنی اگر قرار است $\hat{\beta}_0$ خصوصیت بهترین تخمین زنده ناریب خطی را داشته باشد، باید برابر $\hat{\beta}_{OLS}$ باشد.

روش فوق بسیار عمومی است، زیرا؛ اولاً، خصوصیت بهترین تخمین زنده ناریب خطی را تضمین می کند؛ ثانیاً، مقدار $\hat{\beta}$ را نیز معین می نماید؛ و ثالثاً می تواند واریانس $\hat{\beta}$ را نیز محاسبه کند. برای به دست آوردن واریانس $\hat{\beta}_0$ ، ابتدا معادله ۲-۷۷ را می نویسیم،

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \sum d_i^2 .$$

برای محاسبه $\sum d_i^2$ از معادله ۲-۸۷ استفاده می شود. دو طرف این معادله را در d_i ضرب کرده و برای تمام مشاهدات جمع می کنیم. خواهیم داشت

$$\sum d_i^2 = \frac{-(\sum X_i) \sum d_i + n (\sum d_i X_i)}{n \sum X_i^2 - (\sum X_i)^2} .$$

از دو معادله آخر سیستم معادله های ۲-۸۳ داریم

$$\sum d_i = 0 \quad , \quad \sum d_i X_i = 1 \quad ,$$

بنابراین $\sum d_i^2$ را می توان به صورت زیر نوشت،

$$\sum d_i^2 = \frac{n}{n \sum X_i^2 - (\sum X_i)^2} = \frac{1}{\sum x_i^2} .$$

بدین ترتیب واریانس $\hat{\beta}_0$ عبارت است از

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{\sum x_i^2} ,$$

که دقیقاً برابر $\text{Var}(\hat{\beta}_{OLS})$ است که در معادله ۲-۱۶ به دست آوردیم.

مسائل فصل دوم

۲-۱ برای مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i .$$

مشاهدات زیر مفروض است.

$X_i:$	۲	۳	۱	۵	۹
$Y_i:$	۴	۷	۳	۹	۱۷

- با استفاده از روش حداقل مربعات معمولی، پارامترهای α و β را تخمین بزنید.
- مقادیر پسماند را حساب کنید و نشان دهید $\sum e_i = 0$.
- تغییرات توضیح داده شده (ESS) و تغییرات توضیح داده نشده (RSS) را حساب کنید.
- نشان دهید که تغییرات توضیح داده نشده + تغییرات توضیح داده شده = کل تغییرات (TSS = ESS + RSS).
- r^2 را از چهار روش مختلف به دست آورید.
- با استفاده از r^2 ، مقدار تغییرات توضیح داده نشده را حساب کنید.
- واریانس U_i را تخمین بزنید (محاسبه s^2).
- واریانس $\hat{\alpha}$ و واریانس $\hat{\beta}$ را تخمین بزنید.
- انحراف معیار $\hat{\alpha}$ و انحراف معیار $\hat{\beta}$ را تخمین بزنید.
- فاصله اطمینان ۹۵ درصد را برای α حساب کنید.
- فاصله اطمینان ۹۵ درصد را برای β حساب کنید.
- فرضیه $H_0: \alpha = 0$ را در مقابل فرضیه $H_1: \alpha \neq 0$ در سطح معنی دار ۵ درصد آزمون کنید.

۱۳. فرضیه $H_0: \beta = 0$ را در مقابل فرضیه $H_1: \beta \neq 0$ در سطح معنی داری ۵ درصد آزمون کنید.

۱۴. برای σ_u^2 فاصله اطمینان ۹۵ درصد بسازید.

۱۵. برای این مسأله جدول آنالیز واریانس را تشکیل دهید.

۱۶. فرضیه $H_0: \beta = 0$ را در مقابل فرضیه $H_1: \beta \neq 0$ با استفاده از توزیع F آزمون کنید.

۱۷. معنی داری بودن مدل رگرسیون را در سطح معنی داری ۵ درصد آزمون کنید (آزمون F^2).

۱۸. آماره t را از سه راه مختلف زیر برای فرضیه $H_0: \beta = 0$ محاسبه کرده و تساوی نتایج را تفسیر کنید.

اولاً، فقط با استفاده از F^2 ؛

ثانیاً، با استفاده از استاندارد کردن $\hat{\beta}$ ؛

ثالثاً، فقط با استفاده از F.

۱۹. با داشتن آماره آزمون t ، برای فرضیه $H_0: \beta = 0$ ، مقدار ضریب تعیین R^2 را حساب کنید.

۲-۲ در مدل رگرسیون

$$Y_t = \alpha + \beta X_t + U_t ,$$

نشان دهید که اگر $\bar{X} = 0$ باشد، کوواریانس بین $\hat{\alpha}$ و $\hat{\beta}$ برابر صفر است. آیا می‌توانید این نتیجه را توجیه کنید.

۲-۳ مدل‌های رگرسیون زیر را ملاحظه کنید،

$$Y_t = \alpha_1 + \beta_1 C_t + U_{1t} ,$$

$$C_t = \alpha_2 + \beta_2 Y_t + U_{2t} ,$$

که در آن Y_t و C_t به ترتیب درآمد کل و مصرف کل است. این دو مدل را به صورت زیر

تخمین زده‌ایم.

$$\hat{y}_t = 1/2 c_t ,$$

$$\hat{c}_t = 0/6 y_t ,$$

که در آن \hat{y}_t و c_t و y_t بر حسب انحراف از میانگین هستند. اگر داشته باشیم

$$Y_t \equiv C_t + S_t ,$$

که در آن S_t پس‌انداز است، آنگاه

۱. نسبت انحراف معیار S_t به انحراف معیار Y_t را حساب کنید.

۲. ضریب همبستگی بین S_t و Y_t را به دست آورید.

۳. ضریب همبستگی بین S_t و C_t را محاسبه کنید.

۴-۲. مشاهدات زیر در مورد X_t و Y_t مفروض است. X_t نرخ بیکاری و Y_t درصد

کارگرانی است که از صنایع تولیدی خارج می‌شوند.

t	Y_t	X_t	t	Y_t	X_t
۱	۱/۳	۶/۲	۸	۲/۳	۳/۶
۲	۱/۲	۷/۸	۹	۲/۵	۲/۳
۳	۱/۴	۵/۸	۱۰	۲/۷	۲/۳
۴	۱/۴	۵/۷	۱۱	۲/۱	۵/۶
۵	۱/۵	۵	۱۲	۱/۸	۶/۸
۶	۱/۹	۴	۱۳	۲/۲	۵/۶
۷	۲/۶	۳/۲			

۱. مدل زیر را تخمین یزنید،

$$Y_t = \alpha + \beta X_t + U_t .$$

۲. فاصله اطمینان ۹۵ درصد برای β بسازید.

۳. فرضیه $H_0: \beta = 0$ را در مقابل فرضیه $H_1: \beta \neq 0$ در سطح معنی دار ۵ درصد آزمون کنید.

۴. فاصله اطمینان ۹۰ درصد برای σ^2 بسازید.

۲.۵ مدل رگرسیون زیر مفروض است،

$$Y_t = \alpha + \beta X_t + U_t$$

که در آن U_t جمله اختلال با میانگین صفر و واریانس ثابت است. برای تخمین α و β دو پژوهشگر، مستقل از یکدیگر دو نمونه ۸ تایی به طور تصادفی انتخاب کرده‌اند و با روش حداقل مربعات معمولی پارامترها را تخمین زده‌اند. نتایج تخمین یا مشاهدات هر یک از نمونه‌ها در جدول زیر ارائه شده است.

پژوهشگر اول		پژوهشگر دوم	
Y_t	X_t	Y_t	X_t
۴	۲	۲	۱
۴/۵	۲	۲/۵	۱
۴/۵	۳	۲/۵	۱
۳/۵	۳	۱/۵	۱
۴/۵	۴	۱۱/۵	۱۰
۴/۵	۴	۱۰/۵	۱۰
۵/۵	۴	۱۰/۵	۱۰
۵	۴	۱۱	۱۰
$\hat{Y}_t = 1/875 + 0/750 X_t$ (۱/۲۰) (۰/۳۳۹)		$\hat{Y}_t = 1/5 + 0/970 X_t$ (۰/۲۷) (۰/۰۳۸)	
$r^2 = 0/45$		$r^2 = 0/99$	
$\sigma^2 = 0/48$		$\sigma^2 = 0/48$	

این نکته را توضیح دهید که چرا انحراف معیار $\hat{\beta}$ در مدلی که پژوهشگر اول تخمین زده است از مقدار انحراف معیار $\hat{\beta}$ در مدل دوم بیشتر است.

۲-۶. با توجه به اینکه واریانس $\hat{\beta}$ تابع معکوسی از واریانس X_t است، بهتر است مشاهدات مرکزی X_t را در نمونه حذف کنیم تا به واریانس بیشتری برای X_t و در نتیجه به واریانس کمتری برای $\hat{\beta}$ برسیم. بدین ترتیب به تخمینهای دقیقتری از β خواهیم رسید. آیا موافق چنین نظری هستید؟ چرا.

۲-۷. مدل رگرسیون زیر مفروض است،

$$Y_t = \beta X_t + U_t ,$$

که در آن U_t همه فرضهای کلاسیک را دربر دارد. جمله‌های پسماند را به صورت زیر تعریف می‌کنیم،

$$e_t = Y_t - \hat{\beta} X_t .$$

نشان دهید که

$$E(e_t^2) = \sigma^2 \left(1 - \frac{X_t^2}{\sum X_t^2} \right) , \quad 1.$$

که در آن σ^2 واریانس جمله اختلال است.

$$E(e_t e_s) = \frac{-\sigma^2 X_t X_s}{\sum X_t^2} , \quad s \neq t . \quad 2.$$

۳. در حالت کلی $\sum e_t \neq 0$. این نتیجه را با مدلی که ضریب ثابت دارد، مقایسه

کنید. آیا $E(\sum e_t)$ برابر صفر است؟

۴. نشان دهید که

$$E \left(\frac{\sum e_t^2}{n-1} \right) = \sigma^2 .$$

۲-۸. مدل رگرسیون زیر مفروض است،

$$Y_t = \alpha + \beta X_t + U_t ,$$

که در آن U_t شامل تمام فرضهای کلاسیک است.

۱. نشان دهید که شکل عمومی تخمین‌زنده‌های خطی نااریب به صورت زیر

است،

$$\tilde{\beta} = \frac{\sum z_i y_i}{\sum z_i x_i},$$

Z_i در آن یک متغیر برونزا است که در نمونه گیریهای تکراری ثابت فرض می‌شود. همچنین می‌دانیم $z_i = Z_i - \bar{Z}$.

۲. واریانس صورت عمومی تخمین زنده‌های خطی ناریب را به دست آورید و با واریانس تخمین زنده‌های حداقل مربعات معمولی مقایسه کنید.

۳. نشان دهید کارایی نسبی $\tilde{\beta}$ (در مقایسه با $\hat{\beta}_{OLS}$) تابعی از همبستگی بین x_i و z_i است، به گونه‌ای که وقتی این همبستگی کامل باشد، واریانس $\tilde{\beta}$ با واریانس $\hat{\beta}$ برابر می‌شود.

حل مسائل فصل دوم

۲-۱ برای پاسخ به این سؤال داریم:

۱. تخمینهای α و β را یک بار با استفاده از مشاهدات اصلی و بار دیگر با روش انحراف از میانگین به دست می‌آوریم. ابتدا جدول زیر را تشکیل می‌دهیم.

جدول ۱

X_i	Y_i	$X_i Y_i$	X_i^2	\hat{Y}_i	$e_i = Y_i - \hat{Y}_i$
۲	۴	۸	۴	۴/۵۰	-۰/۵۰
۳	۷	۲۱	۹	۶/۲۵	۰/۷۵
۱	۳	۳	۱	۲/۷۵	۰/۲۵
۵	۹	۴۵	۲۵	۹/۷۵	-۰/۷۵
Σ ۹	۱۷	۱۵۲	۸۱	۱۶/۷۵	۰/۲۵
Σ ۲۰	۴۰	۲۳۰	۱۲۰	۴۰	۰

$$\hat{\beta} = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n \sum X_i^2 - (\sum X_i)^2}$$

$$= \frac{9(230) - 20(40)}{9(120) - (20)^2} = \frac{1150 - 800}{600 - 400} = 1/70$$

$$\hat{\alpha} = \frac{\sum X_i^2 \sum Y_i - \sum X_i \sum X_i Y_i}{n \sum X_i^2 - (\sum X_i)^2}$$

$$= \frac{120(40) - 20(230)}{9(120) - (20)^2} = \frac{4800 - 4600}{600 - 400} = 1$$

در اینجا $\hat{\alpha}$ و $\hat{\beta}$ را با استفاده از مشاهداتی به دست می‌آوریم که برحسب انحراف از

میانگین محاسبه شده است. ابتدا جدول زیر را تشکیل می دهیم.

جدول ۲

x_i	y_i	$x_i y_i$	x_i^2	y_i^2
-۲	-۴	۸	۴	۱۶
-۱	-۱	۱	۱	۱
-۳	-۵	۱۵	۹	۲۵
۱	۱	۱	۱	۱
Σ	$\frac{0}{0}$	$\frac{40}{70}$	$\frac{20}{40}$	$\frac{41}{124}$

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{40}{40} = 1/10 .$$

$$\bar{X} = \frac{0}{0} = 0 \quad , \quad \bar{Y} = \frac{40}{0} = 10 ,$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} = 10 - 1/10 (0) = 10 ,$$

بدین ترتیب تخمین مدل رگرسیون به شرح زیر خواهد بود،

$$\hat{Y}_i = 10 + 1/10 X_i .$$

۲. با استفاده از جدول ۱ ملاحظه می شود که $\sum e_i = 0$.

۳. با استفاده از معادله ۲-۴۵ داریم

$$\begin{aligned} ESS &= \sum \hat{y}_i^2 = \hat{\beta}^2 \sum x_i y_i , \\ &= 1/10 (40) = 4/10 . \end{aligned}$$

مجموع مربعات پسماند (RSS) را از معادله ۲-۴۶ به دست می آوریم،

$$\begin{aligned} RSS &= \sum e_i^2 = \sum y_i^2 - \hat{\beta}^2 \sum x_i y_i , \\ &= 124 - 1/10 (40) = 124 - 4/10 = 120/10 . \end{aligned}$$

۴. برای ثبات: تغییرات توضیح داده نشده + تغییرات توضیح داده شده = کل تغییرات کافی است مقادیر مربوط را جایگزین کنیم،

$$\sum y_i' = \sum \hat{y}_i' + \sum e_i'$$

$$۱۲۴ = ۱۲۲/۵ + ۱/۵ .$$

۵. ابتدا از فرمول ۱-۳۲ استفاده می کنیم. می دانیم ضریب تعیین برابر معذور ضریب همبستگی است،

$$r_{x,y} = \frac{\sum x_i y_i}{\sqrt{\sum x_i'^2 \cdot \sum y_i'^2}}$$

$$= \frac{۷۰}{\sqrt{۴۰ \cdot (۱۲۴)}} = ۰/۹۹۳۹۳۳۲۰۹ ,$$

در نتیجه:

$$r^2 \cong ۰/۹۸۸ .$$

یا استفاده از فرمول ۱-۳۳ داریم

$$r^2 = \frac{ESS}{TSS} = \frac{\sum \hat{y}_i'^2}{\sum y_i'^2}$$

$$= \frac{۱۲۲/۵}{۱۲۴} = ۰/۹۸۸ .$$

روش سوم این است که از فرمول ۱-۳۵ استفاده کنیم،

$$r^2 = \hat{\beta} \frac{\sum x_i y_i}{\sum y_i'^2}$$

$$= ۱/۷۵ \frac{۷۰}{۱۲۴} = \frac{۱۲۲/۵}{۱۲۴} = ۰/۹۸۸ .$$

روش چهارم این است که مبتنی بر فرمول ۱-۴۳ عمل کنیم،

$$r^2 = 1 - \frac{RSS}{TSS} ,$$

$$= 1 - \frac{1/5}{124} = \frac{122/5}{124} = 0.988 .$$

۶. با استفاده از فرمول ۱-۴۳ داریم

$$r^2 = 1 - \frac{RSS}{TSS} ,$$

$$0.988 = 1 - \frac{RSS}{124} ,$$

$$\sum e_i^2 = RSS = 0.012 (124) = 1/488 .$$

ملاحظه می‌شود مقدار دقیق تغییرات توضیح داده نشده (RSS) برابر ۱/۵ است؛ در حالی که این روش به مقدار ۱/۴۸۸ می‌رسیم. دلیل این امر وجود تقریبی است که معمولاً در r^2 وجود دارد.

۷. می‌دانیم

$$s^2 = \hat{\sigma}^2 = \frac{\sum e_i^2}{n-2} ,$$

$$= \frac{1/5}{5-2} = 0.5 .$$

۸. ابتدا واریانس $\hat{\alpha}$ را تخمین می‌زنیم. از فرمول ۲-۲۱ داریم

$$\text{Var}(\hat{\alpha}) = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] ,$$

$$= 0.5 \left[\frac{1}{5} + \frac{16}{40} \right] = \frac{12}{40} = 0.3 .$$

برای تخمین واریانس $\hat{\beta}$ کافی است از معادله ۲-۱۶ استفاده کنیم،

$$\text{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2} ,$$

$$= \frac{0.5}{40} = 0.0125 .$$

۹.

$$SE(\hat{\alpha}) = \sqrt{\text{Var}(\hat{\alpha})} = \sqrt{0/3} = 0/0477,$$

$$SE(\hat{\beta}) = \sqrt{\text{Var}(\hat{\beta})} = \sqrt{0/0125} = 0/1118.$$

۱۰. برای فاصله اطمینان ۹۵ درصد باید از سطح معنی دار ۵ درصد استفاده کنیم؛ بنابراین مساحت ناحیه بحرانی در هر طرف برابر ۰/۰۲۵ خواهد بود. مقدار به دست آمده از جدول t در سطح معنی دار ۲/۵ درصد و با ۳ درجه آزادی برابر است با ۳/۱۸۲؛ بنابراین خواهیم داشت

$$\hat{\alpha} - t_{\alpha/2} SE(\hat{\alpha}) < \alpha < \hat{\alpha} + t_{\alpha/2} SE(\hat{\alpha}),$$

$$1 - 3/182(0/0477) < \alpha < 1 + 3/182(0/0477),$$

$$-0/74 < \alpha < 2/74.$$

۱۱. فاصله اطمینان ۹۵ درصد برای β عبارت است از

$$\hat{\beta} - t_{\alpha/2} SE(\hat{\beta}) < \beta < \hat{\beta} + t_{\alpha/2} SE(\hat{\beta}),$$

$$1/75 - 3/182(0/1118) < \beta < 1/75 + 3/182(0/1118),$$

$$1/39 < \beta < 2/11.$$

۱۲. برای آزمون فرضیه $H_0: \alpha = 0$ ، در مقابل $H_1: \alpha \neq 0$ مقدار آماره آزمون را به دست می آوریم.

$$t = \frac{\hat{\alpha} - 0}{SE(\hat{\alpha})} = \frac{1 - 0}{0/0477} = 1/826.$$

ملاحظه می شود که چون آماره آزمون از مقدار جدول t کمتر است، این آماره معنی دار نبوده، در نتیجه فرضیه H_0 رد نمی شود.

مقدار به دست آمده از جدول $t < t$ آماره آزمون

$$1/826 < 3/182.$$

بنابراین ضریب ثابت در مدل رگرسیون مفروض در سطح معنی دار ۵ درصد معتبر نیست. البته به محاسبه آماره آزمون نیازی نبود؛ زیرا قبلاً فاصله اطمینان را ساخته ایم. در واقع چون $\alpha = 0$ در فاصله اطمینان ۹۵ درصد برای α قرار دارد، فرضیه $H_0: \alpha = 0$ را نمی توان رد کرد.

۱۳. ابتدا آماره آزمون را به دست می آوریم،

$$t = \frac{\hat{\beta} - 0}{SE(\hat{\beta})} = \frac{1/70 - 0}{0/1118} = 10/603 ,$$

و چون این آماره از مقدار جدول $t(3/182)$ بیشتر شده است، بنابراین در ناحیه بحرانی قرار گرفته و فرضیه H_0 رد می شود. مانند بحث بند ۱۲، می توان گفت به محاسبه t نیازی نیست؛ زیرا قبلاً فاصله اطمینان β را ساخته ایم. چون $\beta = 0$ در فاصله اطمینان ۹۵ درصد برای β قرار نمی گیرد؛ پس فرضیه $\beta = 0$ قابل قبول نیست.

۱۴. برای ساختن فاصله اطمینان ۹۵ درصد برای واریانس جامعه σ^2 ابتدا باید دو نقطه در روی محور χ^2 تعیین کنیم که مساحت سطح زیر منحنی از مبدأ مختصات تا آن نقاط به ترتیب برابر $0/025$ و $0/975$ باشد. این نقاط را می توان با ۳ درجه آزادی از جدول χ^2 به دست آورد. خواهیم داشت

$$\chi^2_{0/025}(3) = 9/35 , \quad \chi^2_{0/975}(3) = 0/216 .$$

با توجه به معادله ۲-۴۴ داریم

$$\frac{(n-2)\hat{\sigma}^2}{\chi^2_{0/975}} < \sigma^2 < \frac{(n-2)\hat{\sigma}^2}{\chi^2_{0/025}} .$$

می دانیم

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-2} ,$$

در نتیجه

$$(n-2)\hat{\sigma}^2 = \sum e_i^2 = 1/0 ,$$

بنابراین

$$\frac{1/0}{9/35} < \sigma^2 < \frac{1/0}{0/216} ,$$

$$0/16 < \sigma^2 < 6/16 .$$

۱۵. برای تشکیل جدول آنالیز واریانس به کمیتهای زیر نیاز داریم،

$$ESS , RSS , TSS , ESS/1 , RSS/(n - 2)$$

ملاحظه می شود تمام این کمیتهای را در بندهای قبلی محاسبه کرده ایم؛ بنابراین جدول تجزیه واریانس را به شرح زیر می نویسیم

منبع تغییرات	مجموع مربعات	درجات آزادی	میانگین مربعات یا واریانس
X_i	$ESS = 122/5$	۱	$ESS/1 = 122/5$
e_i	$RSS = 1/5$	۳	$RSS/3 = 0/5$
کل	$TSS = 123$	۴	

۱۶. از معادله ۲-۵۱ مقدار آماره آزمون F را به دست می آوریم، داریم

$$F = \frac{ESS/1}{RSS/(n-2)}$$

$$= \frac{122/5}{0/5} = 240 .$$

در اینجا باید مقدار F از جدول پیوست F را با درجات آزادی ۱ (متعلق به صورت) و ۳ (متعلق به مخرج) و در سطح معنی دار ۵ درصد به دست آورد. خواهیم داشت

$$F_{0/05} (1 و 3) = 10/1 .$$

ملاحظه می شود آماره F در ناحیه بحرانی قرار می گیرد و فرضیه $H_0: \beta = 0$ رد می شود. ۱۷. با استفاده از معادله ۲-۵۴ داریم

$$F = \frac{(n-2) r^2}{1-r^2} ,$$

$$= \frac{(3) \cdot 0/988}{1 - 0/988} = 247 .$$

مقدار F که با روش t^2 به دست آمده است کمی با مقدار دقیق F که در بند ۱۶ محاسبه شده متفاوت است. علت این امر، وجود تقریب در محاسبه t^2 است. کافی است که آماره آزمون F را با مقدار جدول F مقایسه کنیم. می‌دانیم

$$F_{0/10}(1 و 3) = 10/1 ,$$

نتیجه می‌گیریم که آماره آزمون در ناحیه بحرانی قرار می‌گیرد و فرضیه H_0 رد می‌شود؛ یعنی تخمین مدل رگرسیون معنی‌دار است.

۱۸. اولاً، برای محاسبه t با روش t^2 باید از معادله 2.55 استفاده کنیم،

$$t^2 = \frac{(n-2) r^2}{1-r^2} ,$$

$$= \frac{3(0/988)}{1-0/988} = \frac{2/974}{0/012} = 247 ,$$

$$t = \pm \sqrt{247} = \pm 15/71 .$$

ثانیاً، محاسبه t با روش استاندارد کردن $\hat{\beta}$ ، در بند ۱۳ انجام شده است.

$$t = \frac{\hat{\beta}}{SE(\hat{\beta})} = \frac{1/70}{0/1118} = 15/60 .$$

ثالثاً، می‌دانیم

$$t = \pm \sqrt{F} .$$

در بند ۱۶ مقدار F را برای فرضیه $\beta = 0$ برابر $F = 245$ به دست آوردیم؛ بنابراین

$$t = \sqrt{245} = \pm 15/60 .$$

اگر تقریبهای موجود در محاسبه t^2 را نادیده بگیریم، مقادیری که برای t با سه

روش فوق به دست آمده است با یکدیگر برابر است تفسیر این نتایج بدین ترتیب است که سه راه حل فوق در واقع صورتهای مختلف بیان یک واقعیت است؛ به عبارت دیگر، راه حل اول - که از r^2 استفاده می شود - بر آزمون ضریب تعیین مبتنی است. در یک مدل رگرسیون ساده، آزمون فرضیه معنی دار بودن مدل رگرسیون چیزی نیست جز آزمون $\beta = 0$ ؛ زیرا تمام ساختار مدل رگرسیون ساده به β بستگی دارد. اگر $H_0: \beta = 0$ رد نشود بدین معنی است که کل مدل قابل قبول نیست؛ راه حل دوم نیز که دقیقاً همان فرضیه $H_0: \beta = 0$ است؛ و راه حل سوم که بر آزمون F مبتنی است، در حقیقت معنی دار بودن تغییرات توضیح داده شده را آزمون می کند. می دانیم در یک مدل رگرسیون ساده، معنی دار بودن تغییرات توضیح داده شده دقیقاً بر فرض معنی داری شیب مدل متکی است؛ به همین دلیل آزمون دوم و سوم نیز از نظر ماهیت یکی هستند.

۱۹. با استفاده از معادله ۲-۵۶ داریم

$$\begin{aligned} r^2 &= \frac{t^2}{t^2 + (n - 2)} \\ &= \frac{(10/603)^2}{(10/603)^2 + 3} = \frac{240/016}{248/016} \\ &\cong 0/988 . \end{aligned}$$

۲-۲ با توجه به معادله ۲-۲۵ می دانیم

$$\text{Cov}(\hat{\alpha}, \hat{\beta}) = \frac{-\bar{X} \sigma^2}{\sum x_i^2} .$$

بنابراین وقتی $\bar{X} = 0$ ، آنگاه $\text{Cov}(\hat{\alpha}, \hat{\beta}) = 0$ است. این نتیجه را می توان چنین توجیه کرد که می دانیم

$$\hat{\alpha} = \bar{Y} - \beta \bar{X} ,$$

وقتی $\bar{X} = 0$ در نتیجه $\hat{\alpha} = \bar{Y}$. حال می گوئیم که اگر به Y_i مقدار ثابت k را اضافه کنیم

یا از آن کم کنیم، مقدار $\hat{\beta}$ تغییری نمی‌کند؛ زیرا

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2} = \frac{\sum x_t Y_t}{\sum x_t^2} .$$

اگر به Y_t مقدار k را اضافه کنیم، خواهیم داشت

$$\begin{aligned} \hat{\beta} &= \frac{\sum x_t (Y_t + k)}{\sum x_t^2} = \frac{\sum x_t Y_t + k \sum x_t}{\sum x_t^2} , \\ &= \frac{\sum x_t y_t}{\sum x_t^2} . \end{aligned}$$

بنابراین می‌توان به راحتی نتیجه گرفت که $\hat{\beta}$ ، تابعی از \bar{Y} نیست. در نتیجه اگر $\bar{X} = 0$ باشد، از یک طرف $\hat{\alpha}$ دقیقاً برابر \bar{Y} می‌شود و از طرف دیگر $\hat{\beta}$ از \bar{Y} مستقل خواهد بود؛ بنابراین باید انتظار داشت که کوواریانس $\hat{\alpha}$ و $\hat{\beta}$ برابر صفر شود.

۲-۳ الف) از رابطه $Y_t \equiv C_t + S_t$ ، داریم $y_t \equiv c_t + s_t$ واریانس S_t را حساب می‌کنیم.

$$\begin{aligned} \text{Var}(S_t) &= \frac{\sum (S_t - \bar{s})^2}{n} , \\ &= \frac{\sum (s_t)^2}{n} = \frac{\sum (y_t - c_t)^2}{n} = \frac{\sum y_t^2}{n} + \frac{\sum c_t^2}{n} - \frac{2 \sum y_t c_t}{n} . \quad (1) \end{aligned}$$

دو طرف رابطه فوق را بر واریانس Y_t تقسیم می‌کنیم،

$$\begin{aligned} \frac{\text{Var}(S_t)}{\text{Var}(Y_t)} &= 1 + \frac{\text{Var}(C_t)}{\text{Var}(Y_t)} - \frac{2 \sum y_t c_t}{\sum y_t^2} , \\ &= 1 + \frac{\text{Var}(C_t)}{\text{Var}(Y_t)} - 2 \hat{\beta}_1 . \quad (2) \end{aligned}$$

در یک مدل رگرسیون $Y_t = \alpha + \beta X_t + U_t$ ، ضریب همبستگی بین X_t و Y_t

برابر است با

$$r = \frac{\sum x_t y_t}{\sqrt{\sum x_t^2 \sum y_t^2}}$$

عبارت فوق را می توان به صورت زیر نوشت،

$$r = \frac{\sum x_t y_t}{\sum x_t^2} \cdot \frac{\sqrt{\sum x_t^2}}{\sqrt{\sum y_t^2}}$$

کسر اول برابر $\hat{\beta}$ است و اگر صورت و مخرج کسر دوم را بر \sqrt{n} تقسیم کنیم، به انحراف معیار X و انحراف معیار Y خواهیم رسید؛ بنابراین

$$r = \hat{\beta} \frac{SE(X)}{SE(Y)}$$

یا:

$$\frac{SE(Y)}{SE(X)} = \frac{\hat{\beta}}{r} \quad (۳)$$

فرمول کلی فوق را با مدل رگرسیون $C_t = \alpha_1 + \beta_1 Y_t + U_{1t}$ تطبیق می دهیم،

$$\frac{SE(C)}{SE(Y)} = \frac{\hat{\beta}_1}{r} = \frac{0/6}{r} \quad (۴)$$

به همین ترتیب، با توجه به معادله (۳)، برای مدل رگرسیون $Y_t = \alpha_1 + \beta_1 C_t + U_{1t}$ خواهیم داشت

$$\frac{SE(Y)}{SE(C)} = \frac{\hat{\beta}_1}{r} = \frac{1/2}{r} \quad (۵)$$

دو طرف معادله (۴) را بر (۵) تقسیم می کنیم،

$$\frac{[SE(C)]^2}{[SE(Y)]^2} = \frac{\hat{\beta}_1}{\hat{\beta}_1} = \frac{0/6}{1/2} = 0/۵$$

یا

$$\frac{Var(C)}{Var(Y)} = 0/۵$$

اگر این رابطه را در معادله (۲) قرار دهیم، خواهیم داشت

$$\frac{\text{Var}(S_t)}{\text{Var}(Y_t)} = 1 + 0/0 - 2(0/6) = 0/3 ,$$

و در نتیجه

$$\frac{\text{SE}(S_t)}{\text{SE}(Y_t)} = \sqrt{0/3} . \quad (6)$$

ب) برای به دست آوردن ضریب همبستگی بین S_t و Y_t و نیز بین C_t و S_t به ترتیب زیر عمل می‌کنیم. ابتدا نسبت $\frac{\text{SE}(S_t)}{\text{SE}(C_t)}$ را به دست می‌آوریم. دو طرف معادله (۱) را بر $\text{Var}(C_t)$ تقسیم می‌کنیم؛

$$\begin{aligned} \frac{\text{Var}(S_t)}{\text{Var}(C_t)} &= \frac{\text{Var}(Y_t)}{\text{Var}(C_t)} + 1 - 2\hat{\beta}_1 , \\ &= 2 + 1 - 2(1/2) = 3 - 2/4 = 0/6 , \end{aligned}$$

در نتیجه

$$\frac{\text{SE}(S_t)}{\text{SE}(C_t)} = \sqrt{0/6} . \quad (7)$$

حال به محاسبه $r_{s,y}$ و $r_{s,c}$ می‌پردازیم. مطابق فرمول (۳) داریم

$$r_{s,y} = \hat{\beta}_{s,y} \frac{\text{SE}(Y_t)}{\text{SE}(S_t)} ,$$

$$r_{s,c} = \hat{\beta}_{s,c} \frac{\text{SE}(C_t)}{\text{SE}(S_t)} .$$

قبلاً در معادله‌های (۶) و (۷) نسبت انحراف معیارها را حساب کرده‌ایم. فقط باید $\hat{\beta}_{s,y}$ و $\hat{\beta}_{s,c}$ را به دست آورد،

$$\hat{\beta}_{s,y} = \frac{\sum s_t y_t}{\sum y_t^2} = \frac{\sum (y_t - c) y_t}{\sum y_t^2} ,$$

$$= 1 - \frac{\sum c_t y_t}{\sum y_t^2} = 1 - \hat{\beta}_{c,y} ,$$

$$= 1 - \hat{\beta}_r = 1 - 0/6 = 0/4 .$$

به همین ترتیب

$$\hat{\beta}_{s,c} = \frac{\sum s_t c_t}{\sum c_t^2} = \frac{\sum (y_t - c_t) c_t}{\sum c_t^2} ,$$

$$= \frac{\sum y_t c_t}{\sum c_t^2} - 1 = \hat{\beta}_{y,c} - 1 ,$$

$$= \hat{\beta}_r - 1 = 1/2 - 1 = 0/2 .$$

بدین ترتیب، خواهیم داشت

$$r_{s,y} = 0/4 \left(\frac{1}{\sqrt{0/3}} \right) = 0/73 ,$$

$$r_{s,c} = 0/2 \left(\frac{1}{\sqrt{0/6}} \right) = 0/26 ,$$

۲-۴ الف) تخمین مدل به صورت زیر خواهد بود،

$$\hat{Y}_t = 3/366 - 0/286 X_t ,$$

$$(0/221) \quad (0/063)$$

$$r^2 = 0/603 ,$$

$$RSS = \sum e_t^2 = 1/1430 .$$

ب) فاصله اطمینان ۹۵ درصد برای β ، با توجه به مقدار t با ۱۱ درجه آزادی، یعنی $t = \pm 2/201$ ، عبارت است از

$$-0/286 - 2/201 (0/063) < \beta < -0/286 + 2/201 (0/063) ,$$

$$-0.147 < \beta < -0.425 .$$

ج) با توجه به اینکه $\beta = 0$ در فاصله اطمینان ۹۵ درصد برای β واقع نمی شود، بنابراین فرضیه H_0 را در سطح معنی دار ۵ درصد رد می کنیم.

د) برای ساختن یک فاصله اطمینان ۹۰ درصدی برای σ^2 ، به این نکته توجه می کنیم که $\frac{\sum e_i^2}{\sigma^2}$ دارای یک توزیع χ^2 با ۱۱ درجه آزادی است. با استفاده از معادله ۲-۴۴ داریم

$$\frac{(n-2) \hat{\sigma}^2}{\chi^2_{0.95}} < \sigma^2 < \frac{(n-2) \hat{\sigma}^2}{\chi^2_{0.05}} ,$$

$$\frac{1/1435}{19/7} < \sigma^2 < \frac{1/1435}{4/57} ,$$

$$0.058 < \sigma^2 < 0.250 .$$

۲-۵ می دانیم واریانس $\hat{\beta}$ از فرمول زیر به دست می آید،

$$\text{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2} = \frac{\hat{\sigma}^2}{\sum (X_i - \bar{X})^2} .$$

ملاحظه می شود واریانس $\hat{\beta}$ با $\sum (X_i - \bar{X})^2$ نسبت معکوس دارد. به موازات افزایش واریانس یا پراکندگی متغیر برون زای X_i ، مقدار $\sum (X_i - \bar{X})^2$ زیاد می شود و واریانس $\hat{\beta}$ کمتر می شود. در نمونه ۸ تایی پژوهشگر دوم پراکندگی مشاهدات X_i بسیار زیاد است (۱ و ۱۰)؛ در حالی که این پراکندگی برای نمونه اول بسیار کم است (۳ و ۴). بنابراین باید انتظار داشت که واریانس متغیر X_i در نمونه دوم بسیار کمتر از مقدار مشابه برای نمونه اول باشد. در نتیجه واریانس $\hat{\beta}$ و نیز انحراف معیار آن برای نمونه دوم بسیار کمتر از انحراف معیار $\hat{\beta}$ برای نمونه اول خواهد بود.

۲-۶ این مسأله شبیه نکته ای است که در مسأله ۲-۵ مطرح شده است. می دانیم

$$\text{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{\sum x_i^2} = \frac{\hat{\sigma}^2}{\sum (X_i - \bar{X})^2} .$$

همچنین

$$\text{Var}(X_i) = \frac{\sum (X_i - \bar{X})^2}{n} ,$$

یا

$$\sum (X_i - \bar{X})^2 = n \text{Var}(X_i) .$$

با جایگزینی در فرمول $\text{Var}(\hat{\beta})$ خواهیم داشت

$$\text{Var}(\hat{\beta}) = \frac{\hat{\sigma}^2}{n \text{Var}(X_i)} .$$

با حذف مشاهدات مرکزی X_i ، با اینکه واریانس X_i زیاد و $\text{Var}(\hat{\beta})$ کم می‌شود، اما اولاً، با کم شدن حجم نمونه و کاهش n ، مقدار $\text{Var}(\hat{\beta})$ زیاد می‌شود. که با هدف اصلی مغایر است.

ثانیاً، با کم شدن n ، ممکن است مقدار $\hat{\sigma}^2$ نیز افزایش یابد؛ زیرا

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-2} ,$$

در نتیجه صورت کسر $\text{Var}(\hat{\beta})$ زیاد شده و مقدار $\text{Var}(\hat{\beta})$ افزایش نشان خواهد داد که این نیز خلاف هدف اولیه است. البته با کاهش n ، باید به جهت حرکت e_i^2 نیز توجه داشت؛ اما معمولاً به ازای کاهش n انتظار می‌رود که واریانس U_i زیاد شود.

ثالثاً، با حذف مشاهدات مرکزی X_i ، نمونه ما، دیگر یک نمونه تصادفی نخواهد بود و این امر زیربنای تمام محاسبات ما را متزلزل می‌کند؛ زیرا می‌دانیم از فرضهای اولیه تخمین‌زنده‌های روش حداقل مربعات معمولی تصادفی بودن نمونه مشاهدات X_i و Y_i است.

۲.۷ الف) می‌دانیم $e_i = Y_i - \hat{\beta} X_i$. به جای Y_i و $\hat{\beta}$ مقادیر آن را قرار می‌دهیم.

$$\begin{aligned}
 e_i &= \beta X_i + U_i - \frac{X_i \sum X_j Y_j}{\sum X_j^2} \\
 &= U_i + \beta X_i - \frac{X_i \sum (\beta X_j + U_j) X_j}{\sum X_j^2} \\
 &= U_i - \frac{X_i \sum U_j X_j}{\sum X_j^2} \quad (1)
 \end{aligned}$$

$$\begin{aligned}
 E(e_i) &= E \left[U_i + \frac{X_i (\sum U_j X_j)^2}{(\sum X_j^2)^2} - \frac{X_i U_i \sum U_j X_j}{\sum X_j^2} \right] \\
 &= \sigma^2 - \frac{\sigma^2 X_i^2}{\sum X_j^2} = \sigma^2 \left(1 - \frac{X_i^2}{\sum X_j^2} \right) \quad (2)
 \end{aligned}$$

ب) می‌دانیم، $e_i = U_i - \frac{X_i \sum U_j X_j}{\sum X_j^2}$ به همین ترتیب خواهیم داشت

$$e_i = U_i - \frac{X_i \sum U_j X_j}{\sum X_j^2}$$

در نتیجه داریم

$$E(e_i e_j) = E \left[U_i U_j - \frac{X_i U_i \sum X_k U_k}{\sum X_k^2} - \frac{X_j U_j \sum X_k U_k}{\sum X_k^2} + \frac{X_i X_j (\sum U_k X_k)^2}{(\sum X_k^2)^2} \right]$$

می‌دانیم بنا بر فرضهای کلاسیک، $E(U_i U_j) = 0$. امید ریاضی سه جمله دیگر را باید حساب کنیم؛ برای مثال، امید ریاضی یکی از آنها را نشان می‌دهیم،

$$\begin{aligned}
 E \left[- \frac{X_i U_i \sum X_k U_k}{\sum X_k^2} \right] &= - \frac{X_i E[U_i \sum X_k U_k]}{\sum X_k^2} \\
 &= \frac{- X_i E[U_i (X_1 U_1 + X_2 U_2 + \dots + X_i U_i + \dots + X_n U_n)]}{\sum X_k^2}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{-X_t [X_1 E(U_1 U_t) + X_2 E(U_2 U_t) + \dots + X_s E(U_s U_t) + \dots + X_s E(U_s U_t)]}{\sum X_t^2} \\
 &= \frac{-X_t [0 + 0 + \dots + X_s \sigma^2 + 0 + \dots + 0]}{\sum X_t^2} \\
 &= \frac{-X_t X_s \sigma^2}{\sum X_t^2}
 \end{aligned}$$

به همین ترتیب برای دو جمله دیگر عمل می‌کنیم. خواهیم داشت

$$E(e_t e_s) = 0 - \frac{X_t X_s \sigma^2}{\sum X_t^2} - \frac{X_t X_s \sigma^2}{\sum X_t^2} + \frac{X_t X_s \sigma^2}{\sum X_t^2}$$

در نتیجه

$$E(e_t e_s) = -\frac{\sigma^2 X_t X_s}{\sum X_t^2}, \quad t \neq s.$$

(ج) برای نشان دادن $\sum e_t \neq 0$ می‌نویسیم

$$\begin{aligned}
 \sum e_t &= \sum (Y_t - \hat{\beta} X_t) \\
 &= \sum \left[Y_t - \frac{X_t \sum X_t Y_t}{\sum X_t^2} \right].
 \end{aligned}$$

هیچ دلیلی وجود ندارد که رابطه فوق در حالت کلی برابر صفر باشد. اما در مدل‌هایی که ضریب ثابت دارند مقدار $\sum e_t$ برابر است با

$$\sum e_t = \sum (Y_t - \hat{\alpha} - \hat{\beta} X_t).$$

طرف راست رابطه فوق در واقع معادله اول نرمال است که حاصل مشتق‌گیری مجموع مربعات پسماند نسبت به $\hat{\alpha}$ است و می‌دانیم که باید صفر باشد تا $\hat{\alpha}$ بتواند مجموع مربعات پسماند را حداقل کند. بنابراین، چون چنین شرطی در مدل‌های رگرسیون بدون ضریب ثابت موجود نیست، در حالت کلی هیچ دلیلی بر صفر بودن $\sum e_t$ در این گونه

مدلها وجود ندارد.

البته می توان نشان داد که $E(\sum e_i)$ در مدل‌هایی که ضریب ثابت ندارند همواره برابر صفر است. از دو طرف رابطه (۱)، \sum می گیریم،

$$\sum e_i = \sum \left[U_i - \frac{X_i \sum X_i U_i}{\sum X_i^2} \right],$$

و با گرفتن امید ریاضی خواهیم داشت

$$E[\sum e_i] = \sum \left[E(U) - \frac{X_i \sum E(X_i U_i)}{\sum X_i^2} \right],$$

= ۰ .

د) رابطه (۲) را یک بار دیگر می نویسیم،

$$E(e_i^2) = \sigma^2 - \frac{\sigma^2 X_i^2}{\sum X_i^2}.$$

بنابراین

$$\begin{aligned} E\left(\frac{\sum e_i^2}{n-1}\right) &= \sum \left(\sigma^2 - \frac{\sigma^2 X_i^2}{\sum X_i^2}\right) / (n-1), \\ &= \frac{n\sigma^2 - \sigma^2}{n-1} = \sigma^2. \end{aligned}$$

۲.۸ الف) چون $\tilde{\beta}$ یک تخمین زنده خطی است، داریم

$$\tilde{\beta} = w_1 Y_1 + w_2 Y_2 + \dots + w_n Y_n,$$

$$\tilde{\beta} = \sum w_i Y_i. \quad (1)$$

در اینجا باید شرایطی را پیدا کنیم که $\tilde{\beta}$ نااریب شود. از رابطه (۱) امید ریاضی می گیریم.

$$\begin{aligned} E(\tilde{\beta}) &= E[\sum w_i Y_i], \\ &= E[\sum w_i (\alpha + \beta X_i + U_i)], \end{aligned}$$

$$= \alpha \sum w_i + \beta \sum w_i X_i + 0 .$$

برای $E(\tilde{\beta}) = \beta$ ، ضروری است شرایط زیر صادق باشد،

$$\sum w_i = 0 . \tag{۲}$$

$$\sum w_i X_i = 1 .$$

اگر متغیری مانند Z_i داشته باشیم، به گونه‌ای که $z_i = Z_i - \bar{Z}$ ، سپس w_i را به صورت زیر تعریف کنیم

$$w_i = \frac{z_i}{\sum z_i X_i} .$$

آنگاه خواهیم داشت

$$w_i = \frac{z_i}{\sum z_i X_i} . \tag{۳}$$

واضح است که $\sum w_i = 0$ و $\sum w_i X_i = 1$. بنابراین w_i که در رابطه (۳) تعریف شده است شرایط (۲) را تأمین می‌کند. اگر w_i را در رابطه (۱) قرار دهیم، تخمین زنده‌های خطی ناریب به دست می‌آیند،

$$\tilde{\beta} = \frac{\sum z_i Y_i}{\sum z_i X_i} , \tag{۴}$$

یا

$$\tilde{\beta} = \frac{\sum z_i y_i}{\sum z_i x_i} . \tag{۵}$$

ب) از دو طرف رابطه (۱) واریانس می‌گیریم،

$$\text{Var}(\tilde{\beta}) = \text{Var}(w_1 Y_1 + w_2 Y_2 + \dots + w_n Y_n) .$$

با توجه به اینکه $\text{Var}(Y_i) = \text{Var}(U_i) = \sigma^2$ و بنا بر فرض $\text{Cov}(U_i, U_j) = 0$ ،

بنابراین:

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= w_1^t \sigma^t + w_2^t \sigma^t + \dots + w_n^t \sigma^t, \\ &= \sigma^t \sum w_i^t. \end{aligned}$$

با جایگزینی (۳) در رابطه فوق خواهیم داشت

$$\text{Var}(\tilde{\beta}) = \frac{\sigma^t \sum z_i^t}{(\sum z_i x_i)^t}. \quad (6)$$

می‌دانیم واریانس $\hat{\beta}_{OLS}$ برابر است با

$$\text{Var}(\hat{\beta}) = \frac{\sigma^t}{\sum x_i^t}. \quad (7)$$

ج) برای به دست آوردن کارآیی نسبی، کافی است رابطه (۷) را بر (۶) تقسیم کنیم،

$$\frac{\text{Var}(\hat{\beta})}{\text{Var}(\tilde{\beta})} = \frac{(\sum z_i x_i)^t}{\sum x_i^t \cdot \sum z_i^t},$$

که با توجه به نامساوی، کوچکی - شوارتزیین صفر و یک است. ملاحظه می‌شود به ازای $z_i = k x_i$ داریم

$$\frac{\text{Var}(\hat{\beta})}{\text{Var}(\tilde{\beta})} = \frac{(\sum k x_i^t)^t}{\sum x_i^t \cdot k^t \sum x_i^t} = \frac{k^t (\sum x_i^t)^t}{k^t (\sum x_i^t)^t} = 1.$$

پیش‌بینی و مباحث تکمیلی در مدل رگرسیون خطی ساده

۳-۱ مقدمه

در فصلهای اول و دوم دیدیم که با تخمین یک مدل رگرسیون می‌توان به کمک تغییرات متغیر برون‌زا، تغییرات متغیر درون‌زا را توضیح داد. یکی از مهمترین موارد کاربرد تخمین مدل‌های رگرسیون این است که بتوانیم مقادیر آینده متغیر درون‌زا را پیش‌بینی کنیم. همان‌گونه که خواهیم دید، مقدار پیش‌بینی شده متغیر درون‌زا، یک متغیر تصادفی است و باید تابع توزیع احتمال با میانگین و واریانس آن دقیقاً مطالعه شود. با داشتن تابع توزیع احتمال و خصوصیات آن می‌توان فرضیه‌های مختلف را در مورد مقادیر پیش‌بینی شده متغیر درون‌زا به راحتی آزمون کرد. این نکات موضوع قسمت ۳-۲ خواهد بود.

هدف ما، در فصلهای اول و دوم، ارائه مفاهیم کلیدی در تخمین یک مدل رگرسیون خطی ساده و آزمون فرضیه‌های مختلف در مورد پارامترهای آن بوده است، بنابراین بسیاری از مباحث جانبی، ناگفته ماند. طرح این نکته‌های حاشیه‌ای در متن مباحث اولیه، ممکن بود توجه خواننده را از مسیر اصلی ارائه مطالب منحرف سازد. به نظر رسید که بهتر است این موضوعها در فصلی مجزا مطرح شود. در این فصل بعد از پایان مبحث پیش‌بینی، به ذکر چند عنوان در این مورد می‌پردازیم.

توجه به مسأله مقیاس اندازه‌گیری متغیرهای برون‌زا و درون‌زا اهمیت فراوان دارد. سؤال این است: اگر مقیاس اندازه‌گیری یک متغیر را تغییر دهیم، آیا تخمین پارامترهای یک مدل رگرسیون نیز تغییر خواهد کرد؟ این نکته‌ای است که در قسمت ۳-۳ به آن می‌پردازیم.

مطالعات اجمالی در رگرسیون معکوس، موضوع قسمت ۳-۴ خواهد بود. تا اینجا

ما، در مدل‌های رگرسیون، Y_t را تابعی از تغییرات X_t فرض می‌کردیم. آیا می‌توان جهت تأثیرپذیری متغیرها را عوض کرده، این بار X_t را تابعی از تغییرات Y_t بگیریم؟ اگر جواب، مثبت باشد، سؤال این است که چه رابطه‌ای بین تخمین پارامترهای دو مدل وجود دارد.

در مباحث دو فصل گذشته، بویژه در بعضی از مسائل پایانی هر فصل، احتمال دارد خواننده با این نکته مواجه شده باشد که بعضی از مقادیر «دور افتاده» در متغیرهای برون‌زا و درون‌زا می‌تواند بر تخمین پارامترهای یک مدل رگرسیون تأثیر بسیار قابل ملاحظه‌ای داشته باشد. آیا می‌توان این متغیرهای دور افتاده را حذف کرد؟ تأثیر این متغیرها دقیقاً به چه صورت است؟ این نکته‌ها را به طور مختصر در قسمت ۳-۵ بررسی خواهیم کرد.

مدل رگرسیون ساده خطی که در فصل‌های گذشته مطالعه شد در واقع مدلی است که هم برحسب متغیرها و هم برحسب پارامترها خطی است. در بعضی از مدل‌های رگرسیون، متغیرها به صورت غیرخطی وارد می‌شوند؛ و در بعضی دیگر، پارامترها صورت غیرخطی دارند. آیا می‌توان با تبدیلهای مناسب، مدل‌های غیرخطی برحسب متغیرها و مدل‌های غیرخطی برحسب پارامترها را به مدل‌های خطی تبدیل کرد؟ این سؤالی است که در قسمت ۳-۶ مطرح خواهد شد.

۳-۲ پیش‌بینی در مدل‌های رگرسیون ساده

یکی از هدفهای اساسی در تخمین یک مدل رگرسیون، این است که بتوان تغییرات متغیر درون‌زا را به ازای مقدار معینی از متغیر برون‌زا پیش‌بینی کرد. مدل رگرسیون ساده خطی زیر را ملاحظه کنید.

$$Y_t = \alpha + \beta X_t + U_t \quad (3-1)$$

فرض کنید می‌خواهیم برای سال f در آینده پیش‌بینی کنیم. ابتدا باید مقدار متغیر برون‌زا

در سال t را بدانیم. مقداری را که متغیر X در سال t خواهد داشت، با X_t نشان می‌دهیم. پیش‌بینی مقدار متغیر درون‌زا را به ازای X_t به دو صورت می‌تواند انجام شود: پیش‌بینی نقطه‌ای^۱ و پیش‌بینی فاصله‌ای^۲.

پیش‌بینی نقطه‌ای در واقع تخمین یک نقطه از فضای Y در آینده است که به ازای مقدار معین X_t صورت می‌پذیرد. اگر دوره‌ی زمانی مشاهدات را سالانه فرض کنیم، مدل رگرسیون مفروض برای سال t عبارت است از

$$Y_t = \alpha + \beta X_t + U_t \quad (3-2)$$

ملاحظه می‌شود فرض بر عدم تغییر ساختاری در مدل رگرسیون است؛ به عبارت دیگر، پارامترها و شکل ریاضی مدل رگرسیون مفروض - در آینده‌ای که می‌خواهیم برای آن پیش‌بینی کنیم - ثابت مانده است. منظور از پیش‌بینی نقطه‌ای محاسبه مقدار Y_t به ازای X_t است.

تخمین رگرسیون (۳-۱) عبارت است از

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t \quad (3-3)$$

مدل ۳-۳ را برای سال t در آینده می‌نویسیم،

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t \quad (3-4)$$

به ازای X_t ، مقدار \hat{Y}_t به منزله پیش‌بینی نقطه‌ای از متغیر درون‌زا به راحتی از معادله ۳-۴ به دست می‌آید. این پیش‌بینی قاعدتاً یک پیش‌بینی شرطی است، بدین معنی که تحقق \hat{Y}_t اولاً، مشروط به وقوع X_t است، و ثانیاً، فرض بر این است که ترکیب و ساختار رابطه بین X و Y در آینده، (تا دوره‌ای که می‌خواهیم پیش‌بینی کنیم) کوچکترین تغییری نکرده است و در واقع مقادیر α و β و شکل ریاضی مدل مفروض، ثابت خواهد بود.

۱. پیش‌بینی فاصله‌ای برای Y_f

فرض کنید پیش‌بینی نقطه‌ای (\hat{Y}_f) را برای یک دوره زمانی در آینده به دست آورده‌ایم. بدیهی است در سال f ، متغیر درون‌زا مقدار معینی خواهد داشت که آن را با Y_f نشان می‌دهیم. اینکه Y_f چه مقداری دارد، تابعی است از X_f ، مقادیر α و β و مهمتر از همه مقداری که جمله اختلال در سال f ، یعنی U_f ، خواهد داشت. بجز X_f ، از عوامل دیگر هیچ اطلاعاتی نداریم. یعنی α و β و U_f برای ما مجهول هستند. تنها چیزی که از Y_f می‌دانیم پیش‌بینی نقطه‌ای آن، یعنی \hat{Y}_f ، است. سؤال این است که آیا می‌توان با داشتن \hat{Y}_f برای مقدار واقعی Y_f یک فاصله اطمینان ساخت؟ به عبارت دیگر، آیا می‌توان به کمک \hat{Y}_f مثلاً در سطح احتمال ۹۵ درصد گفت که Y_f بین دو مقدار معلوم a و b قرار خواهد گرفت؟ پیش‌بینی فاصله‌ای دقیقاً همین مسأله است؛ یعنی ساختن فاصله‌ای که با احتمال معینی بتواند مقدار متغیر درون‌زا را در آینده مورد نظر شامل شود.

برای ساختن این فاصله اطمینان، ابتدا می‌گوییم \hat{Y}_f یک متغیر تصادفی است؛ زیرا تابعی از $\hat{\alpha}$ و $\hat{\beta}$ است که هر دو تصادفی هستند. در سال f ، مقداری که Y_f می‌گیرد نیز یک متغیر تصادفی خواهد بود؛ زیرا تابعی از U_f تصادفی است. اما قاعدتاً انتظار این است که مقدار پیش‌بینی ما از متغیر درون‌زا در آینده (\hat{Y}_f) دقیقاً با مقداری که واقعاً این متغیر در آینده خواهد داشت (Y_f) برابر نباشد. تفاوت یا انحراف پیش‌بینی نقطه‌ای \hat{Y}_f از مقدار واقعی متغیر درون‌زا (Y_f) را اصطلاحاً خطای پیش‌بینی^۱ می‌نامیم و آن را با e_f نشان می‌دهیم؛

$$e_f = Y_f - \hat{Y}_f \quad (۳-۵)$$

مقادیر Y_f و \hat{Y}_f را از معادله‌های ۳-۲ و ۳-۳ در معادله ۳-۵ جایگزین می‌کنیم؛

$$e_f = (\alpha + \beta X_f + U_f) - (\hat{\alpha} + \hat{\beta}) X_f \quad \text{یا}$$

$$e_f = U_f - (\hat{\alpha} - \alpha) - (\hat{\beta} - \beta) X_f \quad (۳-۶)$$

در سمت راست معادله ۳-۶ متغیرهای $\hat{\alpha}$ ، $\hat{\beta}$ و U_i متغیرهای تصادفی هستند، بنابراین e_i نیز یک متغیر تصادفی خواهد بود. می‌دانیم توابع توزیع احتمال $\hat{\alpha}$ ، $\hat{\beta}$ و U_i نرمال هستند؛ در نتیجه e_i نیز - که یک تابع خطی از آنهاست - توزیع احتمال نرمال خواهد داشت. در اینجا باید میانگین و واریانس خطای پیش‌بینی را به دست آوریم. از معادله ۳-۶ امید ریاضی می‌گیریم،

$$E(e_i) = E(U_i) - E(\hat{\alpha} - \alpha) - X_i E(\hat{\beta} - \beta) .$$

با توجه به نااریب بودن $\hat{\alpha}$ ، $\hat{\beta}$ و فرض صفر بودن امید ریاضی U_i خواهیم داشت

$$E(e_i) = 0 . \quad (3-7)$$

بنابراین نخستین نتیجه‌ای که به دست آوردیم این است که میانگین خطای پیش‌بینی، صفر است. قبل از وارد شدن به بحث واریانس e_i ، به یک اصطلاح، اشاره می‌کنیم. می‌دانیم $\hat{\alpha}$ یا $\hat{\beta}$ را تخمین زنده‌های روش حداقل مربعات معمولی از α ، β می‌گوییم. به همین ترتیب می‌توان \hat{Y}_i را پیش‌بینی کننده^۱ روش حداقل مربعات معمولی از Y_i نامید.

برای به دست آوردن واریانس e_i ، از تعریف آن شروع می‌کنیم،

$$\text{Var}(e_i) = E[e_i - E(e_i)]^2 .$$

با استفاده از معادله ۳-۷ داریم

$$\text{Var}(e_i) = E(e_i)^2 .$$

دو طرف معادله ۳-۶ را مجذور کرده خواهیم داشت

$$e_i^2 = U_i^2 + (\hat{\alpha} - \alpha)^2 + (\hat{\beta} - \beta)^2 X_i^2 - 2 U_i (\hat{\alpha} - \alpha) \\ - 2 U_i (\hat{\beta} - \beta) X_i + 2 (\hat{\alpha} - \alpha) (\hat{\beta} - \beta) X_i .$$

از معادله فوق امید ریاضی می‌گیریم،

$$\begin{aligned} \text{Var}(e_p) &= \text{Var}(U_p) + \text{Var}(\hat{\alpha}) + X_p' \text{Var}(\hat{\beta}) \\ &\quad - 2 \text{Cov}[U_p, (\hat{\alpha} - \alpha)] - 2 \text{Cov}[U_p, (\hat{\beta} - \beta) X_p] \\ &\quad + 2 X_p \text{Cov}(\hat{\alpha}, \hat{\beta}) . \end{aligned}$$

بنا بر فرض واریانس همسانی، می‌دانیم: $\text{Var}(U_p) = \text{Var}(U_i) = \sigma^2$. می‌توان مقادیر واریانس $\hat{\alpha}$ و واریانس $\hat{\beta}$ و کوواریانس $\hat{\alpha}$ و $\hat{\beta}$ را به ترتیب از فرمولهای ۲-۱۶، ۲-۲۱ و ۲-۲۵ به دست آورد. می‌توان درباره $\text{Cov}[U_p, (\hat{\beta} - \beta) X_p]$ ، توجه را فقط به کوواریانس U_p و $\hat{\beta}$ محدود کرد، زیرا β و X_p مقادیر ثابتی هستند. از معادله ۲-۱۰ می‌دانیم که $\hat{\beta}$ ، یک تابع خطی از مقادیر U_1, U_2, \dots, U_n است و چون بنا بر فرض عدم خودهمبستگی، این مقادیر از U_p مستقل هستند؛ بنابراین $\hat{\beta}$ نیز از U_p مستقل بوده و کوواریانس آنها صفر می‌شود:

$$\text{Cov}[U_p, (\hat{\beta} - \beta) X_p] = 0 .$$

به همین ترتیب می‌توان نشان داد که

$$\text{Cov}[U_p, (\hat{\alpha} - \alpha)] = 0 .$$

بنابراین خواهیم داشت

$$\begin{aligned} \text{Var}(e_p) &= \sigma^2 + \sigma^2 \left[\frac{1}{n} + \frac{\bar{X}'^2}{\sum x_i'^2} \right] + \frac{X_p' \sigma^2}{\sum x_i'^2} - \frac{2 X_p \bar{X} \sigma^2}{\sum x_i'^2} , \\ &= \sigma^2 \left[1 + \frac{1}{n} + \frac{\bar{X}'^2}{\sum x_i'^2} + \frac{X_p'^2}{\sum x_i'^2} - \frac{2 X_p \bar{X}}{\sum x_i'^2} \right] , \end{aligned}$$

در نتیجه

$$\text{Var}(e_p) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum x_i'^2} \right] . \quad (۳-۸)$$

با ملاحظه معادله ۳-۶ دیدیم که e_f توزیع نرمال دارد. معادله‌های ۳-۷ و ۳-۸ به ترتیب میانگین و واریانس خطای پیش‌بینی را مشخص می‌کند. بدین ترتیب به راحتی می‌توان e_f را استاندارد کرد. کافی است e_f را از میانگین آن کم کرده و بر انحراف معیارش تقسیم کنیم، آنگاه توزیع Z با میانگین صفر و واریانس یک خواهد داشت،

$$\frac{e_f - E(e_f)}{\sqrt{\text{Var}(e_f)}} \sim Z(0, 1),$$

$$\frac{e_f}{\sigma_u \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}}} \sim Z(0, 1).$$

اگر به جای σ_u مقدار تخمین آن را قرار دهیم، توزیع کسر فوق از Z به t تبدیل می‌شود:

$$\frac{e_f}{\hat{\sigma}_u \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}}} \sim t(0, 1), \quad (3-9)$$

که در آن $\hat{\sigma}_u = \sqrt{\frac{\sum e_i^2}{(n-2)}}$

سوالی که در مقدمه مطرح کردیم، این بود که آیا می‌توان با داشتن پیش‌بینی نقطه‌ای (\hat{Y}_f) به فاصله اطمینانی برای مقدار واقعی Y_f رسید؟ تاکنون بحث ما بیشتر روی خطای پیش‌بینی و محاسبه میانگین و واریانس آن متمرکز بوده است. با جایگزینی معادله ۳-۵ در معادله ۳-۹ به راحتی می‌توان به چنین فاصله اطمینانی دست یافت،

$$\frac{Y_f - \hat{Y}_f}{\hat{\sigma}_u \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}}} \sim t.$$

اگر در سطح معنی‌دار α ، مقدار t را با $(n-2)$ درجه آزادی از جدول t به دست آوریم ($\pm t_{\alpha/2}$)، دقیقاً می‌توان به روش گذشته فاصله اطمینان برای Y_f را به دست آورد. یادآوری می‌کنیم که در رابطه فوق، بجز Y_f ، تمام مقادیر دیگر را می‌دانیم. بدین ترتیب

فاصله اطمینان مطلوب عبارت است از

$$\hat{Y}_f - t_{\alpha/2} \hat{\sigma}_U \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}} < Y_f < \hat{Y}_f + t_{\alpha/2} \hat{\sigma}_U \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}}, \quad (3-10)$$

یا به صورت ساده‌تر

$$\hat{Y}_f - t_{\alpha/2} SE(e_f) < Y_f < \hat{Y}_f + t_{\alpha/2} SE(e_f). \quad (3-11)$$

مثال ۳-۱ با استفاده از یک نمونه دوازده‌تایی، تابع مصرف را به صورت زیر تخمین زده‌ایم،

$$\hat{C}_t = 10 + 0/90 Y_t,$$

که در آن C_t و Y_t به ترتیب مصرف و درآمد است (محاسبات کاملاً فرضی است). همچنین این کمیتها را داریم،

$$\hat{\sigma}^2 = 0/01, \quad \bar{Y} = 200, \quad \sum (Y_t - \bar{Y}) = 4000.$$

سطح درآمد در پنج سال دیگر $X_f = 250$ فرض می‌شود.

اولاً، پیش‌بینی نقطه‌ای از سطح مصرف در آن سال (\hat{C}_f) را به دست آورید.

ثانیاً، یک فاصله اطمینان ۹۵ درصد برای مقدار واقعی مصرف در سال f

پیش‌بینی کنید.

بهتر است علامت‌گذاربهای استفاده شده در مثال را با علامتهای به کار رفته در

متن کتاب، هماهنگ کنیم،

$$\hat{Y}_t = 10 + 0/90 X_t,$$

$$\hat{\sigma}^2 = 0/01, \quad \bar{X} = 200, \quad \sum (X_t - \bar{X}) = 4000, \quad Y_f = 250.$$

اولاً،

$$\hat{Y}_f = \hat{\alpha} + \hat{\beta} X_f,$$

$$= 10 + 0/90 (250) = 235,$$

یعنی مصرف در پنج سال دیگر به ازای $Y = 250$ ، برابر ۲۳۵ خواهد بود. ثانیاً، برای پیش‌بینی فاصله اطمینان ۹۵ درصد برای مصرف در سال t ، ابتدا واریانس خطای پیش‌بینی را به دست می‌آوریم. از معادله ۳-۸ داریم

$$\begin{aligned} \text{Var}(e_t) &= \hat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2} \right], \\ &= 0.01 \left[1 + \frac{1}{12} + \frac{(250 - 200)^2}{4000} \right] = 0.017. \end{aligned}$$

برای رسیدن به انحراف معیار خطای پیش‌بینی از واریانس e_t جذر می‌گیریم،

$$\text{SE}(e_t) = \sqrt{0.017} = 0.131.$$

برای محاسبه فاصله اطمینان ۹۵ درصد باید مقدار t را با $10 = 12 - 2$ درجه آزادی و $\frac{\alpha}{2} = 0.025$ از جدول t به دست آورد. خواهیم داشت $t = \pm 2/228$. با استفاده از نامساوی ۳-۱۱ فاصله اطمینان مورد نظر براحتی قابل محاسبه است،

$$235 - 2/228(0.131) < Y_t < 235 + 2/228(0.131),$$

$$235 - 0.29 < Y_t < 235 + 0.29,$$

یا

$$234.71 < Y_t < 235.29.$$

در پایان این مثال به دو اصطلاح اشاره می‌کنیم. نمونه مورد استفاده در این مثال شامل دوازده مشاهده است؛ بنابراین محدوده تغییرات متغیر برون‌زا (X_t) کاملاً مشخص است. اگر X_t ، یعنی مقداری از X - که می‌خواهیم به ازای آن پیش‌بینی کنیم - در محدوده تغییرات X_t در نمونه قرار گیرد، آنگاه می‌گوییم «پیش‌بینی درون - نمونه‌ای»^۱ داریم. در غیر این صورت پیش‌بینی ما «برون - نمونه‌ای»^۲ خواهد بود؛ بنابراین اصطلاح پیش‌بینی در اقتصادسنجی، همواره ناظر به مقادیر آینده متغیرها نیست.

1. Within-Sample Prediction

2. Out-of-Sample Prediction

۲. پیش‌بینی فاصله‌ای برای $E(Y_t)$

در قسمت قبل، برای Y_t (یک نقطه در آینده)، فاصله اطمینان پیش‌بینی کردیم. در بسیاری از مواقع، آنچه مورد توجه ماست، میانگین یا امید ریاضی Y_t است؛ بنابراین می‌خواهیم برای میانگین یا حد متوسط Y_t یک فاصله اطمینان بسازیم. از تعریف میانگین یا امید ریاضی Y_t آغاز می‌کنیم. از معادله ۳-۲ امید ریاضی می‌گیریم،

$$E(Y_t) = \alpha + \beta X_t . \quad (3-12)$$

برای تخمین مقدار $E(Y_t)$ به ازای X_t ، کافی است به جای α و β مقادیر $\hat{\alpha}$ و $\hat{\beta}$ را قرار دهیم،

$$\hat{E}(Y_t) = \hat{\alpha} + \hat{\beta} X_t .$$

اما در معادله ۳-۴، دیدیم $\hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t$ در نتیجه

$$\hat{E}(Y_t) = \hat{Y}_t . \quad (3-13)$$

بنابراین فرقی نمی‌کند که Y_t یا $E(Y_t)$ را پیش‌بینی کنیم، زیرا در هر دو حالت به یک جواب می‌رسیم. با اینکه پیش‌بینیهای نقطه‌ای برای Y_t و میانگین Y_t یکی هستند، فاصله‌های اطمینان این دو کاملاً از یکدیگر متفاوت است؛ زیرا همان‌گونه که خواهیم دید خطای پیش‌بینی و واریانس خطای پیش‌بینی آنها متساوی نیست.

اگر خطای پیش‌بینی را در این مورد با e_t^m نشان دهیم - که منعکس‌کننده خطای پیش‌بینی برای میانگین Y_t است - خواهیم داشت

$$e_t^m = E(Y_t) - \hat{E}(Y_t) ,$$

یا

$$e_t^m = E(Y_t) - \hat{Y}_t . \quad (3-14)$$

با جایگزینی معادله‌های ۳-۱۲ و ۳-۴ در معادله (۳-۱۴) داریم

$$e_t^m = \alpha + \beta X_t - (\hat{\alpha} + \hat{\beta} X_t) ,$$

در نتیجه

$$e_i^m = -(\hat{\alpha} - \alpha) - (\hat{\beta} - \beta) X_i \quad (3-15)$$

با مقایسه خطای پیش‌بینی در این حالت، یعنی معادله ۳-۱۵، با خطای پیش‌بینی در حالت اول، یعنی معادله ۳-۶، نتیجه می‌گیریم که خطاهای پیش‌بینی از یکدیگر متفاوت است.

در اینجا باید میانگین و واریانس خطای پیش‌بینی e_i^m را به دست آورد. از معادله

۳-۱۵ امید ریاضی می‌گیریم. خواهیم داشت

$$E(e_i^m) = 0 \quad (3-16)$$

برای محاسبه واریانس e_i^m می‌گوییم که

$$\text{Var}(e_i^m) = E[e_i^m - E(e_i^m)]^2 = E(e_i^m)^2$$

باید دو طرف معادله ۳-۱۵ را مجذور کنیم و امید ریاضی بگیریم،

$$\text{Var}(e_i^m) = E(e_i^2) = E[-(\hat{\alpha} - \alpha) - (\hat{\beta} - \beta) X_i]^2$$

اگر به ترتیبی کاملاً مشابه با آنچه در مورد واریانس e_i انجام دادیم، عمل کنیم، خواهیم داشت

$$\text{Var}(e_i^m) = \sigma^2 \left[\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum x_i^2} \right] \quad (3-17)$$

می‌توان نشان داد که e_i^m همانند e_i توزیعی نرمال با میانگین صفر و واریانس معلوم دارد که می‌توان آن را از معادله ۳-۱۷ به دست آورد. اگر به جای σ^2 مقدار تخمین آن ($\hat{\sigma}^2$) را قرار دهیم و e_i^m را استاندارد کنیم، توزیع t با $(n-2)$ درجه آزادی خواهد داشت،

$$\frac{e_i^m - 0}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum x_i^2}}} \sim t(n-2)$$

یا استفاده از معادله ۳-۱۴ داریم

$$\frac{E(Y_t) - \hat{Y}_t}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}}} \sim t(n-2)$$

یا به دست آوردن مقدار t از جدول در سطح معنی‌دار α درصد، فاصله اطمینان $(1 - \alpha)$ درصد برای $E(Y_t)$ به دست می‌آید،

$$\hat{Y}_t - t_{\alpha/2} \hat{\sigma}_0 \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}} < E(Y_t) < \hat{Y}_t + t_{\alpha/2} \hat{\sigma}_0 \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}}, \quad (3-18)$$

یا به صورت ساده‌تر

$$\hat{Y}_t - t_{\alpha/2} SE(e_t) < E(Y_t) < \hat{Y}_t + t_{\alpha/2} SE(e_t) \quad (3-19)$$

مثال ۳-۲ مدل رگرسیون زیر برای پیش‌بینی سطح فروش یک شرکت تجاری مفروض است،

$$Y_t = \alpha + \beta X_t + U_t$$

که در آن Y_t درآمد حاصل از فروش و X_t هزینه تبلیغات است. جدول ۳-۱ هزینه تبلیغات و درآمد حاصل از فروش ماهانه ۵ ماه را نشان می‌دهد.

جدول ۳-۱

t	درآمد حاصل از فروش Y_t (به صد هزار تومان)	هزینه تبلیغات X_t (به ده هزار تومان)
۱	۳	۱
۲	۴	۲
۳	۲	۳
۴	۶	۴
۵	۸	۵

الف) مدل رگرسیون را تخمین بزنید.

ب) اگر در آینده، هزینه تبلیغات ماهی ۶۰ هزار تومان باشد، میزان فروش چقدر خواهد بود.

ج) برای سطح فروش به دست آمده در بند «ب»، یک فاصله اطمینان ۹۰ درصدی بسازید.

د) اگر هزینه تبلیغات ماهی ۶۰ هزار تومان را برای ۱۰ ماه ادامه دهیم، متوسط میزان فروش ماهانه چقدر خواهد بود.

ه) برای متوسط فروش ماهانه که در بند «د» به دست آورده‌اید فاصله اطمینان ۹۰ درصدی بسازید.

الف) بر اساس کمیتهای مندرج در جدول ۳-۱ می‌توان مدل رگرسیون مفروض را به صورت زیر تخمین زد،

$$\hat{Y}_t = 1 + 1/2 X_t$$

همچنین می‌دانیم

$$\bar{X} = 3 \quad , \quad RSS = \sum e_t^2 = 1/8 \quad , \quad \sum x_t^2 = 10$$

ب) پیش‌بینی نقطه‌ای به ازای $X_T = 6$ برابر است با

$$\hat{Y}_T = 1 + 1/2 X_T \quad ,$$

$$= 1 + 1/2 (6) = 4/2 \quad ,$$

بنابراین به ازای ماهی ۶۰ هزار تومان هزینه تبلیغات، پیش‌بینی می‌شود، فروش ماهانه این شرکت برابر ۸۲۰ هزار تومان باشد.

ج) برای به دست آوردن فاصله اطمینان ۹۰ درصد برای Y_T ، ابتدا انحراف معیار خطای پیش‌بینی را حساب می‌کنیم. براساس معادله ۳-۸ داریم:

$$\text{Var}(e_T) = \hat{\sigma}^2 \left[1 + \frac{1}{n} + \frac{(X_T - \bar{X})^2}{\sum x_t^2} \right] \quad ,$$

با توجه به

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-2} = \frac{1/8}{5-2} = 2/93 ,$$

خواهیم داشت

$$\text{Var}(e_f) = 2/93 \left[1 + \frac{1}{5} + \frac{(6-3)^2}{10} \right] = 2/93 (2/1) = 6/103$$

بنابراین انحراف معیار خطای پیش‌بینی عبارت است از

$$\text{SE}(e_f) = \sqrt{\text{Var}(e_f)} = \sqrt{6/103} = 2/48 .$$

با استفاده از نامساوی ۳-۱۱ و با توجه به اینکه مقدار t با ۳ درجه آزادی و سطح معنی‌دار ۱۰ درصد برابر با ۲/۳۵۳ است، فاصله اطمینان ۹۰ درصدی برای Y_f به صورت زیر خواهد بود،

$$\hat{Y}_f - t_{\alpha/2} \text{SE}(e_f) < Y_f < \hat{Y}_f + t_{\alpha/2} \text{SE}(e_f) ,$$

$$1/2 - 2/353 (2/48) < Y_f < 1/2 + 2/353 (2/48) ,$$

یا

$$2/36 < Y_f < 14/04 ,$$

یعنی به ازای ماهانه ۶۰ هزار تومان هزینه تبلیغات، پیش‌بینی می‌شود که با ۹۰ درصد احتمال سطح فروش بین ۲۳۶ هزار تومان تا ۱ میلیون و ۴۰۴ هزار تومان نوسان داشته باشد.

د) در این حالت به جای Y_f باید $E(Y_f)$ را پیش‌بینی کرد. با توجه به معادله ۳-۱۳ می‌دانیم تخمین نقطه‌ای Y_f یا $E(Y_f)$ متساویند:

$$\hat{E}(Y_f) = \hat{\alpha} + \hat{\beta} X_f$$

$$= 1 + 1/2 (6) = 1/2$$

ه) همان‌گونه که در معادله ۳-۱۷ ملاحظه می‌شود واریانس خطای پیش‌بینی در

این حالت با حالت قبلی تفاوت دارد،

$$\text{Var}(e_f^m) = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2} \right]$$

با استفاده از کمیت‌های داده شده خواهیم داشت

$$\text{Var}(e_f^m) = 2/93 \left[\frac{1}{5} + \frac{(6-3)^2}{10} \right] = 2/93 (1/1) = 3/223 ,$$

بنابراین انحراف معیار خطای پیش‌بینی عبارت است از

$$\text{SE}(e_f^m) = \sqrt{\text{Var}(e_f)} = \sqrt{3/223} = 1/795 .$$

با استفاده از نامساوی ۳-۱۹ و با توجه به اینکه مقدار t با $2 = 3 - 0$ درجه آزادی و سطح معنی‌دار ۱۰ درصد برابر با ۲/۳۵۳ است؛ بنابراین فاصله اطمینان ۹۰ درصدی برای $E(Y_f)$ به صورت زیر خواهد بود،

$$\hat{Y}_f - t_{\alpha/2} \text{SE}(e_f^m) < E(Y_f) < \hat{Y}_f + t_{\alpha/2} \text{SE}(e_f^m) ,$$

$$8/2 - 2/353 (1/795) < E(Y_f) < 8/2 + 2/353 (1/795) ,$$

$$3/98 < E(Y_f) < 12/42 .$$

یعنی به ازای ماهانه ۶۰ هزار تومان هزینه تبلیغات و با این فرض که این هزینه تا ۱۰ ماه ادامه یابد، پیش‌بینی می‌شود که با ۹۰ درصد احتمال متوسط سطح فروش در این دوره ۱۰ ماهه بین ۳۹۸ هزار تومان تا ۱ میلیون و ۲۴۲ هزار تومان در نوسان باشد. ملاحظه می‌شود که فاصله اطمینان در این حالت از حالت قبلی کمتر است. دلیل این امر کوچکتر بودن واریانس e_f^m نسبت به e_f می‌باشد که به معنای دقیقتر بودن تخمین و در نتیجه کوچکتر شدن فاصله اطمینان است.

۳. دقت پیش‌بینی

در قسمتهای قبل دیدیم که می‌توان فاصله‌های اطمینان را برای Y_f و $E(Y_f)$ به ترتیب از

نامساویهای ۳-۱۰ و ۳-۱۸ به دست آورد. مطلوب این است که دقت پیش‌بینی حداکثر باشد. فاصله‌های اطمینان بزرگ نشان‌دهنده دقت‌های کم در پیش‌بینی است، برعکس هر چه دقت در پیش‌بینی بیشتر باشد، فاصله‌های اطمینان کوچکتر می‌شود. سؤال این است که در چه شرایطی می‌توان به حداکثر دقت در پیش‌بینی رسید؟
نامساویهای ۳-۱۰ و ۳-۱۸ را یک بار دیگر می‌نویسیم،

$$\hat{Y}_f - t_{\alpha/2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}} < Y_f < \hat{Y}_f + t_{\alpha/2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}},$$

یا به طور خلاصه

$$\hat{Y}_f - t \text{ SE}(e_f) < Y_f < \hat{Y}_f + \hat{Y}_f - t \text{ SE}(e_f).$$

و برای $E(Y_f)$ داریم

$$\hat{Y}_f - t_{\alpha/2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}} < E(Y_f) < \hat{Y}_f + t_{\alpha/2} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}},$$

یا به طور خلاصه

$$\hat{Y}_f - t \text{ SE}(e_f^m) < E(Y_f) < \hat{Y}_f + t \text{ SE}(e_f^m).$$

ملاحظه می‌شود مهمترین عاملی که می‌تواند فاصله‌های اطمینان برای Y_f و $E(Y_f)$ را حداقل کند، این است که انحراف معیار e_f یا e_f^m حداقل شود. سؤال این است که چگونه می‌توان این انحراف معیارها را حداقل کرد.

بدیهی است با افزایش حجم نمونه n ، مقدار $\frac{1}{n}$ کاهش می‌یابد و در نتیجه $\text{SE}(e_f)$

و $\text{SE}(e_f^m)$ نیز کمتر می‌شود. با وجود این افزایش حجم نمونه، حدی دارد و معمولاً محدودیتهای آماری به ما اجازه نمی‌دهد که n را به دلخواه و در ابعاد بزرگ افزایش دهیم؛ بنابراین باید توجه خود را به کسر $\frac{(X_f - \bar{X})^2}{\sum x_i^2}$ معطوف کنیم. برای رسیدن به

مقدار حداقل این کسر، صورت و مخرج آن باید به ترتیب حداکثر و حداقل شود. مخرج کسر در واقع پراکنندگی مشاهدات X_f را اندازه‌گیری می‌کند. قاعدتاً نمی‌توان

به طور دقیق شرایطی را تعریف کرد که این پراکندگی حداکثر شود. تنها کاری که می توان انجام داد این است که با نمونه گیریهای صحیح و تصادفی به سمت پراکندگیهای بیشتر در X_t میل کنیم. می دانیم افزایش پراکندگی X_t به همراه افزایش حجم نمونه، شرایط عمومی افزایش دقت در تخمین پارامترهاست که قبلاً در معادله ۲-۴۲ به طور خلاصه مطرح شد. با مدلی که دقیقتر تخمین زده شده است می توان بهتر و دقیقتر پیش بینی کرد.

حداقل کردن صورت کسر اهمیت فراوان دارد با توجه به این شرط میتوان گفت که به ازای مقادیری از X_t که به \bar{X} نزدیکتر باشد، فاصله اطمینان از Y_t یا $E(Y_t)$ کوچکتر خواهد بود. می دانیم پراکندگی X_t با توجه به مقدار میانگین آن، یعنی \bar{X} ، مشخص می شود. بنابراین اگر پیش بینی به ازای مقادیری از متغیر برونزا صورت گیرد که در محدوده تغییرات X_t قرار دارد، مکانیسم پیش بینی در محدوده ای انجام خواهد شد که دقت تخمینها بسیار بالاست، و در نتیجه پیش بینی ما نیز از دقت خوبی برخوردار خواهد بود، یعنی واریانس کمتری خواهد داشت. هر چه X_t به سمت \bar{X} میل کند، بیشتر در محورهای مرکزی تغییرات X_t قرار می گیرد و پیش بینی ما دارای واریانس کمتر و دقت بیشتر خواهد بود. آوردن مثالی در این مورد مفید است.

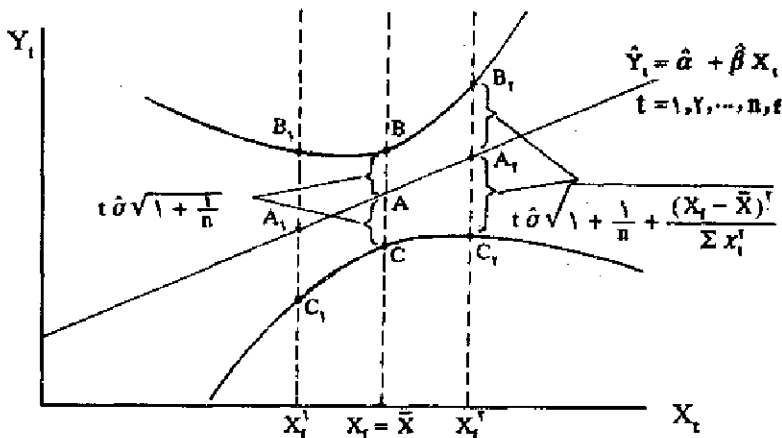
فرض کنید حجم واردات را تابعی از سطح صادرات مواد نفتی گرفته و می خواهیم یک مدل رگرسیون بر این مبنا تخمین بزنیم. یک نمونه ۱۰ تایی از متغیر برونزا داریم که صادرات نفت را طی ۱۰ سال گذشته نشان می دهد. یادآوری می کنیم که این ۱۰ عدد و درجه پراکندگی آنها (واریانس X_t)، از مهمترین عوامل تخمین پارامترهای مدل است؛ زیرا به کمک تغییرات X_t است که می خواهیم تغییرات Y_t را توضیح دهیم. فرض کنید میانگین صادرات نفت \bar{X} برابر ۲۰ میلیارد دلار است، سؤال این است که اگر در صادرات آینده نفت، مثلاً در پنجمین سال - X_t برابر ۲۰۰ میلیارد دلار باشد، فاصله اطمینان ۹۵ درصد برای حجم واردات در آن سال چقدر خواهد بود؟

بدیهی است با چنین مقدار از X_t - که بسیار دور از \bar{X} است - مقدار خطای پیش بینی دارای واریانسی بسیار بزرگ بوده و پیش بینیهای ما از دقت بسیار کمی

برخوردار خواهد بود. وقتی می‌گوییم X_t بسیار دور از \bar{X} است، در واقع می‌خواهیم بگوییم که X_t از واقعیت بسیار دور است. توجه داریم که «دور از واقعیت» بدین معنی نیست که بعد از ۵ سال صادرات نفت مطلقاً نمی‌تواند به مرز ۲۰۰ میلیارد دلار برسد. ممکن است این طور باشد یا نباشد. به هر حال مدل رگرسیون مفروض برای پیش‌بینی صادرات مواد نفتی طراحی نشده است. منظور ما از واقعیتها، در واقع اوضاع حاکم بر ساختار مدلی است که بر اساس مشاهدات گذشته، ساختار مدلی را تخمین می‌زند که می‌خواهد برای آینده پیش‌بینی کند. بنابراین هر چه مقدار X_t به \bar{X} نزدیکتر شود، پیش‌بینی مدل با شرایط حاکم بر ساختار مدل هماهنگی بیشتری خواهد داشت. به ازای $X_t = \bar{X}$ ، مقدار پیش‌بینی دقیقاً در محور تغییرات X_t صورت می‌گیرد و بهترین مقادیر پیش‌بینی را نتیجه خواهد داد.

از نظر ریاضی، وقتی X_t با \bar{X} برابر شود، صورت کسر صفر شده و مقدار حداقل کسر حاصل می‌شود. با حذف این کسر، مقدار $SE(e_t)$ یا $SE(e_t^2)$ حداقل می‌شود و فاصله‌های اطمینان برای Y_t و $E(Y_t)$ نیز حداقل خواهد شد. البته در عمل هیچگاه کوشش ما این نیست که این فاصله‌ها حداقل شود؛ زیرا ما نمی‌توانیم X_t را به دلخواه انتخاب کنیم. مقدار X_t به عنوان متغیر برون‌زا در خارج از سیستم تعیین می‌شود و ما فقط می‌توانیم مقادیر Y_t یا $E(Y_t)$ را به ازای آن پیش‌بینی کنیم. معیار $X_t = \bar{X}$ تنها شاخصی است که با توجه به آن می‌توان، قبل از محاسبه فاصله‌های اطمینان، گفت که آیا پیش‌بینی‌های خوبی خواهیم داشت یا خیر. هر چقدر X_t به \bar{X} نزدیکتر باشد، فاصله‌های اطمینان برای Y_t و $E(Y_t)$ کوچکتر و دقت پیش‌بینی بیشتر خواهد شد و برعکس نمودار (۳-۱).

تخمین مدل رگرسیون ($\hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t$) را در روی صفحه مختصات رسم می‌کنیم. t می‌تواند مقدار f را نیز بگیرد. ابتدا نقطه \bar{X} را در روی محور X_t مشخص نموده و می‌گوییم اگر بخواهیم به ازای $X_t = \bar{X}$ برای Y_t پیش‌بینی کنیم، پیش‌بینی نقطه‌ای ما از Y_t ، یعنی \hat{Y}_t ، برابر مقدار $A\bar{X}$ خواهد بود. برای رسیدن به پیش‌بینی فاصله‌ای، باید مقدار $SE(e_t)$ را یک بار به $A\bar{X}$ اضافه و بار دیگر از آن کم کنیم تا به ترتیب نقاط B و



نمودار ۳.۱ فاصله‌های اطمینان برای Y_t

C به دست آید. بنابراین پیش‌بینی فاصله‌ای برای Y_t به ازای $X_t = \bar{X}$ برابر است با

$$A\bar{X} - AC < Y_t < A\bar{X} + AB ,$$

که در آن $AB = AC$ چون در نقطه $X_t = \bar{X}$ مقدار $\frac{(X_t - \bar{X})^2}{\sum x_i^2}$ صفر می‌شود، در این نقطه مقدار $t\hat{\sigma}\sqrt{1 + \frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}}$ برابر $t\hat{\sigma}\sqrt{1 + \frac{1}{n}}$ خواهد شد که در واقع کمترین مقدار ممکن آن است. در نتیجه AB یا AC کمترین مقدار خود را در این نقطه خواهند داشت. به ازای مقادیری از X_t که از \bar{X} بزرگتر یا کوچکتر باشد، مقادیر AB یا AC بزرگتر خواهد شد. برای مثال، به ازای X_t^1 ، مقدار $A_1B_1 = A_1C_1$ به دست می‌آید، که قطعاً $A_1B_1 > AB$ و $A_1C_1 > AC$. با توجه به اینکه انحراف X_t از \bar{X} در $SE(e)$ به صورت مربع ظاهر می‌شود، قرینگی دارد؛ یعنی فرقی نمی‌کند که X_t از \bar{X} بزرگتر باشد یا کوچکتر. بنابراین اگر $X_t^1 \bar{X}$ برابر $X_t^2 \bar{X}$ باشد، خواهیم داشت $B_1C_1 = B_2C_2$. یا حرکت X_t بر روی محور X ، مجموعه نقاط B_1, B_2, \dots و نیز C_1, C_2, \dots به دست می‌آید که دو منحنی را در صفحه مختصات نشان می‌دهد. فضای بین این دو منحنی منعکس‌کننده تغییرات فاصله‌های اطمینان Y_t به ازای تغییر مکان X_t در روی محور X است.

حال که با مسأله پیش‌بینی در مدل‌های رگرسیون خطی آشنا شدیم، به ذکر مباحث تکمیلی می‌پردازیم. این مباحث در واقع مجموعه‌ای از نکات مهم در تخمین مدل‌های خطی است که فرصت مناسب برای طرح آنها در فصل‌های قبل نبوده است. تغییر مقیاس در اندازه‌گیری متغیرها، رگرسیون معکوس، تأثیر مقادیر دور افتاده متغیرهای درون‌زا یا برون‌زا و تبدیل مدل‌های رگرسیونی از جمله عناوین مهمی است که در ادامه این فصل به آنها اشاره خواهد شد.

۳-۳ تغییر مقیاس در متغیرها

در مباحثی که تا به حال داشتیم، هیچ‌گونه اشاره‌ای به مقیاس متغیرهای موجود در یک مدل رگرسیون نشده است. در عمل، می‌توان متغیرهای برون‌زا و درون‌زا را با مقیاس‌های مختلف وارد مدل رگرسیون کرد. برای مثال، وقتی متغیرهای ما تعداد بسیاری صفر در سمت راست دارند معمولاً بهتر است صفرها را تا حد ممکن حذف کرد و در عوض مقیاس متغیرها را بالا برد و آنها را با مقیاس ۱۰ هزار یا میلیون و مانند آن بیان کرد. در این قسمت پاسخ به این سؤال را مطرح می‌کنیم که اگر مقیاس اندازه‌گیری متغیرهای برون‌زا و درون‌زا را تغییر دهیم، آیا تخمین پارامترهای مدل، تغییر می‌کند؟ بحث را به صورت یک مثال مطرح کرده و سه حالت را در نظر می‌گیریم، تغییر مقیاس در متغیر برون‌زا، تغییر مقیاس در متغیر درون‌زا و تغییر مقیاس در هر دو متغیر.

۱. تغییر مقیاس در متغیر برون‌زا

تابع تقاضا برای نان را در نظر می‌گیریم که در آن Y_t ، مقدار تقاضا شده برحسب کیلوگرم در هفته و X_t ، قیمت نان برحسب ریال برای هر کیلوگرم است

$$Y_t = \alpha + \beta X_t + U_t \quad (۳-۲۰)$$

حال مقیاس قیمت نان (X_t) را تغییر داده و بجای ریال آن را با تومان اندازه‌گیری می‌کنیم. اگر X_t^* قیمت نان برحسب تومان باشد، مدل جدید به صورت زیر خواهد بود،

$$Y_t = \alpha^* + \beta^* X_t^* + U_t \quad (3-21)$$

توجه داریم که پیش فرض مدل ۳-۲۱ این است که ورود X_t^* به مدل، هم تخمین α و هم تخمین β را تغییر می دهد، در حالی که ممکن است اینچنین نباشد؛ با وجود این، برای اینکه در این مقطع از بحث، حالت عمومی را حفظ کرده باشیم، فرض را بر تغییر هر دو می گذاریم. تخمین β^* و α^* با روش حداقل مربعات معمولی عبارت است از

$$\hat{\beta}^* = \frac{\sum x_t^* y_t}{\sum x_t^{*2}} \quad (3-22)$$

$$\hat{\alpha}^* = \bar{Y} - \beta^* \bar{X}^* \quad (3-23)$$

با توجه به تغییر مقیاس در متغیر X_t ، می توان نوشت

$$X_t^* = \lambda X_t \quad (3-24)$$

که بر فرض تغییر مقیاس قیمت نان از ریال به تومان، مقدار λ برابر ۱۰۰ خواهد بود. می توان معادله ۳-۲۴ را برحسب انحراف از میانگین نیز نوشت،

$$x_t^* = \lambda x_t \quad (3-25)$$

با جایگزینی در معادله ۳-۲۲ خواهیم داشت

$$\hat{\beta}^* = \frac{\sum (\lambda x_t y_t)}{\sum (\lambda x_t)^2} = \frac{\lambda \sum x_t y_t}{\lambda^2 \sum x_t^2} = \frac{1}{\lambda} \cdot \frac{\sum x_t y_t}{\sum x_t^2}$$

در نتیجه داریم

$$\hat{\beta}^* = \frac{1}{\lambda} \hat{\beta} \quad (3-26)$$

می توان از معادله ۳-۲۴ نتیجه گرفت که $\bar{X}^* = \lambda \bar{X}$ با جایگزینی آن و همراه با جایگزینی ۳-۲۶ در معادله ۳-۲۳ خواهیم داشت

$$\hat{\alpha}^* = \bar{Y} - \frac{1}{\lambda} \hat{\beta} (\lambda \bar{X}) ,$$

$$= \bar{Y} - \hat{\beta} \bar{X} .$$

بنابراین نتیجه می‌گیریم که

$$\hat{\alpha}^* = \hat{\alpha} , \quad (۳-۲۷)$$

یعنی تغییر مقیاس متغیر برون‌زا به اندازه λ باعث می‌شود که تخمین ضریب آن متغیر به اندازه $\frac{1}{\lambda}$ تغییر کند؛ در حالی که تخمین ضریب ثابت مدل رگرسیون بدون تغییر باقی می‌ماند.

نتیجه دیگری که از بحث فوق می‌توان به دست آورد این است که اگر یک مدل رگرسیون را تخمین زده و $\hat{\alpha}$ و $\hat{\beta}$ را محاسبه کرده باشیم و سپس تصمیم به تغییر مقیاس متغیر برون‌زا بگیریم به تخمین مجدد پارامترها نیازی نیست. کافی است $\hat{\beta}$ را در معکوس مقیاس تغییر ضرب کنیم تا تخمین جدیدی از شیب مدل داشته باشیم. در $\hat{\alpha}$ نیز هیچگونه تغییری ایجاد نمی‌شود.

می‌توان تأثیر تغییر مقیاس در متغیر برون‌زا بر r^2 و بر مجموع مربعات پسماند را به صورت یک بحث تکمیلی بررسی کرد. با توجه به معادله ۱-۳۴ ضریب تعیین برای مدل اولیه عبارت است از

$$r^2 = \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2} ,$$

که برای مدل جدید به صورت زیر خواهد بود،

$$r^{*2} = \frac{(\sum x_i^* y_i)^2}{\sum x_i^{*2} \sum y_i^2} .$$

معادله ۳-۲۵ را در فرمول فوق جایگزین می‌کنیم،

$$r^{*2} = \frac{\sum (\lambda x_i y_i)^2}{\sum (\lambda x_i)^2 \sum y_i^2} = \frac{\lambda^2 (\sum x_i y_i)^2}{\lambda^2 \sum x_i^2 \sum y_i^2}$$

در نتیجه می توان نوشت

$$r^{*2} = r^2, \quad (3-28)$$

یعنی تغییر مقیاس متغیر برونزا بر ضریب تعیین تأثیری ندارد.

برای ارزیابی تأثیر متغیر مقیاس متغیر برونزا بر مجموع مربعات پسماند (RSS)، ابتدا با توجه به معادله ۲-۴۷ مقدار مجموع مربعات پسماند را برای مدل اولیه می نویسیم،

$$RSS = \sum e_i^2 = \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2},$$

که برای مدل جدید به صورت زیر خواهد بود

$$RSS^* = \sum y_i^{*2} - \frac{(\sum x_i^* y_i^*)^2}{\sum x_i^{*2}}.$$

معادله ۳-۲۵ را در فرمول فوق قرار می دهیم.

$$\begin{aligned} RSS^* &= \sum y_i^2 - \frac{(\sum \lambda x_i y_i)^2}{\sum (\lambda x_i)^2} = \sum y_i^2 - \frac{\lambda^2 (\sum x_i y_i)^2}{\lambda^2 \sum x_i^2} \\ &= \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2}, \end{aligned}$$

در نتیجه

$$RSS^* = RSS, \quad (3-29)$$

یعنی تغییر مقیاس متغیر برونزا، مجموع مربعات پسماند را تغییر نمی دهد. بدیهی است تخمین واریانس U_1 نیز تغییر نمی کند و در نتیجه آزمون فرضیه های مختلف در مورد پارامترهای مدل نیز متأثر نخواهد شد.

۲. تغییر مقیاس در متغیر درون‌زا
مانند گذشته مدل رگرسیون زیر را در نظر گرفته

$$Y_i = \alpha + \beta X_i + U_i ,$$

و این بار مقیاس متغیر درون‌زا (Y) را به اندازه μ تغییر می‌دهیم، به گونه‌ای که

$$Y_i^* = \mu Y_i . \quad (۳-۳۰)$$

سؤال این است که آیا با تغییر Y_i به Y_i^* ، تخمین پارامترهای α و β تغییر خواهد کرد؟
همچنین آیا این مسأله بر ضریب تعیین r^2 و نیز بر مجموع مربعات پسماند $\sum e_i^2 = RSS$
تأثیر می‌گذارد؟

برای پاسخ به این سؤال دقیقاً مطابق حالت اول عمل می‌کنیم. ابتدا دو طرف مدل
اولیه را در μ ضرب می‌کنیم،

$$\mu Y_i = \mu \alpha + \mu \beta X_i + \mu U_i ,$$

یا

$$Y_i^* = \alpha^* + \beta^* X_i + U_i^* ,$$

β^* را تخمین می‌زنیم،

$$\hat{\beta}^* = \frac{\sum x_i y_i^*}{\sum x_i^2} .$$

از معادله ۳-۳۰ می‌دانیم که $y_i^* = \mu y_i$ با جایگزینی در رابطه فوق خواهیم داشت

$$\hat{\beta}^* = \frac{\sum x_i \mu y_i}{\sum x_i^2} = \mu \frac{\sum x_i y_i}{\sum x_i^2} .$$

نتیجه می‌گیریم که

$$\hat{\beta}^* = \mu \hat{\beta} , \quad (۳-۳۱)$$

یعنی تغییر مقیاس متغیر درون‌زا به اندازه μ ، باعث می‌شود که تخمین جدید از β

به اندازه μ برابر تخمین قبلی باشد.

برای تخمین α^* رابطه زیر را در نظر می‌گیریم

$$\hat{\alpha}^* = \bar{Y}^* - \hat{\beta}^* \bar{X} ,$$

که با استفاده از معادله ۳-۳۱ خواهیم داشت

$$\hat{\alpha}^* = \mu \bar{Y} - \mu \hat{\beta} \bar{X} ,$$

$$= \mu (\bar{Y} - \hat{\beta} \bar{X}) .$$

نتیجه می‌گیریم که

$$\hat{\alpha}^* = \mu \hat{\alpha} , \quad (3-32)$$

یعنی تغییر مقیاس متغیر درون‌زا به اندازه μ ، باعث می‌شود که تخمین جدید از α به اندازه μ برابر تخمین قبلی باشد.

با توجه به معادله ۱-۳۴ می‌دانیم ضریب تعیین برای مدل اولیه عبارت است از

$$r^2 = \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2} ,$$

که برای مدل جدید به صورت زیر خواهد بود

$$\begin{aligned} r^{*2} &= \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2} \\ &= \frac{[\sum x_i (\mu y_i)]^2}{\sum x_i^2 \sum (\mu y_i)^2} = \frac{\mu^2 (\sum x_i y_i)^2}{\mu^2 \sum x_i^2 \sum y_i^2} , \end{aligned}$$

در نتیجه خواهیم داشت

$$r^{*2} = r^2 , \quad (3-33)$$

یعنی تغییر مقیاس متغیر درون‌زا بر ضریب تعیین تأثیری ندارد.

برای ارزیابی تأثیر تغییر مقیاس متغیر درون‌زا بر مجموع مربعات پسماند (RSS) ابتدا با توجه به معادله ۲-۴۷، مقدار مجموع مربعات پسماند را برای مدل اولیه می‌نویسیم،

$$RSS = \sum e_i^2 = \sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2},$$

که برای مدل جدید به صورت زیر خواهد بود

$$\begin{aligned} RSS^* &= \sum y_i^{*2} - \frac{(\sum x_i y_i^*)^2}{\sum x_i^2}, \\ &= \sum (\mu y_i)^2 - \frac{(\sum \mu x_i y_i)^2}{\sum x_i^2}, \\ &= \mu^2 \sum y_i^2 - \mu^2 \frac{(\sum x_i y_i)^2}{\sum x_i^2}, \end{aligned}$$

در نتیجه

$$= RSS^* = \mu^2 RSS, \quad (3-34)$$

یعنی تغییر مقیاس متغیر درون‌زا، مجموع مربعات پسماند را تغییر می‌دهد و بر تخمین واریانس U_e و نیز آزمون فرضیه‌های مختلف تأثیر خواهد گذاشت؛ به عبارت دیگر

$$\begin{aligned} \hat{\sigma}_U^{*2} &= \frac{\sum e_i^{*2}}{n-2} = \frac{RSS^*}{n-2}, \\ &= \frac{\mu^2 RSS}{n-2} = \mu^2 \hat{\sigma}^2, \end{aligned}$$

در نتیجه خطای معیار U_e یا خطای معیار تخمین (SEE) در مدل جدید دقیقاً μ برابر خطای معیار تخمین مدل اولیه است،

$$SEE^* = \sqrt{\hat{\sigma}_U^{*2}} = \mu SEE.$$

۳. تغییر مقیاس در متغیرهای برونزا و درونزا
مدل رگرسیون اولیه را می‌نویسیم،

$$Y_i = \alpha + \beta X_i + U_i .$$

مقیاس متغیر برونزا را به اندازه λ و مقیاس متغیر درونزا را به اندازه μ تغییر می‌دهیم،

$$X_i^* = \lambda X_i ,$$

$$Y_i^* = \mu Y_i .$$

از رابطه اول داریم، $X_i = \frac{1}{\lambda} X_i^*$ ، که اگر آن را با معادله دوم در معادله

$$\mu Y_i = \mu \alpha + \mu \beta X_i + \mu U_i$$

قرار دهیم، خواهیم داشت

$$Y_i^* = \mu \alpha + \mu \beta \left(\frac{1}{\lambda} X_i^* \right) + \mu U_i ,$$

یا

$$Y_i^* = \alpha^* + \beta^* X_i^* + U_i^* , \quad (3-35)$$

که در آن $\alpha^* = \mu \alpha$ و $\beta^* = \frac{\mu}{\lambda} \beta$. سؤال این است که آیا تخمین پارامترها در مدل ۳-۳۵ با مقادیر مشابه در مدل اولیه تفاوت دارد؟

برای پاسخ، کافی است به ترتیب $\hat{\beta}^*$ و $\hat{\alpha}^*$ را به دست آورده با $\hat{\beta}$ و $\hat{\alpha}$ مقایسه کنیم.

می‌دانیم

$$\begin{aligned} \hat{\beta}^* &= \frac{\sum x_i^* y_i^*}{\sum x_i^{*2}} , \\ &= \frac{\sum (\lambda x_i) (\mu y_i)}{\sum (\lambda x_i)^2} = \frac{\mu \lambda \sum x_i y_i}{\lambda^2 \sum x_i^2} , \end{aligned}$$

در نتیجه خواهیم داشت

$$\hat{\beta}^* = \frac{\mu}{\lambda} \hat{\beta} . \quad (۳-۳۶)$$

یعنی برای رسیدن به تخمین جدید از β ، باید تخمین قبلی را در نسبت دو مقیاس $(\frac{\mu}{\lambda})$ ضرب کنیم. برای تخمین $\hat{\alpha}^*$ رابطه زیر را در نظر می‌گیریم

$$\hat{\alpha}^* = \bar{Y}^* - \hat{\beta}^* \bar{X}^* ,$$

که با استفاده از معادله ۳-۳۶ خواهیم داشت

$$\begin{aligned} \hat{\alpha}^* &= \frac{\sum \mu Y_i}{n} - \frac{\mu}{\lambda} \hat{\beta} \frac{\sum \lambda X_i}{n} , \\ &= \mu \bar{Y} - \frac{\mu}{\lambda} \hat{\beta} \lambda \bar{X} , \\ &= \mu \bar{Y} - \mu \hat{\beta} \bar{X} = \mu (\bar{Y} - \hat{\beta} \bar{X}) , \end{aligned}$$

بنابراین

$$\hat{\alpha}^* = \mu \hat{\alpha} . \quad (۳-۳۷)$$

نتیجه کلی این است که تغییر مقیاس، تخمین پارامترها را دقیقاً به همان صورت تغییر می‌دهد که بر رابطه بین پارامترهای واقعی تأثیر دارد؛ به عبارت دیگر، در معادله ۳-۳۵ دیدیم که

$$\beta^* = \frac{\mu}{\lambda} \beta , \quad \alpha^* = \mu \alpha .$$

معادله‌های فوق تأثیر مستقیم تغییر مقیاس بر رابطه بین پارامترهای واقعی را نشان می‌دهد. در معادله‌های ۳-۳۶ و ۳-۳۷ نیز دیدیم که

$$\hat{\beta}^* = \frac{\mu}{\lambda} \hat{\beta} , \quad \hat{\alpha}^* = \mu \hat{\alpha}$$

دقیقاً همان ساختار روابط قبلی را دارد.

برای بررسی آثار تغییر مقیاس بر r^2 و بر مجموع مربعات پسماند و خطای معیار

تخمین کافی است که فرمول هر یک از معیارهای فوق را بعد از تغییر مقیاس بررسی کنیم. برای مدل جدید، یعنی معادله ۳-۳۵ داریم

$$r^{*2} = \frac{(\sum x_i^* y_i^*)^2}{\sum x_i^{*2} \sum y_i^{*2}}$$

با جایگزینی $x_i^* = \lambda x_i$ و $y_i^* = \mu y_i$ در معادله فوق خواهیم داشت

$$\begin{aligned} r^{*2} &= \frac{(\sum \lambda x_i \mu y_i)^2}{\sum (\lambda x_i)^2 \cdot \sum (\mu y_i)^2} \\ &= \frac{\lambda^2 \mu^2 (\sum x_i y_i)^2}{\lambda^2 \sum x_i^2 \cdot \mu^2 \sum y_i^2} \end{aligned}$$

در نتیجه

$$r^{*2} = r^2 \quad (3-38)$$

یعنی تغییر همزمان مقیاس متغیرهای درونزا و برونزا، بر ضریب تعیین تأثیری ندارد. برای بررسی خطای معیار تخمین کافی است معادله ۲-۴۷ را برای مدل ۳-۳۵

بنویسیم:

$$\begin{aligned} RSS^* &= \sum y_i^{*2} - \frac{(\sum x_i^* y_i^*)^2}{\sum x_i^{*2}} \\ &= \sum (\mu y_i)^2 - \frac{(\sum \lambda x_i \mu y_i)^2}{\sum (\lambda x_i)^2} \\ &= \mu^2 \sum y_i^2 - \frac{\lambda^2 \mu^2 (\sum x_i y_i)^2}{\lambda^2 \sum x_i^2} \\ &= \mu^2 \left[\sum y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} \right] \end{aligned}$$

در نتیجه

$$RSS^* = \mu^T RSS, \quad (3.39)$$

یعنی مجموع مربعات پسماند در مدل جدید، فقط می‌تواند به اندازه μ^T ، مقدار مشابه از مدل اولیه را تغییر دهد؛ در حالی که λ مطلقاً در مجموع مربعات پسماند جدید تأثیری ندارد. به همین ترتیب داریم

$$\hat{\sigma}_U^{*2} = \frac{RSS^*}{n-2} = \frac{\mu^T RSS}{n-2} = \mu^T \hat{\sigma}_U^2,$$

و همچنین خواهیم داشت

$$SEE^* = \sqrt{\hat{\sigma}_U^{*2}} = \mu SEE, \quad (3.40)$$

یعنی تغییر همزمان مقیاس متغیرهای برون‌زا و درون‌زا را به اندازه‌های λ و μ ، واریانس U_i را به اندازه μ^2 و خطای معیار تخمین را به اندازه μ تغییر می‌دهد.

۳-۴ رگرسیون معکوس

مدلی را که تا اینجا موضوع مطالعه ما بوده است یک بار دیگر ملاحظه کنید،

$$Y_i = \alpha + \beta X_i + U_i.$$

در این مدل اصطلاحاً می‌گوییم Y_i را روی X_i «رگرس» کرده‌ایم؛ یعنی Y_i را متغیر وابسته و X_i را متغیر توضیحی گرفته‌ایم. به مدل فوق، یک مدل «رگرسیون مستقیم»^۱ نیز می‌گویند. در مواردی ضروری است که رگرسیون X_i روی Y_i را بررسی کنیم؛ یعنی برعکس مدل فوق عمل کرده، در واقع X_i را متغیر وابسته و Y_i را متغیر توضیحی در نظر بگیریم. خواهیم داشت

$$X_i = \alpha' + \beta' Y_i + V_i.$$

مدل فوق را یک مدل «رگرسیون معکوس»^۲ می‌نامیم.

برای مثال، اگر بخواهیم تأثیر مسأله تبعیض جنسیت یا تبعیض نژادی را بر سطح دستمزدها بررسی کنیم، می‌توان از مدل‌های رگرسیون معکوس استفاده کرد. فرض کنید مشاهدات ما به صورت مقطعی است. Y_i سطح دستمزدها و X_i شاخص صلاحیت و تخصص نیروی کار است. اگر مسأله مورد نظر ما تأثیر تبعیض جنسیت بر دستمزدها باشد آنگاه می‌توان دو سؤال زیر را مطرح کرد.

الف) آیا مردان و زنانی که از یک سطح تخصصی معین برخوردارند (مقادیر X_i آنها یکسان است) دستمزدهای یکسان دارند؟ می‌توان این سؤال را با یک مدل رگرسیون مستقیم (مدل رگرسیون Y_i روی X_i) پاسخ داد.

ب) آیا مردان و زنانی که دستمزدهای یکسان دارند (مقادیر Y_i آنها مساوی است) از تخصصهای یکسان بهره‌مند هستند؟ می‌توان پاسخ به این سؤال را با یک مدل رگرسیون معکوس (مدلی که X_i را بر حسب Y_i بیان می‌کند) به دست آورد.

در مورد این مثال خاص، مدل‌های رگرسیون مستقیم و معکوس به ترتیب عبارتند

از

$$Y_i = \alpha + \beta X_i + U_i ,$$

$$X_i = \alpha' + \beta' Y_i + V_i , \quad (3-41)$$

که در آن جمله اختلال در رگرسیون معکوس بوده و شامل تمام فرضهای کلاسیک است.

با روشی همانند تخمین پارامترها در مدل رگرسیون مستقیم، می‌توان $\hat{\alpha}'$ و $\hat{\beta}'$ را به دست آورد. خواهیم داشت

$$\hat{\beta}' = \frac{\sum x_i y_i}{\sum y_i} , \quad (3-42)$$

$$\hat{\alpha}' = \bar{X} - \hat{\beta}' \bar{Y} . \quad (3-43)$$

همانند فرمول ۲-۴۷ - که برای به دست آوردن مجموع مربعات پسماند (RSS) در مدل رگرسیون مستقیم داشتیم - می‌توان مجموع مربعات پسماند را در مدل رگرسیون معکوس

(RSS') از فرمول زیر به دست آورد،

$$\sum e_i'^2 = RSS' = \sum x_i'^2 - \frac{(\sum x_i y_i)^2}{\sum y_i'^2} \quad (3-44)$$

اگر $\hat{\beta}$ مانند قبل، تخمین β در مدل رگرسیون مستقیم باشد، آنگاه

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i'^2}$$

در نتیجه به سهولت ملاحظه می‌شود که

$$\hat{\beta} \hat{\beta}' = \frac{\sum x_i y_i}{\sum x_i'^2} \cdot \frac{\sum x_i y_i}{\sum y_i'^2}$$

و خواهیم داشت

$$\hat{\beta} \hat{\beta}' = r_{x,y}^2 \quad (3-45)$$

یعنی حاصلضرب تخمین شیبه‌های دو خط رگرسیون مستقیم و معکوس برابر ضریب تعیین است؛ بنابراین شیب رگرسیون معکوس، هنگامی معکوس شیب رگرسیون مستقیم است که ضریب تعیین مدل برابر یک باشد.

مثال ۳-۳ مدل رابطه بین تولید و نیروی کار، موضوع مثالهای ۱-۲ و ۲-۱ و جدول ۱-۳ را یک بار دیگر ملاحظه می‌کنیم.

Q_i	۱۱	۱۰	۱۲	۶	۱۰	۷	۹	۱۰	۱۱	۱۰
L_i	۱۰	۷	۱۰	۵	۸	۸	۶	۷	۹	۱۰

با فرض $Q_i = Y_i$ و $L_i = X_i$ ، نتایج محاسبات قبلی را دوباره می‌نویسیم،

$$\bar{X} = 8, \quad \bar{Y} = 9/6,$$

$$\sum x_i'^2 = 28, \quad \sum y_i'^2 = 30/4,$$

$$\sum x_i y_i = 21, \quad r^2 = 0/52,$$

$$\hat{\beta} = 0/75, \quad \hat{\alpha} = 3/6,$$

بنابراین تخمین مدل رگرسیون مستقیم، عبارت است از

$$\hat{Y}_i = 3/6 + 0/75 X_i.$$

برای تخمین مدل رگرسیون معکوس، معادله ۳-۴۲ را می نویسیم،

$$\hat{\beta}' = \frac{\sum x_i y_i}{\sum y_i^2},$$

$$= \frac{21}{30/8} = 0/69.$$

برای تخمین $\hat{\alpha}'$ از معادله ۳-۴۳ استفاده می کنیم،

$$\hat{\alpha}' = \bar{X} - \hat{\beta}' \bar{Y},$$

$$= 8 - 0/69 (9/6) = 1/37.$$

در نتیجه مدل رگرسیون معکوس عبارت خواهد بود از

$$\hat{X}_i = 1/37 + 0/69 Y_i.$$

در محاسبات قبلی دیدیم که $r^2 = 0/52$. این نتیجه را می توان با توجه به معادله ۳-۴۵

نیز به دست آورد. کافی است تخمین شیبهای دو خط رگرسیون مستقیم و معکوس را در

هم ضرب کنیم،

$$\hat{\beta} \hat{\beta}' = r^2,$$

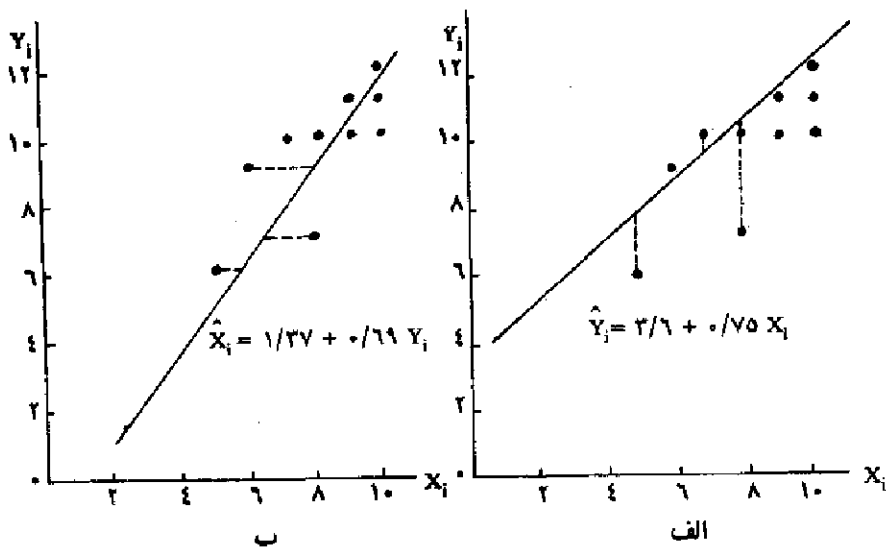
$$0/75 (0/69) = 0/52.$$

تیین بیشتر مدل‌های رگرسیون مستقیم و معکوس، مستلزم نمایش هندسی

آن‌هاست. در نمودار ۳-۲ الف تخمین مدل رگرسیون مستقیم و در نمودار ۳-۲ ب مدل

رگرسیون معکوس نشان داده شده است. هر دو نمودار در صفحه مختصات XOY رسم

شده است.



نمودار ۳-۲ نمودار هندسی رگرسیونهای مستقیم و معکوس

ملاحظه می‌شود که در نمودار ۳-۲ الف، وقتی می‌خواهیم مجموع مربعات پسماند، یعنی $\sum (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2$ برای رگرسیون مستقیم را حداقل کنیم، در واقع مجموع مربع خطوط قائمی حداقل می‌شود که فاصله نقاط مشاهده شده تا خط تخمین رگرسیون را نشان می‌دهد. برای مثال، سه خط از این فاصله‌ها به صورت نقطه‌چین مشخص شده است. اما مجموع مربعات پسماند برای مدل رگرسیون معکوس عبارت است از $\sum (X_i - \hat{\alpha}' - \hat{\beta}' Y_i)^2$ ، که حداقل کردن آن، به معنای به حداقل رساندن مجموع مربع خطوط افقی است که فاصله نقاط مشاهده شده را تا خط تخمین رگرسیون نشان می‌دهد. در نمودار ۳-۲ ب سه خط از این فاصله‌ها به صورت نقطه‌چین مشخص شده است. البته می‌توان فاصله نقاط مشاهده شده تا خط تخمین رگرسیون را، به جای خطوط قائم یا افقی، با خطوط عمود بر خط تخمین رگرسیون تعریف کرد، در این صورت می‌گوییم مدل ما یک «رگرسیون متعامد»^۱ است. که بحث آن از موضوع این کتاب خارج است. ناگفته نماند که اگر بخواهیم تخمین یک مدل رگرسیون معکوس را در صفحه مختصات

نشان دهیم، مطابق قاعده باید X_i را روی محور عمودی و Y_i را روی محور افقی اندازه گیری کنیم. اما در نمودار ۲-۳ ب، Y_i را روی محور عمودی برده ایم. علت تنها این بوده است که می خواستیم تفاوت دو مفهوم «پسماند در جهت Y_i » و «پسماند در جهت X_i » را در یک نمودار نشان دهیم. در واقع اگر e_i را به صورت زیر تعریف کنیم:

$$Y_i \text{ پسماند در جهت } = e_i = Y_i - \hat{Y}_i ,$$

یک رگرسیون مستقیم خواهیم داشت که نمایش هندسی آن نمودار ۲-۳ الف و مقادیر e_i به صورت خطوط عمودی است. اما اگر داشته باشیم

$$X_i \text{ پسماند در جهت } = e_i = X_i - \hat{X}_i ,$$

در آن صورت مدل رگرسیون ما معکوس بوده و همان گونه که در نمودار ۲-۳ ب ملاحظه می شود، مقادیر e_i به صورت خطوط افقی خواهد بود.

آخرین نکته این است که آیا معیاری وجود دارد که بتوان گفت در چه مواردی باید از رگرسیون معکوس استفاده کرد؟ پاسخ منفی است؛ اما می توان ملاحظات کلی زیر را در انتخاب نوع مدل رگرسیون در نظر گرفت.

الف) اگر نظریه های اقتصادی، جهت رابطه علت و معلولی را بین دو متغیر X و Y ، مشخص کرده باشد، به راحتی می توان نسبت به نوع مدل رگرسیون تصمیم گرفت؛ یعنی متغیری که علت تغییرات است را به عنوان متغیر برونزا در سمت راست مدل قرار می دهیم و معلول را در سمت چپ می بریم. برای مثال، فرض کنید براساس نظریه های اقتصادی و ملاحظات دیگر، مطمئن هستیم که در واحد تولیدی i ، حجم نیروی کار در زمان t بر سطح تولیدات در زمان t تأثیر می گذارد، نه برعکس. در این صورت یک مدل رگرسیون خواهیم داشت که در آن حجم نیروی کار، متغیر توضیحی یا برونزا، و سطح تولیدات، یک متغیر وابسته یا درونزا است، خواهیم داشت

$$Q_i = \alpha + \beta L_i + U_i .$$

در چنین شرایطی، ساختن یک مدل رگرسیون معکوس، مفهوم چندان روشنی ندارد و

توصیه نمی‌شود.

بنابراین نتیجه می‌گیریم که در تمام مواردی که علت تغییرات را می‌دانیم، ضروری است در مدل رگرسیون، علت به عنوان متغیر توضیحی و معلول به صورت متغیر وابسته در نظر گرفته شود. یادآوری می‌کنیم که در چهارچوب این مثال مدل رگرسیون Q_t بر I_t می‌تواند به دو سؤال زیر پاسخ دهد.

۱. به ازای مقدار معینی از حجم نیروی کار در آینده، چگونه می‌توان سطح تولیدات را پیش‌بینی کرد؟

۲. برای رسیدن به سطح معینی از تولیدات در آینده، چگونه می‌توان حجم نیروی کار لازم را تخمین زد؟

به عبارت دیگر نیازی نیست که برای پاسخ به سؤال دوم، یک مدل رگرسیون معکوس بسازیم؛ بلکه رگرسیون مستقیم کافی است.

ب) در مواردی که جهت رابطه علت و معلولی بین دو متغیر X و Y چندان روشن نیست، معمولاً هر دو مدل رگرسیون Y روی X و X روی Y می‌تواند مفید باشد. در این موارد بهتر است ابتدا هر دو مدل رگرسیون را تخمین زده، سپس با توجه به نتایج تخمین و شرایط مسأله مفروض، نسبت به انتخاب نهایی یکی از آن دو تصمیم بگیریم.

ج) مواردی نیز وجود دارد که از مفروضات خود مسأله می‌توان دریافت که آیا ساختن یک مدل رگرسیون معکوس ضروری است یا خیر؛ برای مثال، در بحث مسأله تبعیض جنسیت و سطح دستمزدها در معادله ۳-۴۱ دیدیم که طراحی مدل رگرسیون معکوس می‌تواند روشنگر باشد.

د) و سرانجام این مسأله که آیا رگرسیون مستقیم بر رگرسیون معکوس برتری دارد یا برعکس، در بسیاری موارد به چگونگی تولید آمار و کیفیت عرضه آن بر می‌گردد؛ برای مثال، جدول مثال ۳-۳ را دوباره ملاحظه کنید. I_t یا X_t ساعتهای کار انجام شده و Q_t یا Y_t حجم تولیدات بوده است. تعداد ۱۰ مشاهده مندرج در این جدول مربوط به ۱۰ کارگر است. اینکه آیا Y_t را تابعی از X_t بگیریم یا برعکس، تا حد بسیاری به این نکته بر می‌گردد که آمار مندرج در این جدول چگونه ایجاد شده است. این دو حالت را

در نظر می‌گیریم:

۱. اگر ساعت‌های معین و متفاوتی از زمان کار را به کارگران مختلف بدهیم (X_i)، سپس حجم تولیدات حاصل از این ساعت‌های کار را مشاهده و ثبت کنیم (Y_i)، بهتر است که از مدل رگرسیون Y_i بر X_i استفاده کنیم.

۲. اگر به هر یک از کارگران مقادیر مختلفی از تولیدات (Y_i) را سفارش بدهیم و سپس ساعت‌های کار انجام شده برای تولید مقادیر مفروض (X_i) را مشاهده و ثبت کنیم آنگاه مفید است که مدل رگرسیون ما X_i را بر حسب Y_i بیان کند.

یکی از معیارهایی که می‌تواند در انتخاب بین مدل‌های رگرسیون مستقیم و معکوس استفاده شود این است که متغیر تحت کنترل ما کدام است. اگر توانستیم چنین متغیری را شناسایی کنیم، آن را به عنوان متغیر توضیحی یا برون‌زا در سمت راست معادله قرار می‌دهیم؛ برای مثال، در حالت اول، متغیر تحت کنترل ما ساعت‌های کار انجام شده است که آن را به کارگران ابلاغ کرده‌ایم، بنابراین به عنوان متغیر توضیحی یا برون‌زا وارد مدل رگرسیون می‌شود. اما در حالت دوم، با تعیین حجم تولیدات و درخواست از کارگران برای تولید آنها، در واقع حجم تولیدات را کنترل می‌کنیم و ساعت‌های کار را فقط کارگران تعیین می‌کنند. به همین دلیل است که در حالت دوم (Y_i) به عنوان متغیر توضیحی وارد مدل رگرسیون می‌شود؛ در حالی که X_i متغیر درون‌زاست.

۳-۵ مشاهدات دورافتاده

در بعضی از مدل‌های رگرسیون ممکن است چند مشاهده «غیرعادی» یا «دورافتاده»^۱ بتواند بر تخمین پارامترهای مدل تأثیر بسیاری بگذارد، به گونه‌ای که از دقت آن کم کند و نتایج را غیرقابل اعتماد نماید. دقت در مشاهدات X و Y برای تفسیر نتایج حاصل از تخمین یک مدل رگرسیون همواره ضروری است. باید همیشه به دنبال این سؤال باشیم که آیا مجموعه مشاهدات ما شامل موارد دور افتاده هست یا خیر. تشخیص اینکه آیا

یک یا چند مشاهده دورافتاده است با بررسی جمله‌های پسماند صورت می‌پذیرد. اگر X_i شامل موارد دورافتاده باشد، اثر آن با معادله $\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i$ در \hat{Y}_i ظاهر می‌شود؛ در نتیجه با توجه به $e_i = Y_i - \hat{Y}_i$ ، در e_i منعکس خواهد شد. اگر مشاهدات Y_i دارای موارد دورافتاده باشد، با رابطه $e_i = Y_i - \hat{Y}_i$ ، به طور قطع این موارد بر e_i تأثیر می‌گذارد؛ بنابراین بررسی سیر تغییرات e_i برای کشف موارد دورافتاده در مشاهدات X و Y ضروری است.

معمولاً بعد از تخمین یک مدل رگرسیون ساده، فقط به $\hat{\alpha}$ ، $\hat{\beta}$ ، و σ^2 خطای معیار تخمین زنده‌های $\hat{\alpha}$ و $\hat{\beta}$ ، یعنی $SE(\hat{\alpha})$ و $SE(\hat{\beta})$ توجه می‌شود؛ در حالی که به نظر می‌رسد مفید است که مقادیر مختلف پسماندها را نیز محاسبه و به روند تغییرات آن توجه کامل کنیم. بنابراین مسأله «آنالیز پسماندها»^۱ یکی از مباحث مهم اقتصادسنجی است که باید به طور مستقل بررسی شود. در این قسمت فقط به طرح مطلب اکتفا می‌شود.

یک مشاهده دورافتاده، در واقع مشاهده‌ای است که با مشاهدات دیگر هماهنگی لازم را ندارد. دلایل مختلفی برای به وجود آمدن این مشاهدات وجود دارد که موضوع تحلیلها و نظریه‌های اقتصادی است. اما مسأله مهم این است که روش حداقل مربعات معمولی در تخمین پارامترها حتی از یک مشاهده دورافتاده نیز می‌تواند بسیار متأثر باشد. در مورد رگرسیون ساده، به راحتی می‌توان با استفاده از نمودار هندسی مشاهدات در صفحه مختصات، وجود موارد دورافتاده را کشف کرد؛ اما وقتی رگرسیون ما چند متغیره باشد، یعنی چند متغیر توضیحی را شامل شود، نمودار هندسی مشاهدات، معمولاً غیرممکن است و حتماً باید از آنالیز جمله اختلال استفاده کرد.

مثال ۳-۴ این مثال نشان می‌دهد چگونه فقط تخمین پارامترها همراه با σ^2 و انحراف معیار تخمین پارامترها، نمی‌تواند تمام خصوصیات موجود در تخمین یک مدل رگرسیون را بیان کند.^۱ جدول ۳-۲ چهار سری مشاهدات X_i و Y_i را نشان می‌دهد.

1. Analysis of Residuals

۲. به مقاله (۱۹۷۳) F. J. Anscombe مراجعه شود.

جدول ۳.۲

متغیر	X_i	Y_i	Y_i	Y_i	X_i	Y_i	
سری مشاهدات	(۱-۳)	۱	۲	۳	۴	۴	
مشاهدات	۱	۱۰	۸/۰۴	۹/۱۴	۷/۴۶	۸	۶/۵۸
	۲	۸	۶/۹۵	۸/۱۴	۶/۷۷	۸	۵/۷۶
	۳	۱۳	۷/۵۸	۸/۷۴	۱۲/۷۴	۸	۷/۷۱
	۴	۹	۸/۸۱	۸/۷۷	۷/۱۱	۸	۸
	۵	۱۱	۸/۳۳	۹/۲۶	۷/۸۱	۸	۸/۴۷
	۶	۱۴	۹/۹۶	۸/۱۰	۸/۸۴	۸	۷/۰۴
	۷	۶	۷/۲۴	۶/۱۳	۶/۰۸	۸	۵/۲۵
	۸	۴	۴/۲۶	۳/۱۰	۵/۳۹	۱۹	۱۲/۵
	۹	۱۲	۱۰/۸۴	۹/۱۳	۸/۱۵	۸	۵/۵۶
	۱۰	۷	۴/۸۲	۷/۲۶	۶/۴۲	۸	۷/۹۱
	۱۱	۵	۵/۶۸	۴/۷۴	۵/۷۳	۸	۶/۸۹

برای سه سری مشاهدات اول و دوم و سوم مربوط به Y_i ، مشاهدات X_i مندرج در ستون دوم ثابت فرض شده است. دو ستون آخر جدول منعکس کننده سری چهارم مشاهدات X_i و Y_i است. نکته جالب این است که هر چهار سری این مشاهدات در مورد X_i و Y_i به یک تخمین واحد از مدل رگرسیون $Y_i = \alpha + \beta X_i + U_i$ می رسد.

از هر چهار سری مشاهدات فوق کمیتهای زیر را داریم

$$n = 11, \quad \bar{X} = 9, \quad \bar{Y} = 7/5,$$

$$\sum x_i^2 = 110, \quad \sum y_i^2 = 41/25, \quad \sum x_i y_i = 55,$$

در نتیجه تخمین مدل رگرسیون به صورت زیر خواهد بود

$$\hat{Y}_i = 3 + 0/5 X_i \quad (0/118)$$

ملاحظه می شود که $SE(\hat{\beta}) = 0/118$. همچنین برای مدل فوق داریم:

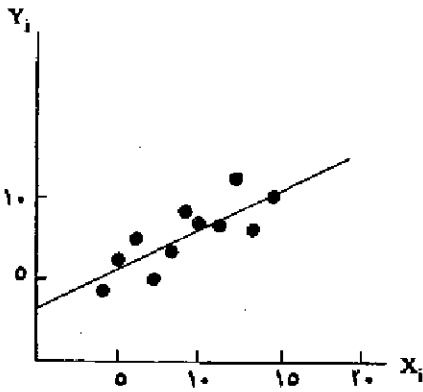
$$r^2 = 0/667,$$

$$\sum \hat{y}_i^2 = ESS = 27/5 ,$$

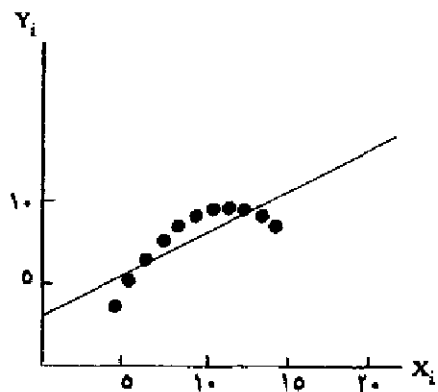
$$\sum e_i^2 = RSS = 13/50 .$$

با اینکه چهار سری مشاهدات فوق به تخمین واحدی از مدل رگرسیون مفروض منتهی می‌شود، اما به وضوح ملاحظه می‌گردد که این چهار سری مشاهده، خصوصیات کاملاً متفاوتی با یکدیگر دارند. برای تبیین این نکته کافی است که آنها را در چهار نمودار ۳-۳ منعکس کنیم. نمودار سری اول از مشاهدات مسأله خاصی ندارد؛ اما نمودار هندسی سری دوم از مشاهدات، نشان می‌دهد که یک مدل رگرسیون خطی برای تبیین تغییرات Y_1 صلاحیت کافی را ندارد زیرا همان گونه که در نمودار ملاحظه می‌کنیم تغییرات Y_2 غیرخطی است. در نمودار سری سوم از مشاهدات ملاحظه می‌کنیم که چگونه یک مشاهده دورافتاده ($Y = 12/74, X = 13$) شیب تخمین مدل رگرسیون را نسبتاً زیاد تغییر داده و آن را به اصطلاح تندتر کرده است. نکته مهم این است که اگر مقادیر مختلف پسماندها را حساب کنیم، مقدار پسماند مربوط به این مشاهده - که در واقع مشاهده سوم است - یعنی e_3 ، نسبت به مقادیر دیگر e_i ، مقدار کاملاً متفاوت و بیشتری را نشان می‌دهد. ممکن است استدلال شود که باید مشاهده سوم را حذف کرد تا دیگر مشاهدات که هماهنگی بیشتری با یکدیگر دارند تخمینهای بهتری را نتیجه دهد این نکته را در ادامه بحث دوباره مطرح خواهیم کرد. در نمودار سری چهارم از مشاهدات، تأثیر یک مشاهده دورافتاده در نمونه به وضوح منعکس شده است. اگر هشتم ($Y = 12/5, X = 19$) را از نمونه حذف کنیم، تخمین مدل رگرسیون مفروض به صورت یک خط عمودی ظاهر خواهد شد. بنابراین معلوم می‌شود که چگونه تأثیر یک مشاهده دورافتاده توانسته است وضعیت تخمین را کاملاً تغییر دهد.

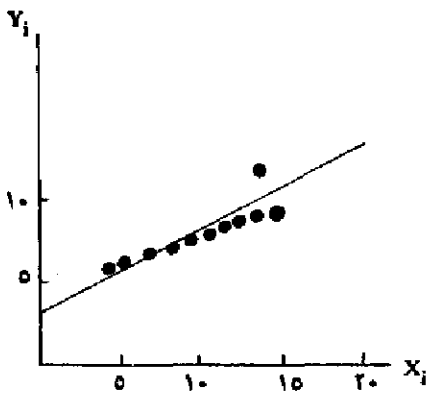
سؤال این است که آیا می‌توان مشاهدات دورافتاده را از نمونه حذف کرد؟ به این سؤال نمی‌توان پاسخ مثبت یا منفی داد؛ شرایط هر مسأله در مجموع می‌تواند رهگشا باشد، با وجود این، می‌توان این ملاحظات کلی را مطرح کرد: در مواردی که شرایط حاکم بر تعیین مقادیر مشاهدات دورافتاده کاملاً با شرایط مشاهدات دیگر متفاوت است، می‌توان تصمیم بر حذف آن مشاهدات گرفت. برای تبیین این نکته دو مثال ذکر می‌کنیم.



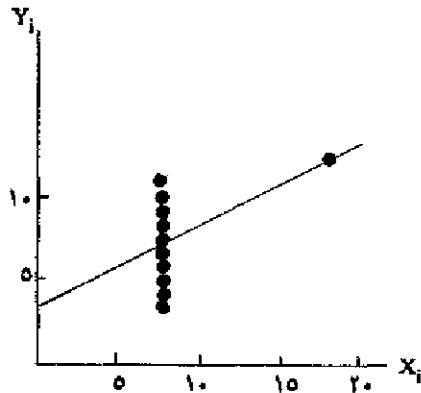
(۱)



(۲)



(۳)



(۴)

نمودار ۳-۳ نمودار هندسی چهارسری مشاهدات

فرض کنید مشاهدات ما به صورت مقطعی و مربوط به تابع تولید در n واحد صنعتی است. همچنین فرض کنید که m مشاهده دورافتاده در نمونه ملاحظه شده است. اگر این m مشاهده دورافتاده به واحدهای تولیدی دولتی و در بخش صنایع نظامی محدود باشد، قطعاً شرایط حاکم بر نظام تولیدی این واحدها، مانند نیروی کار تخصصی، بازار فروش، استفاده از ارز دولتی با واحدهای صنعتی دیگر، متفاوت است. در چنین شرایطی

می‌توان این m مشاهده را حذف کرد و برای $(n - m)$ مشاهده دیگر، یک مدل رگرسیون به صورت مستقل تخمین زد.

مثالی دیگر؛ فرض کنید مشاهدات ما به صورت سری زمانی از تابع تولید در یک بخش صنعتی برای n سال مختلف است. همچنین فرض می‌شود که m مشاهده دورافتاده در این نمونه ملاحظه شده است اگر فرض کنیم که این m مشاهده متعلق به m سالی است که شرایط استثنایی، مانند جنگ، بر نظام تولیدی این بخش صنعتی حاکم بوده است، چه بسا بتوان این چند نمونه دورافتاده را حذف کرد و برای سایر مشاهدات که هماهنگی بهتری با یکدیگر دارند تخمین جداگانه‌ای به دست آورد. البته در فصلهای آینده خواهیم دید که مشاهدات دورافتاده را در اکثر موارد می‌توان با استفاده از «متغیرهای مجازی»^۱ در مدل حفظ کرد، به نحوی که در آنالیز پسماند مشاهدات دورافتاده دیگر دیده نشود. در بسیاری موارد، ظهور مشاهدات دورافتاده بویژه در آنالیز پسماند، نه به علت شرایط استثنایی حاکم بر دوره‌های خاصی از مشاهدات است و نه اینکه عدم هماهنگی (مثلاً) واحدهای تولیدی زمینه‌ساز ایجاد آنها بوده است، بلکه علت اصلی، سادگی مبانی نظری مدل رگرسیون خطی است که نتوانسته است قانونمندی موجود در مشاهدات عینی را به خوبی منعکس کند. در این موارد تخمین مدل رگرسیون در واقع تقریب بسیار ضعیفی از عینیت اقتصادی بوده، در نتیجه مقادیر پسماند در مواردی کاملاً غیرعادی است. راه حل این مسأله قطعاً نمی‌تواند حذف مشاهداتی باشد که عامل ظهور چنین رفتار غیرعادی در پسماند است؛ بلکه برعکس باید این گونه مشاهدات را حفظ کرد و در عوض کوشید تا مدل رگرسیون مفروض را کاملتر کرده و قدرت توضیحی آن را افزایش داد. این کار همان‌گونه که بعد خواهیم دید - با اضافه کردن متغیرهای توضیحی جدید، تبدیل مدل رگرسیون خطی به غیرخطی، تخمین سیستم معادلات به جای یک تک معادله، یا حتی استفاده از روشهای دیگر تخمین، می‌تواند عملی شود. بنابراین نتیجه کلی این است که حذف مشاهدات دورافتاده همیشه باید نامناسبترین و در عین حال

آخرین راه حل تلقی شود.

۶-۳ رابطه های غیرخطی و تبدیل متغیرها

تاکنون بحث ما درباره مدل های رگرسیون خطی ساده بود. در قسمت ۱-۱ به طور خلاصه به این نکته اشاره کردیم که منظور از خطی بودن رگرسیون این است که متغیرها (Y, X) و پارامترها (α) و (β) به صورت خطی وارد مدل می شوند. یک معادله برحسب مثلاً متغیر X خطی است اگر X فقط با توان یک در آن معادله ظاهر شده و در متغیر دیگری ضرب یا بر متغیر دیگری تقسیم نشده باشد؛ بنابراین معادله های

$$Y_i = \alpha + \beta X_i + U_i ,$$

$$Y_i = \alpha + \beta \sqrt{X_i} + U_i ,$$

$$Y_i = \alpha + \beta X_i Z_i + U_i ,$$

که در آن Z_i یک متغیر است، دیگر برحسب X_i خطی نیستند. همچنین یک معادله برحسب مثلاً پارامتر β خطی است، اگر β فقط با توان یک در آن معادله ظاهر شده و در پارامتر دیگری ضرب یا بر پارامتر دیگری تقسیم نشده باشد. بنابراین معادله های

$$Y_i = \alpha + \sqrt{\beta} X_i + U_i ,$$

$$Y_i = \alpha + \ln \beta X_i + U_i ,$$

$$Y_i = \frac{\alpha}{\beta} \ln X_i + U_i ,$$

که در آن \ln لگاریتم بر مبنای e یا لگاریتم طبیعی است، دیگر برحسب β خطی نیستند. ناگفته نماند که معادله آخر، هم برحسب β و هم برحسب X_i غیرخطی است.

از دو حالت غیرخطی برحسب پارامترها و غیرخطی برحسب متغیرها، حالت اول اهمیت فراوان دارد. اگر در یک مدل رگرسیون، پارامترها غیرخطی باشد و هیچ «تبدیلی» برای خطی کردن آنها نتوان یافت، روش حداقل مربعات معمولی دیگر

کاربردی نخواهد داشت. اما معمولاً می‌توان تبدیلهایی یافت که متغیرهای غیرخطی را خطی کند. به همین دلیل می‌توان از واژه «خطی» در عبارت «مدل رگرسیون خطی ساده»، خطی بودن پارامترها را نتیجه گرفت. در این قسمت به بررسی تبدیلهایی می‌پردازیم که می‌تواند یک مدل غیرخطی برحسب متغیرها را به یک مدل خطی تبدیل کند.

۱. مدل‌های خطی - لگاریتمی

مدل رگرسیون زیر را ملاحظه کنید،

$$Y_t = \alpha X_t^\beta e^{U_t} \quad (3-46)$$

که برحسب X_t غیرخطی است. اگر از دو طرف این معادله لگاریتم بگیریم خواهیم داشت

$$\ln Y_t = \ln \alpha + \beta \ln X_t + U_t \quad ,$$

و با فرض $\alpha_0 = \ln \alpha$ داریم

$$\ln Y_t = \alpha_0 + \beta \ln X_t + U_t \quad . \quad (3-47)$$

مدل ۳-۴۷ برحسب پارامترها و همچنین برحسب لگاریتم متغیرهای X_t و Y_t خطی است. به همین دلیل به آن مدل «خطی - لگاریتمی»^۱ می‌گویند. با توجه به اینکه در مدل ۳-۴۷ متغیرهای سمت راست و سمت چپ هر دو برحسب لگاریتم است، به آن مدل «Log - log» یا مدل «لگاریتم مضاعف»^۲ نیز گفته میشود.

برای تخمین پارامترهای α و β این روابط را تعریف می‌کنیم،

$$Y_t^* = \ln Y_t \quad , \quad X_t^* = \ln X_t \quad ,$$

که با جایگزینی در معادله ۳-۴۷ خواهیم داشت

$$Y_t^* = \alpha + \beta X_t^* + U_t$$

اگر U_t تمام فرضهای کلاسیک را داشته باشد، به راحتی می توان α و β را با روش حداقل مربعات معمولی تخمین زد. باید توجه داشت که در تخمین β مسأله ای نخواهیم داشت زیرا از تخمین مدل فوق مستقیماً $\hat{\beta}$ به دست می آید. اما با توجه به معادله ۳-۴۶ پارامتر دیگر مورد نظر ما α است؛ در حالی که از تخمین مدل فوق به $\hat{\alpha}$ می رسیم نه $\hat{\alpha}$ ؛ بنابراین با توجه به فرض $\alpha_0 = \ln \alpha$ ، باید از $\hat{\alpha}_0$ آنتی لگاریتم بگیریم تا $\hat{\alpha}$ به دست آید. می توان نشان داد که تخمین $\hat{\alpha}$ اریب دارد؛ در حالی که $\hat{\beta}$ تخمین نااریبی از β است. با اینکه در این روش $\hat{\alpha}$ اریب دارد، اما چندان باعث تضعیف تبدیل «خطی - لگاریتمی» نیست؛ زیرا مهم این است که توانسته ایم به یک تخمین نااریب از شیب معادله رگرسیون برسیم. می دانیم بین مفاهیم کشش و لگاریتم رابطه بسیاری نزدیکی وجود دارد. در واقع در معادله $Y_t = f(X_t)$ ، کشش نقطه ای Y_t نسبت به X_t (η) برابر است با^۱

۱. برای اثبات، به تعریف کشش مراجعه می کنیم. در معادله $Y_t = f(X_t)$ ، اگر ΔX بتواند ΔY را ایجاد کند آنگاه رابطه

$$\frac{\Delta Y_t}{Y_t} \div \frac{\Delta X_t}{X_t} = \frac{\Delta Y_t}{\Delta X_t} \cdot \frac{X_t}{Y_t}$$

می تواند تغییر نسبی در Y_t را به ازای یک واحد تغییر نسبی از X_t اندازه گیری نماید. کشش Y_t نسبت به X_t در واقع مقدار حدی عبارت فوق است وقتی ΔX_t به سمت صفر میل کند؛

$$\eta = X_t \text{ نسبت به } Y_t \text{ کشش نقطه ای } = \frac{dY_t}{dX_t} \cdot \frac{X_t}{Y_t}$$

حال ثابت می کنیم $\eta = \frac{d(\ln Y_t)}{d(\ln X_t)}$. ابتدا تعریف می کنیم $x_t = \ln X_t$ و نیز $w_t = e^{w_t} = X_t$ و در نتیجه بنابراین $w_t = \ln X_t$

$$\frac{d(\ln Y_t)}{d(\ln X_t)} = \frac{dZ_t}{dX_t} = \frac{dZ_t}{dY_t} \cdot \frac{dY_t}{dX_t} \cdot \frac{dX_t}{dw_t}$$

می دانیم، $\frac{dZ_t}{dY_t} = \frac{1}{Y_t}$ ، $\frac{dZ_t}{dX_t} = \frac{1}{X_t}$ ، $\frac{dX_t}{dw_t} = \frac{1}{d w_t / d X_t} = \frac{1}{1/X_t} = X_t$ ، و در نتیجه با جایگزینی در عبارت فوق خواهیم داشت

$$\frac{d(\ln Y_t)}{d(\ln X_t)} = \frac{dY_t}{dX_t} \cdot \frac{X_t}{Y_t} = \eta$$

$$\eta = \frac{d(\ln Y_t)}{d(\ln X_t)}$$

با ملاحظه معادله ۳-۴۷ می‌توان نتیجه گرفت که مشتق $(\ln Y_t)$ نسبت به $(\ln X_t)$ برابر کشش Y_t نسبت به X_t است که دقیقاً برابر β است:

$$\text{کشش نقطه‌ای } Y_t \text{ نسبت به } X_t = \eta = \frac{d(\ln Y_t)}{d(\ln X_t)} = \beta$$

نتیجه می‌گیریم که اگر مدل ۳-۴۶ را به یک مدل خطی - لگاریتمی تبدیل کنیم، نه تنها پارامترهای α و β به راحتی با روش حداقل مربعات معمولی تخمین زده می‌شوند بلکه $\hat{\beta}$ تخمینی از کشش متغیر درون‌زا نسبت به متغیر برون‌زا نیز هست. با توجه به اینکه این کشش ثابت است، به مدل ۳-۴۷ «مدل کشش ثابت»^۱ نیز می‌گویند.

به نمودار هندسی تغییرات مدل ۳-۴۶ اشاره‌ای می‌کنیم. جمله اختلال (e^{u_t}) را در نظر نگرفته و فقط به رابطه زیر توجه می‌کنیم،

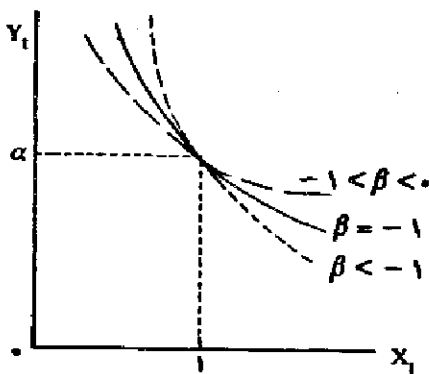
$$Y_t = \alpha X_t^\beta \quad (3-48)$$

اولاً، باید فقط مقادیر مثبت X_t و Y_t را در نظر گرفت؛ چون لگاریتم برای کمیت‌های منفی تعریف نمی‌شود. ثانیاً، بهتر است نمودار تغییرات تابع فوق را یک بار برای $\beta > 0$ و بار دیگر برای $\beta < 0$ بررسی کنیم. ابتدا از معادله ۳-۴۸ مشتق می‌گیریم تا شیب تابع به دست آید،

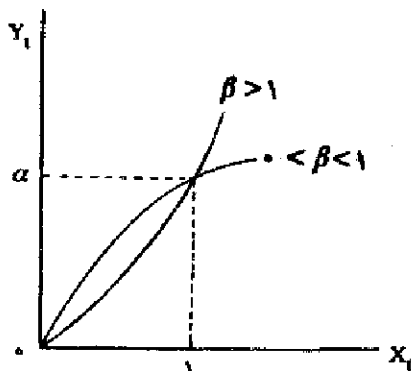
$$\frac{dY_t}{dX_t} = \alpha \beta X_t^{\beta-1} \quad (3-49)$$

اگر β مثبت باشد، با دیدن ۳-۴۹ می‌توان گفت که شیب $(\frac{dY_t}{dX_t})$ همواره مثبت است و به ازای افزایش X_t و میل آن به سمت بی‌نهایت، Y_t نیز زیاد شده و به سمت بی‌نهایت میل می‌کند.

دو حالت را می توان از یکدیگر تمیز داد. اگر $0 < \beta < 1$ ، آنگاه شیب منحنی، به ازای افزایش X_1 ، به طور مرتب کمتر می شود، هر چند همواره مثبت باقی می ماند. اما اگر $\beta > 1$ ، آنگاه به ازای افزایش X_1 ، شیب منحنی به طور مرتب زیاد خواهد شد. این دو حالت را می توان برای $\beta > 0$ در نمودار ۳-۴ مشاهده کرد توجه داریم که در معادله ۳-۴۸ مقدار Y_1 به ازای $X_1 = 1$ برابر α می شود.



نمودار ۳-۵ تابع خطی - لگاریتمی و $\beta < 0$



نمودار ۳-۴ تابع خطی - لگاریتمی و $\beta > 0$

منحنی تغییرات تابع ۳-۴۸ را در حالت $\beta < 0$ بررسی می کنیم. با توجه به معادله ۳-۴۹ می توان گفت که اگر β منفی باشد، شیب همواره منفی خواهد بود. البته می دانیم که α همواره مثبت است؛ زیرا $\ln \alpha$ برای $\alpha > 0$ نمی تواند وجود داشته باشد. مانند نمودار ۳-۴، می توان حالت های مختلفی را برای β در نظر گرفت و نمودار تغییرات تابع ۳-۴۸ را برای آنها ترسیم کرد. در نمودار ۳-۵ سه حالت $-1 < \beta < 0$ ، $\beta = -1$ ، و $\beta < -1$ بررسی شده است.

۲. مدل های نیمه لگاریتمی

در مدل رگرسیون زیر

$$Y_1 = e^{(\alpha + \beta X_1 + U_1)} \quad (3-50)$$

که برحسب X_t غیرخطی است، اگر از دو طرف این مدل لگاریتم بگیریم خواهیم داشت:

$$\ln Y_t = \alpha + \beta X_t + U_t \quad (۳-۵۱)$$

با فرض $Y_t^* = \ln Y_t$ ، معادله ۳-۵۱ را می‌توان به صورت زیر نوشت،

$$Y_t^* = \alpha + \beta X_t + U_t \quad ,$$

که در واقع یک مدل رگرسیون خطی است که می‌توان پارامترهای آن را به راحتی با روش حداقل مربعات معمولی تخمین زد. با توجه به اینکه در مدل ۳-۵۱ فقط متغیر درون‌زا برحسب لگاریتم است، به آن «مدل نیمه‌لگاریتمی»^۱ می‌گویند. مدل‌های رگرسیون از نوع

$$Y_t = \alpha e^{(\beta X_t + U_t)} \quad (۳-۵۲)$$

نیز وقتی به مدل‌های نیمه‌لگاریتمی تبدیل شوند برحسب X_t خطی خواهند بود. اگر از دو طرف معادله ۳-۵۲ لگاریتم بگیریم خواهیم داشت

$$\ln Y_t = \ln \alpha + \beta X_t + U_t \quad ,$$

که با تعریف $\alpha_0 = \ln \alpha$ به صورت زیر خواهد بود،

$$\ln Y_t = \alpha_0 + \beta X_t + U_t \quad ,$$

که دقیقاً با معادله ۳-۵۱ ساختاری مشابه دارد. با فرض $Y_t^* = \ln Y_t$ ، می‌توان مدل فوق را به صورت زیر نوشت،

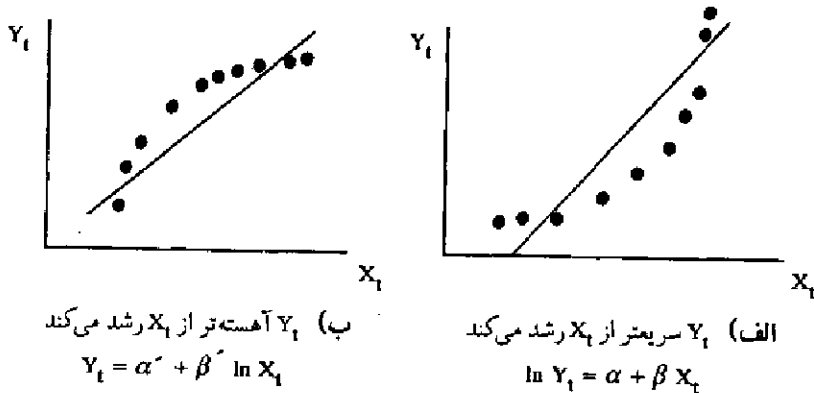
$$Y_t^* = \alpha_0 + \beta X_t + U_t \quad .$$

ناگفته نماند که مدل‌هایی از نوع ۳-۵۱ در تخمین توابع سرمایه‌های انسانی کاربردهای فراوانی دارند، که معمولاً X_t ، سال‌های اشتغال به تحصیل و Y_t ، درآمد فرد است. برای آشنایی بیشتر با توابع نیمه‌لگاریتمی باید در نظر داشت که اگر Y_t سریعتر از

X_t افزایش یابد، از مدل ۳-۵۱ و در حالتی که Y_t آهسته تر از X_t رشد کند، می توان از مدل زیر استفاده کرد

$$Y_t = \alpha' + \beta' \ln X_t + U_t \quad (۳-۵۳)$$

به مدل ۳-۵۳ نیز یک مدل نیمه لگاریتمی می گویند؛ زیرا فقط X_t بر حسب لگاریتم بیان شده است. در نمودار ۳-۶ الف حالتی نشان داده شده است که مدل ۳-۵۱ کاربرد دارد و در نمودار ۳-۶ ب کاربرد مدل ۳-۵۳ منعکس است.



نمودار ۳-۶ مدلهای نیمه لگاریتمی

مدل نیمه لگاریتمی ۳-۵۱ را در نظر می گیریم. برای اینکه بتوانیم تفسیر مناسبی از

β ارائه کنیم، از $\ln Y_t$ نسبت به X_t مشتق می گیریم،

$$\frac{d \ln Y_t}{d X_t} = \beta$$

رابطه فوق را می توان به صورت زیر نوشت،

$$\beta = \frac{1}{Y_t} \cdot \frac{d Y_t}{d X_t} = \frac{d Y_t}{Y_t} \cdot \frac{1}{d X_t}$$

از طرف دیگر می دانیم $\frac{d Y_t}{Y_t}$ با تغییر نسبی در Y_t برابر است و $d X_t$ نیز چیزی جز تغییر

مطلق در X_t نیست به این ترتیب

$$\beta = \frac{\text{تغییر نسبی در } Y_t}{\text{تغییر مطلق در } X_t} .$$

در مواردی که به ازای یک مقدار معین از تغییر مطلق در X_t ، متغیر Y_t به صورت یک درصد ثابتی تغییر می‌کند، می‌توان نتیجه گرفت که بهترین مدل رگرسیون برای آنها، مدل نیمه‌لگاریتمی از این نوع است

$$\ln Y_t = \alpha + \beta X_t + U_t .$$

این گونه مدلها را «مدل رشد ثابت»^۱ نیز می‌گویند. می‌دانیم مدل‌های رشد ثابت کاربرد فراوانی در اندازه‌گیری نرخ رشد روند تغییرات متغیرهایی چون قیمت، بیکاری، صادرات، واردات و مانند آن دارد.

با ترتیبی کاملاً مشابه می‌توان β' را در مدل ۳-۵۳ به صورت زیر تفسیر کرد،

$$\beta' = \frac{\text{تغییر مطلق در } Y_t}{\text{تغییر نسبی در } X_t} .$$

بنابراین در مواردی که مشاهده می‌شود به ازای یک مقدار معین از تغییر نسبی در X_t ، متغیر Y_t به صورت مطلق تغییر می‌کند، نتیجه می‌گیریم که مدل نیمه‌لگاریتمی به صورت زیر

$$Y_t = \alpha' + \beta' \ln X_t + U_t .$$

بهترین مدل برای تخمین رابطه بین آنها خواهد بود.

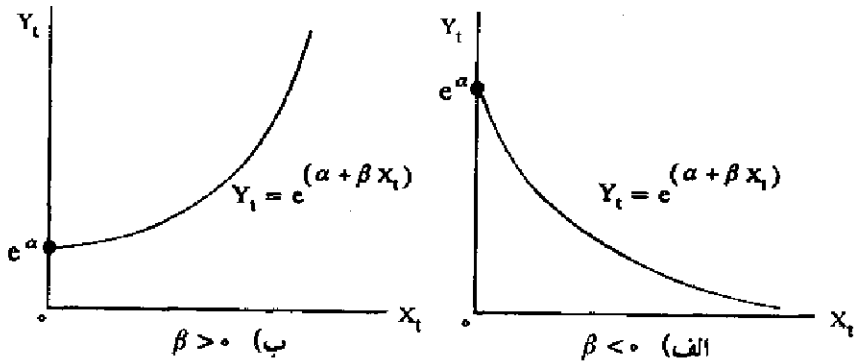
برای بررسی نمودار هندسی مدل‌های نیمه‌لگاریتمی، مدل ۳-۵۱ را یک بار دیگر ملاحظه کنید،

$$\ln Y_t = \alpha + \beta X_t + U_t .$$

می‌دانیم این تابع فقط برای مقادیر مثبت Y_t تعریف می‌شود، زیرا کمیتهای منفی، لگاریتم ندارند. اگر جمله اختلال را نادیده بگیریم، می‌توان مدل فوق را به صورت زیر نوشت،

$$Y_t = e^{(\alpha + \beta X_t)}$$

که در واقع همان معادله ۳-۵۰ و بدون جمله اختلال است. این معادله را می توان برای دو حالت $\beta < 0$ و $\beta > 0$ در نمودار ۳-۷ ملاحظه کرد.



نمودار ۳-۷. مدل های نیمه لگاریتمی برای $\beta > 0$ و $\beta < 0$

به ازای $X = 0$ مقدار Y_t برابر e^α خواهد بود؛ بنابراین عرض از مبدأ در هر دو نمودار برابر e^α است. مشاهده می شود که علامت شیب تابع، از علامت β تبعیت می کند. اگر β منفی باشد شیب منفی و تابع نزولی خواهد بود، که در نمودار ۳-۷ الف نشان داده شده است. برای مواردی که β مثبت است، همچنانکه در نمودار ۳-۷ ب ملاحظه می کنیم منحنی تغییرات تابع ۳-۵۰ صعودی است.

می توان حالت خاصی را از معادله ۳-۵۰ در نظر گرفت که X_t متغیر زمان باشد. در

این صورت خواهیم داشت

$$Y_t = e^{\alpha + \beta t} \quad (3-54)$$

معادله ۳-۵۴ نشان می دهد که Y_t به ازای $\beta > 0$ ، یک نرخ نسبی رشد ثابت داشته و به ازای $\beta < 0$ از یک نرخ رکود نسبی و ثابت برخوردار است.

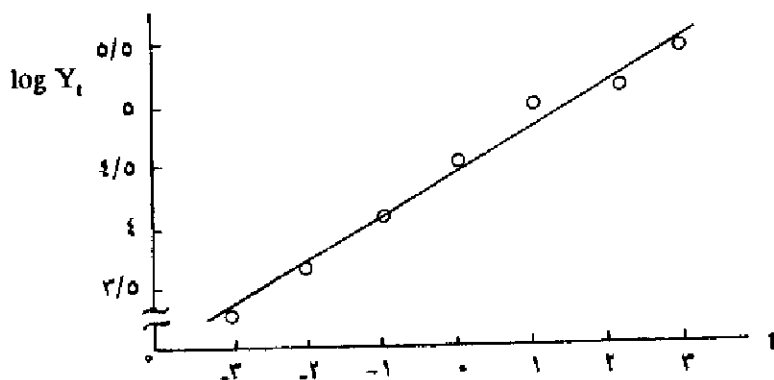
مثال ۳-۵ آمار تولیدات یک بخش صنعتی در ۷ دهه گذشته در جدول زیر ارائه شده

است. به نظر شما برای این مشاهدات چه نوع مدل رگرسیون می‌توان پیشنهاد کرد. آیا می‌توان گفت که این صفت از نرخ رشد ثابتی برخوردار بوده است.

جدول ۳.۳

دهه	Y_t : متوسط تولید	$\log Y_t$	$X=t$
۱	۱۸۳۷	۳/۲۶۴۱	-۳
۲	۴۸۶۸	۳/۶۸۷۳	-۲
۳	۱۲۴۱۱	۴/۰۹۳۷	-۱
۴	۳۲۶۱۷	۴/۵۱۳۵	۰
۵	۸۲۷۷۰	۴/۹۱۷۹	۱
۶	۱۴۸۴۵۷	۵/۱۷۱۸	۲
۷	۳۲۲۹۵۸	۵/۵۰۹۲	۳

لگاریتم متوسط تولید را در دهه‌های مختلف حساب کرده و آن را در یک نمودار و در مقابل t نشان دهیم. به سهولت ملاحظه می‌شود که مشاهدات موجود در نمودار ۳.۸



نمودار ۳.۸

یک رابطه تقریباً خطی را نشان می‌دهد؛ بنابراین مدلی که می‌توان پیشنهاد کرد عبارت است از

$$\ln Y_t = \alpha + \beta t + U_t$$

که دقیقاً همان مدل ۳-۵۴ است. با استفاده از روش حداقل مربعات معمولی می‌توان به تخمین پارامترهای α و β رسید. خواهیم داشت

$$\ln \hat{Y}_t = 4/4510 + 0/3760 t .$$

توجه داریم که در اندازه‌گیری متغیر زمان، دوره دهه چهارم را مینا گرفته و مقدار آن را صفر قرار داده‌ایم.

۳. مدل‌های وارون

به مدل رگرسیون زیر

$$Y_t = \alpha + \beta \left(\frac{1}{X_t} \right) + U_t , \quad (3-55)$$

که در آن Y_t برحسب معکوس X_t بیان شده است یک «مدل وارون»^۱ می‌گویند. برای بررسی منحنی تغییرات این تابع، ابتدا جمله اختلال مدل (U_t) را در معادله ۳-۵۵ در نظر نگرفته، از آن مشتق می‌گیریم:

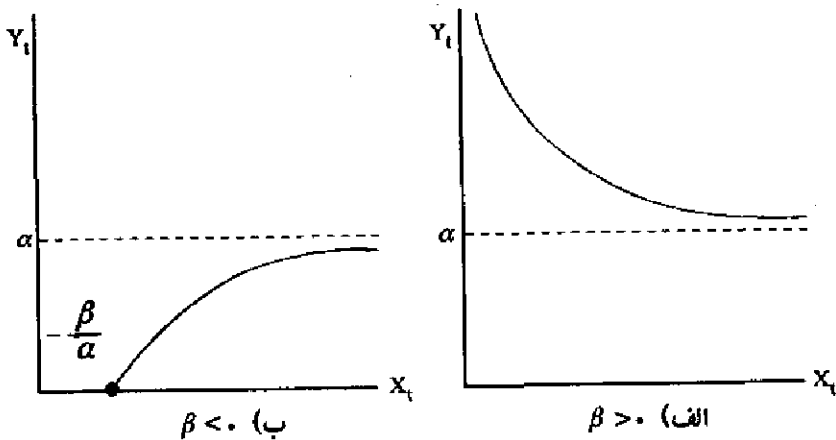
$$\frac{d Y_t}{d X_t} = - \frac{\beta}{X_t^2}$$

چون علامت X_t^2 همواره مثبت است؛ بنابراین علامت مشتق دقیقاً برعکس علامت β است. حال منحنی تغییرات تابع Y_t را در دو حالت بررسی خواهیم کرد:

حالت اول $\beta > 0$

هرگاه β مثبت باشد، علامت مشتق منفی و منحنی تغییرات Y_t نزولی است. اگر X_t به سمت بی‌نهایت میل کند، مقدار $\frac{1}{X_t}$ به صفر میل کرده، در نتیجه مقدار Y_t در معادله ۳-۵۵ به سمت α میل خواهد کرد؛ بنابراین خط $Y_t = \alpha$ در واقع مجانب منحنی تغییرات

Y_t است. نمودار ۳-۹ الف، منحنی تغییرات تابع را به ازای $\beta > 0$ نشان می‌دهد.



نمودار ۳-۹. مدل‌های معکوس

حالت دوم $\beta < 0$

به ازای مقادیر منفی β ، مقدار $\frac{dY_t}{dX_t}$ مثبت شده و منحنی تغییرات Y_t یک سیر صعودی را نشان می‌دهد. در این حالت نیز اگر $X_t \rightarrow \infty$ ، $Y_t = \alpha$ ، بجانب منحنی تغییرات Y_t خواهد بود. به ازای $Y_t = 0$ خواهیم داشت

$$\alpha + \frac{1}{\beta} X_t = 0$$

در نتیجه

$$X_t = \frac{-\beta}{\alpha}$$

یعنی منحنی تغییرات Y_t محور X را در نقطه $\frac{-\beta}{\alpha}$ قطع می‌کند.

مدل‌های رگرسیون وارون کاربردهای بسیاری بویژه در اقتصادسنجی خرد دارد؛ برای مثال، می‌توان گفت که نمودار ۳-۹ الف در واقع مبین شکل عمومی یک «منحنی فیلیس»^۱ است. می‌دانیم در بسیاری موارد می‌توان نرخ تغییر سطح دستمزدها

را با رگرسیون W_i بر U_i به صورت زیر تخمین زد،

$$W_i = \alpha + \beta \frac{1}{U_i} + \varepsilon_i,$$

که در آن W_i و U_i به ترتیب نرخ تغییر سطح دستمزدها و نرخ بیکاری است. مشاهده می‌شود که با تخمین α ، می‌توان حد پایین تغییر در نرخ دستمزدها را به دست آورد، به گونه‌ای که با افزایش بیشتر نرخ بیکاری، دستمزدها کاهش نمی‌یابد.

مدلهای تخمین تابع مخارج خانوار برای کالاها و خدمات خاص، می‌تواند در بسیاری موارد به صورت یک مدل وارون طراحی شود که در آن β منفی است. فرض کنید با استفاده از داده‌های مقطعی روی یک نمونه n تایی از خانوارها می‌خواهیم تابع مخارج خانوار را برای کالا یا خدمت معینی تخمین بزنیم،

$$Y_i = \alpha + \beta \frac{1}{X_i} + U_i,$$

که در آن Y_i مخارج خانوار i برای کالای مفروض و X_i درآمد کل یا مخارج کل این خانوار است. به نمودار ۹-۳ مراجعه می‌کنیم. قبل از اینکه درآمد یا مخارج کل خانوار به نقطه $X_i = -\frac{\beta}{\alpha}$ برسد، هیچ تقاضایی برای این کالا یا خدمت مفروض وجود ندارد. ولی از نقطه $-\frac{\beta}{\alpha}$ به بعد، به ازای افزایش درآمد، مخارج خانوار بشدت افزایش می‌یابد، اما بزودی به نقطه تعادل $Y_i = \alpha$ می‌رسد، به گونه‌ای که با افزایش بیشتر درآمد، خانوار حاضر نیست برای آن کالا یا خدمت هزینه بیشتر بپردازد. بسیاری از مواد غذایی لوکس یا پوشاک وارداتی و مانند آن خصوصیت فوق را دارد. لازم است یادآوری کنیم که هنگامی می‌توان از مدل وارون برای تخمین این گونه توابع مصرف خانوار استفاده کرد که $\hat{\beta}$ منفی و $\hat{\alpha}$ مثبت باشد.

حال که با مفهوم توابع وارون و کاربرد آن در تخمین مدل‌های رگرسیون آشنا شدیم به چگونگی تخمین این گونه مدل‌ها می‌پردازیم. می‌دانیم مدل ۳-۵۵ برحسب X_i غیرخطی است. برای اینکه بتوان α و β را با روش حداقل مربعات معمولی تخمین زد، باید متغیرها و پارامترها همه خطی باشند. با توجه به اینکه در این مدل α و β به صورت

خطی وارد شده‌اند فقط X_t را خطی می‌کنیم. با استفاده از تعریف

$$X_t^* = \frac{1}{X_t} ,$$

و جایگزینی آن در مدل ۳-۵۵ خواهیم داشت

$$Y_t = \alpha + \beta X_t^* + U_t ,$$

که یک مدل خطی برحسب متغیرها و پارامترهاست و می‌توان براحتی آن را با روش حداقل مربعات معمولی تخمین زد.

۴. مدل‌های وارون لگاریتمی

به مدل رگرسیون زیر

$$\ln Y_t = \alpha - \beta \left(\frac{1}{X_t} \right) + U_t \quad (3-56)$$

که در آن Y_t به صورت لگاریتمی و برحسب وارون X_t بیان شده است یک «مدل وارون لگاریتمی»^۱ می‌گویند. ساختار ریاضی این مدل با مدل ۳-۵۵ تشابه دارد و تفاوت آنها فقط این است که در ۳-۵۶ متغیر Y_t به صورت لگاریتمی وارد شده است. می‌توان گفت که مدل ۳-۵۶ در واقع عبارت است از مدل

$$Y_t = e^{\left(\alpha - \beta \frac{1}{X_t} + U_t \right)} , \quad (3-57)$$

زیرا کافی است از معادله ۳-۵۷ لگاریتم گرفته تا معادله ۳-۵۶ به دست آید. برای بررسی منحنی تغییرات Y_t ابتدا U_t را نادیده گرفته خواهیم داشت

$$Y_t = e^{\left(\alpha - \frac{\beta}{X_t} \right)} . \quad (3-58)$$

مقدار Y_t به ازای $X_t = 0$ تعریف نمی‌شود زیرا Y_t در معادله ۳-۵۸ تابعی از $\frac{1}{X_t}$ است. اما

اگر X_1 به سمت صفر میل کند، مقدار Y_1 نیز به سمت صفر میل خواهد کرد؛ بنابراین با صرف نظر کردن از دقت‌های ریاضی، می‌توان به طور خلاصه گفت که $Y_1 = 0$. در نتیجه با این قرارداد، می‌توانیم منحنی تغییرات Y_1 را از مبدأ مختصات شروع کنیم؛ چون Y_1 همواره مثبت است؛ منحنی تغییرات آن از مبدأ مختصات و به طور پیوسته به سمت راست ادامه خواهد داشت.

برای بررسی شیب تغییرات این منحنی، باید از معادله ۳-۵۸ مشتق بگیریم، خواهیم داشت

$$\frac{d Y_1}{d X_1} = \left(\frac{\beta}{X_1^2} \right) e^{\left(\alpha - \frac{\beta}{X_1} \right)}, \quad (3-59)$$

یعنی علامت شیب دقیقاً از علامت β تبعیت می‌کند. اگر β مثبت باشد، شیب تابع مثبت بوده و منحنی تغییرات Y_1 صعودی است و به ازای مقادیر منفی β ، علامت $\frac{d Y_1}{d X_1}$ منفی خواهد شد که در آن صورت، منحنی تغییرات Y_1 نزولی است. حال به بررسی مشتق دوم می‌پردازیم. خواهیم داشت

$$\frac{d^2 Y_1}{d X_1^2} = - \frac{2 \beta X_1}{X_1^4} e^{\left(\alpha - \frac{\beta}{X_1} \right)} + \frac{\beta}{X_1^2} \left[\frac{\beta}{X_1^2} e^{\left(\alpha - \frac{\beta}{X_1} \right)} \right],$$

یا

$$\frac{d^2 Y_1}{d X_1^2} = \left(\frac{\beta^2}{X_1^4} - \frac{2 \beta}{X_1^3} \right) e^{\left(\alpha - \frac{\beta}{X_1} \right)}$$

اگر مشتق دوم را مساوی صفر قرار دهیم، نقطه عطف منحنی تغییرات Y_1 در نقطه‌ای به طول $X_1 = \frac{\beta}{\gamma}$ به دست می‌آید. در سمت چپ این نقطه، یعنی $X_1 < \frac{\beta}{\gamma}$ ، شیب منحنی به ازای افزایش X_1 ، زیاد می‌شود و در سمت راست آن، به ازای افزایش X_1 ، شیب منحنی کاهش خواهد یافت. برای به دست آوردن عرض نقطه عطف کافی است $X_1 = \frac{\beta}{\gamma}$ را در معادله ۳-۵۸ قرار دهیم،

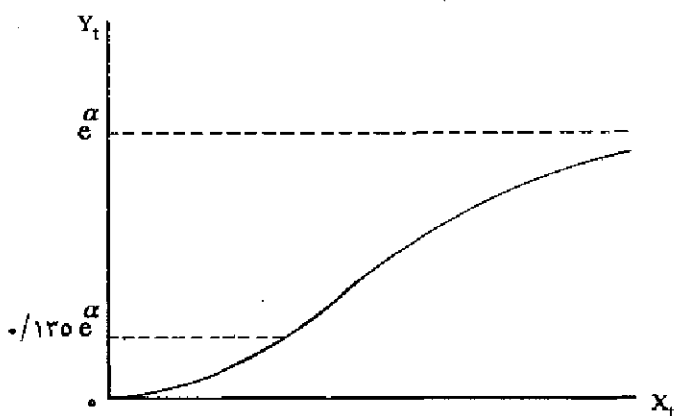
$$Y_1 = e^{\left(\alpha - \frac{\beta}{\beta/\gamma} \right)}$$

$$= e^{\alpha - \beta} = e^{\alpha} \cdot e^{-\beta}$$

که با توجه به $e = 2/718$ خواهیم داشت

$$Y_t = 0/135 e^{\alpha}$$

از طرف دیگر می‌دانیم به ازای افزایش X_t و میل آن به سمت بی‌نهایت، مقدار Y_t به سمت e^{α} میل می‌کند؛ بنابراین خط $Y_t = e^{\alpha}$ در واقع مجانب منحنی تغییرات Y_t است. مشاهده می‌شود که عرض نقطهٔ عطف، یعنی $0/135 e^{\alpha}$ ، حدود ۱۳ درصد مقدار حدی Y_t ، یعنی e^{α} است. با این اطلاعات، می‌توانیم نمودار تغییرات منحنی Y_t را در معادلهٔ ۳-۵۸ در نمودار ۳-۱۰ نشان دهیم.



نمودار ۳-۱۰ مدل وارون لگاریتمی

مشاهده می‌کنیم که نرخ رشد Y_t به موازات افزایش X_t بشدت کاهش می‌یابد این خصوصیت مدل‌های وارون لگاریتمی را باید در کاربرد آنها همواره در نظر داشت. از نظر ریاضی نیز می‌توان این خصوصیت را بسادگی نشان داد. معادلهٔ ۳-۹۵ را یک بار دیگر می‌نویسیم،

$$\frac{dY_t}{dX_t} = \left(\frac{\beta}{X_t}\right) e^{\left(\alpha - \frac{\beta}{X_t}\right)}$$

دو طرف این معادله را بر معادله ۳-۵۸ تقسیم می‌کنیم. جمله $e^{(\alpha - \frac{\beta}{X_t})}$ حذف شده خواهیم داشت

$$\frac{1}{Y_t} \cdot \frac{dY_t}{dX_t} = \frac{\beta}{X_t^2}$$

یعنی نرخ رشد Y_t تابعی از معکوس X_t است و با افزایش X_t شدت کم خواهد شد. بعد از آشنایی با خصوصیات ریاضی و نمودار تغییرات منحنی مدل‌های وارون لگاریتمی به بررسی چگونگی تخمین آنها می‌پردازیم. معادله ۳-۵۶ را یک بار دیگر می‌نویسیم،

$$\ln Y_t = \alpha - \beta \left(\frac{1}{X_t} \right) + U_t$$

با تبدیل متغیرها به

$$\ln Y_t = Y_t^*, \quad \frac{1}{X_t} = X_t^*$$

خواهیم داشت

$$Y_t^* = \alpha - \beta X_t^* + U_t$$

که بر حسب متغیرها و پارامترها خطی است؛ بنابراین با روش حداقل مربعات معمولی براحتی می‌توان پارامترهای آن را تخمین زد.

در پایان، یک بار دیگر نمودار ۳-۱۰ را ملاحظه کنیم. منحنیهای دیگری هم وجود دارند که شبیه نمودار تغییرات منحنی مدل‌های وارون لگاریتمی هستند؛ اما:

اولاً، علاوه بر یک مجانب در بالا، یک مجانب دیگر نیز دارند که همان محور X_t است؛ به عبارت دیگر، معادله‌های $Y_t = k$ و $Y_t = 0$ مجانبهای آن است که k عدد ثابت و معینی است.

ثانیاً، منحنی تغییرات آنها در بین دو مجانب تقریباً قرینگی دارد.

به معادله‌های این گونه منحنیها «مدلهای لجیستیک»^۱ می‌گویند. این مدلها در تخمین

مدلهای رگرسیون مربوط به روند رشد متغیرها کاربرد فراوانی دارند و در جلد دوم این کتاب با نام «مدلهای عکس‌العمل کیفی»^۱ بررسی خواهند شد.

۵. مدل‌های غیرخطی برحسب پارامترها

تا به حال بحث ما عمدتاً دربارهٔ آن دسته از مدل‌های رگرسیون بود که متغیرهای آن به صورت غیرخطی وارد مدل شده‌اند. در این قسمت به بررسی حالت خاصی می‌پردازیم که پارامترها غیرخطی هستند. مدل رگرسیون زیر را ملاحظه کنید،

$$Y_t = \alpha + \frac{\beta}{X_t + \gamma} + U_t \quad (3-60)$$

بدیهی است به علت غیرخطی بودن γ ، نمی‌توان پارامترهای مدل را با روش حداقل مربعات معمولی تخمین زد. در اینجا به روشی اشاره می‌کنیم که حالت عمومی دارد و می‌تواند برای همهٔ مدل‌هایی به کار رود که با ثابت نگهداشتن یک پارامتر، مدل غیرخطی به یک مدل خطی برحسب پارامترها تبدیل می‌شود؛ برای مثال، اگر در مدل ۳-۶۰ مقدار γ را ثابت بگیریم، مدل مزبور برحسب پارامترها خطی شده و به راحتی می‌توان آن را، یا فرض $X_t^* = \frac{1}{X_t + \gamma}$ ، به کمک روش حداقل مربعات معمولی تخمین زد.

می‌دانیم مبنای روش حداقل مربعات معمولی در تخمین مدل ۳-۶۰ این است که مجموع مربعات پسماند را حداقل کند؛ بنابراین باید $\hat{\alpha}$ ، $\hat{\beta}$ و $\hat{\gamma}$ را چنان به دست آوریم که عبارت زیر حداقل شود،

$$\sum e_t^2 = \sum \left(Y_t - \hat{\alpha} - \frac{\hat{\beta}}{X_t + \hat{\gamma}} \right)^2$$

برای هر مقدار ممکن از $\hat{\gamma}$ ، متغیر Z_t را به صورت زیر تعریف می‌کنیم،

$$Z_t = \frac{1}{X_t + \hat{\gamma}} \quad (3-61)$$

با جایگزینی معادله ۳-۶۱ در مجموع مربعات پسماند خواهیم داشت

$$\sum e_i^2 = \sum (Y_i - \hat{\alpha} - \hat{\beta} Z_i)^2 \quad (3-62)$$

از معادله ۳-۶۲ بر حسب $\hat{\alpha}$ و $\hat{\beta}$ مشتق گرفته و مقادیری از $\hat{\alpha}$ و $\hat{\beta}$ که مشتق را صفر می‌کند، تخمینهای روش حداقل مربعات معمولی از α و β خواهد بود.

نکته مهم این است که مقدار Z_i در معادله ۳-۶۱ باید به ازای تمام مقادیر ممکن \hat{Y} تعریف شود. اگر m مقدار برای \hat{Y} فرض شود، در آن صورت m مقدار مختلف برای Z_i خواهیم داشت. مشاهده می‌شود که به ازای m مقدار Z_i می‌توان m معادله از ۳-۶۲ به دست آورد؛ بنابراین m مقدار مختلف $\sum e_i^2$ خواهیم داشت. از بین m مقدار $\sum e_i^2$ ، کوچکترین آن را انتخاب کرده و Z_i آن را مشخص می‌سازیم. آن مقدار از \hat{Y} که Z_i را تعیین کرده است - تخمین زنده روش حداقل مربعات معمولی از γ خواهد بود. مقادیر $\hat{\alpha}$ و $\hat{\beta}$ متناظر با کمترین مقدار $\sum e_i^2$ نیز تخمین زنده‌های روش حداقل مربعات معمولی α و β است.

بدیهی است لازمه اصلی کاربرد روش فوق این است که با استفاده از نظریه‌ها و مشاهدات اقتصادی، محدوده تغییرات γ را از قبل بدانیم تا بتوانیم مقادیر Z_i و $\sum e_i^2$ را در قلمروی محدود محاسبه کنیم. روشهای دیگری نیز برای خطی کردن پارامترهای غیرخطی وجود دارد که از برنامه این کتاب خارج است.

مسائل فصل سوم

۳-۱ مسأله ۲-۱ را یک بار دیگر ملاحظه کنید. مدل و مشاهدات زیر مفروض است،

$$Y_i = \alpha + \beta X_i + U_i .$$

$$X_i: \quad 2 \quad 3 \quad 1 \quad 5 \quad 9$$

$$Y_i: \quad 4 \quad 7 \quad 3 \quad 9 \quad 17$$

می‌دانیم نتایج زیر به دست آمده است،

$$\hat{Y}_i = 1 + 1/70 X_i ,$$

$$\bar{X} = 4 ,$$

$$\hat{\sigma}_U^2 = s^2 = \frac{\sum e_i^2}{n-2} = \frac{1/5}{3} = 0/5 ,$$

$$\sum x_i^2 = 40 .$$

به ازای $X_i = 10$ ، مطلوب است

۱. پیش‌بینی نقطه‌ای برای Y_i ، یعنی \hat{Y}_i ؛

۲. فاصله اطمینان ۹۵ درصد برای Y_i ؛

۳. فاصله اطمینان ۹۵ درصد برای میانگین Y_i ، یعنی $E(Y_i)$.

۳-۲ مدل رگرسیون زیر مفروض است

$$Y_i = \alpha + \beta X_i + \varepsilon_i ,$$

که در آن ε_i جمله اختلال مدل بوده و توزیع نرمال با میانگین صفر و واریانس σ^2 دارد. این کمیته‌ها را از نمونه‌ای با ۲۰ مشاهده محاسبه کرده‌ایم،

$$\sum Y_i = 21/9 \quad , \quad \sum (Y_i - \bar{Y})^2 = 87/9 \quad , \quad \sum (X_i) = 187/2 \quad ,$$

$$\sum (X_i - \bar{X})(Y_i - \bar{Y}) = 107/4 \quad , \quad \sum (X_i - \bar{X})^2 = 215/4 .$$

الف) پارامترهای α و β را تخمین بزنید.

ب) واریانس $\hat{\alpha}$ و $\hat{\beta}$ را تخمین بزنید.

ج) مقدار میانگین شرطی Y_t را به شرط $X_t = 10$ تخمین بزنید،

$$\hat{E}(Y_t | X_t = 10)$$

د) برای میانگین شرطی Y_t که در بند (ج) تعریف شده است - یک فاصله اطمینان ۹۵ درصد بسازید.

۳-۳ مدل رگرسیون زیر را که جمله ثابت ندارد، مشاهده کنید،

$$Y_t = \beta X_t + U_t ,$$

که در آن $t = 1, 2, \dots, n$ و تمام فرضهای کلاسیک را شامل است. به ازای مقدار X_t ، مقدار Y_t تعریف می شود. می خواهیم Y_t^* را پیش بینی کنیم. پیش بینی کننده \hat{Y}_t^* را چنان تعیین کنید که

$$E(\hat{Y}_t^*) = E(Y_t^*) .$$

۳-۴ مدل رگرسیون زیر مفروض است،

$$Y_t = \beta X_t + U_t ,$$

که در آن $t = 1, 2, 3, \dots, T, T+1, T+2, \dots, T+F$. تمام فرضهای کلاسیک در مورد U_t صادق است. برای هر دو متغیر X_t و Y_t مشاهدات $t = 1, 2, \dots, T$ موجود است، اما مشاهدات $t = T+1, T+2, \dots, T+F$ را فقط برای X_t داریم. پیش بینی کننده ما از Y_{T+m} عبارت است از

$$\hat{Y}_{T+m} = \hat{\beta} X_{T+m} ,$$

که در آن $\hat{\beta}$ تخمین روش حداقل مربعات معمولی از β است که با استفاده از مشاهدات $t = 1, 2, \dots, T$ به دست آمده است. نشان دهید که واریانس پیش بینی کننده در حول مقدار واقعی آن برابر است با واریانس پیش بینی کننده در حول مقدار واقعی آن برابر

است با واریانس پیش‌بینی کننده در حول میانگین آن، بعلاوه واریانس جمله اختلال:

$$E(\hat{Y}_{T+F} - Y_{T+F})^2 = E[\hat{Y}_{T+F} - E(Y_{T+F})]^2 + \sigma_u^2.$$

۳-۵ نشان دهید که اگر در مدل رگرسیون ساده، یک تبدیل خطی روی هر دو متغیر برون‌زا و درون‌زا انجام شود؛

$$Y_i^* = p_1 + q_1 Y_i,$$

$$X_i^* = p_2 + q_2 X_i,$$

آنگاه ضریب تعیین، (r^2) بدون تغییر باقی می‌ماند.

۳-۶ نشان دهید که فقط هنگامی تخمین شیب رگرسیون Y_i بر X_i برابر معکوس تخمین شیب رگرسیون X_i بر Y_i است که داشته باشیم $r^2 = 1$.

۳-۷ در مدل رگرسیون $Y_i = \alpha + \beta X_i + U_i$ ، با استفاده از یک نمونه ۲۰۰ تایی مشاهده در مورد X_i و Y_i ، کمیت‌های زیر را محاسبه کرده‌ایم،

$$\sum X_i = 11/34, \quad \sum Y_i = 20/72,$$

$$\sum X_i^2 = 12/16, \quad \sum Y_i^2 = 84/96, \quad \sum X_i Y_i = 22/13.$$

α و β را تخمین زده و آنها را با $\hat{\alpha}$ و $\hat{\beta}$ در مدل رگرسیون زیر مقایسه کنید،

$$X_i = \alpha' + \beta' Y_i + V_i.$$

۳-۸ مدل رگرسیون $Y_i = \beta X_i + U_i$ مفروض است

β ، ضریب تعیین (r^2) ، و تخمین واریانس جمله اختلال $(\hat{\sigma}_u^2)$ را در نظر می‌گیریم. مدل رگرسیون معکوس را بصورت $X_i = \beta' Y_i + V_i$ می‌نویسیم. $\hat{\beta}'$ ، r^2 و $\hat{\sigma}_{V_i}^2$ را برای این مدل به دست آورده، با مقادیر مشابه در رگرسیون مستقیم مقایسه کنید.

۳-۹ آیا مدل‌های رگرسیون زیر را می‌توان به مدل‌هایی تبدیل کرد که برحسب پارامترها خطی باشد؟

$$Y_t = \alpha e^{(\alpha + \beta X_t + U_t)}. \quad (\text{ج}) \quad Y_t = \alpha e^{\beta X_t} U_t. \quad (\text{الف})$$

$$Y_t = \frac{\alpha}{(\beta - X_t)} + U_t. \quad (\text{د}) \quad Y_t = \alpha e^{-\beta X_t} + U_t. \quad (\text{ب})$$

۳-۱۰ فرض کنید آمار X_t و Y_t را داریم. توضیح دهید که چگونه می‌توان پارامتر مدلهای زیر را با روش حداقل مربعات معمولی تخمین زد:

$$Y_t = \frac{e^{\alpha + \beta X_t}}{1 + e^{\alpha + \beta X_t}}. \quad (\text{ه}) \quad Y_t = \alpha X_t^\beta. \quad (\text{الف})$$

$$Y_t = \alpha + \beta \sqrt{X_t}. \quad (\text{و}) \quad Y_t = \alpha e^{\beta X_t}. \quad (\text{ب})$$

$$Y_t = \alpha + e^{\beta X_t}. \quad (\text{ز}) \quad Y_t = \alpha + \beta \ln X_t. \quad (\text{ج})$$

$$Y_t = \alpha + \frac{\beta}{\gamma X_t - \gamma}. \quad (\text{ح}) \quad Y_t = \frac{X_t}{\alpha X_t - \beta}. \quad (\text{د})$$

درباره چگونگی ورود جمله اختلال یعنی U_t ، در هر یک از مدلهای فوق بحث کنید.

حل مسائل فصل سوم

۳-۱ الف) با استفاده از تخمین مدل مفروض داریم،

$$\hat{Y}_f = 1 + 1/70 X_f .$$

به ازای $X_f = 10$ ، پیش‌بینی نقطه‌ای عبارت است از

$$\hat{Y}_f = 1 + 1/70 (10) = 18/5 .$$

ب) نامساوی ۳-۱۰ را می‌نویسیم،

$$\hat{Y}_f - t_{\alpha/2} \hat{\sigma}_v \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum X_i^2}} < Y_f < \hat{Y}_f + t_{\alpha/2} \hat{\sigma}_v \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum X_i^2}}$$

مقدار به دست آمده از جدول t با $n - 2 = 5 - 2 = 3$ درجه آزادی و در سطح معنی‌دار $\alpha_p = 0.025$ برابر است با $3/182$. با استفاده از کمیت‌های داده شده، خواهیم داشت

$$18/5 - 3/182 \sqrt{0/5} \sqrt{1 + \frac{1}{5} + \frac{(10 - 4)^2}{40}} < Y_f < 18/5 + 3/182 \sqrt{0/5} \sqrt{1 + \frac{1}{5} + \frac{(10 - 4)^2}{40}}$$

در نتیجه فاصله اطمینان ۹۵ درصد برای Y_f عبارت خواهد بود از

$$18/5 - 3/26 < Y_f < 18/5 + 3/26 ,$$

یا:

$$15/24 < Y_f < 21/76 .$$

ج) برای به دست آوردن فاصله اطمینان ۹۵ درصد برای میانگین Y_f ، از نامساوی

۱۸-۳ استفاده می‌کنیم،

$$\hat{Y}_t - t_{\alpha/t} \hat{\sigma}_u \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}} < E(Y_t) < \hat{Y}_t + t_{\alpha/t} \hat{\sigma}_u \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}}$$

در نتیجه خواهیم داشت

$$18/0 + 2/182 \sqrt{0/0} \sqrt{1 + \frac{1}{0} + \frac{(10 - 4)^2}{40}} < E(Y_t) < 18/0 + 2/182 \sqrt{0/0} \sqrt{1 + \frac{1}{0} + \frac{(10 - 4)^2}{40}}$$

$$18/0 - 2/36 < E(Y_t) < 18/0 + 2/36$$

یا

$$16/14 < E(Y_t) < 20/16$$

ملاحظه می‌شود فاصله اطمینان برای $E(Y_t)$ از فاصله اطمینان برای Y_t کمتر است.

۲-۳ الف) $\hat{\alpha}$ و $\hat{\beta}$ را می‌توان به ترتیب زیر محاسبه کرد،

$$\hat{\beta} = \frac{\sum x_t y_t}{\sum x_t^2} = \frac{106/4}{210/4} = 0/494$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} = 1/090 - 0/494 (9/31)$$

ب) ابتدا σ_u^2 را تخمین می‌زنیم،

$$\begin{aligned} \hat{\sigma}_u^2 &= \frac{\sum e_t^2}{n-2} = \frac{\sum y_t^2 - \hat{\beta} \sum x_t y_t}{n-2} \\ &= \frac{87/9 - 0/494 (106/4)}{20-2} \\ &= 1/9077 \end{aligned}$$

با داشتن واریانس U_t ، تخمین واریانسهای $\hat{\alpha}$ و $\hat{\beta}$ به راحتی به دست می‌آید.

$$\begin{aligned} \text{Var}(\hat{\alpha}) &= \hat{\sigma}_u^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right), \\ &= 1/90.77 \left(\frac{1}{20} + \frac{9/31}{215/4} \right), \\ &= 0.8630. \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \frac{\hat{\sigma}_u^2}{\sum x_i^2}, \\ &= \frac{1/90.77}{215/4} = 0.0089. \end{aligned}$$

ج) می‌دانیم منظور از میانگین شرطی Y_f به شرط $X_f = 10$ ، دقیقاً همان مفهوم پیش‌بینی است؛ به عبارت دیگر می‌خواهیم بگوییم اگر X مقدار 10 بگیرد، میانگین Y چقدر خواهد بود؟ در معادله 13-3 دیدیم که

$$\hat{E}(Y_f) = \hat{Y}_f,$$

تخمین یا پیش‌بینی میانگین Y_f ، یعنی $\hat{E}(Y_f)$ ، دقیقاً برابر تخمین یا پیش‌بینی Y_f است. در بند (ج) می‌خواهیم مقدار $\hat{E}(Y_f)$ به ازای $X_f = 10$ را به دست آوریم؛ بنابراین کافی است \hat{Y}_f به ازای $X_f = 10$ را حساب کنیم. از بند (الف) می‌دانیم که

$$\hat{Y}_f = -3/50.4 + 0.494 X_f,$$

یا

$$\hat{Y}_f = -3/50.4 + 0.494 X_f.$$

بنابراین به ازای $X_f = 10$ داریم

$$\hat{Y}_f = -3/50.4 + 0.494(10) = 1/436.$$

د) ابتدا مقدار جدول t را با $18 = 20 - 2 = n - 2$ درجه آزادی و در سطح معنی‌دار $\alpha_p = 0.025$ به دست می‌آوریم، خواهیم داشت $t = \pm 2/101$. فاصله اطمینان اطمینان برای $E(Y_f)$ عبارت است از:

$$\hat{Y}_t - t_{\alpha/4} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}} < E(Y_t) < \hat{Y}_t + t_{\alpha/4} \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(X_t - \bar{X})^2}{\sum x_i^2}},$$

$$1/436 - 2/101 \sqrt{1/90.77} \sqrt{\frac{1}{20} + \frac{(10 - 9/31)^2}{210/8}} < E(Y_t) < 1/436$$

$$+ 2/101 \sqrt{1/90.77} \sqrt{\frac{1}{20} + \frac{(10 - 9/31)^2}{210/8}},$$

$$1/436 - 0.663 < E(Y_t) < 1/436 + 0.663,$$

در نتیجه

$$0.773 < E(Y_t) < 2.099.$$

۳-۳ می دانیم که تخمین β در یک مدل رگرسیون - که جمله ثابت ندارد - عبارت

است از

$$\hat{\beta} = \frac{\sum X_t Y_t}{\sum X_t^2}.$$

می خواهیم \hat{Y}_t^* را چنان به دست آوریم که

$$E(\hat{Y}_t^*) = E(Y_t^* | X_t). \quad (1)$$

ابتدا $E(Y_t^* | X_t)$ را بررسی می کنیم. می دانیم

$$Y_t | X_t = \beta X_t + U_t,$$

در نتیجه خواهیم داشت

$$E(Y_t^* | X_t) = E(\beta X_t + U_t) = E(\beta X_t + U_t + \beta X_t U_t).$$

با توجه به استقلال X_t از U_t و نیز $E(U_t) = 0$ و همچنین بنا بر فرض واریانس همسانی

داریم

$$E(Y_i^* | X_i) = \beta^T X_i^* + \sigma_u^2 \quad (2)$$

با توجه به مدل رگرسیون مفروض، ملاحظه می شود که یک مورد مناسب برای \hat{Y}_i^* مقدار $(\hat{\beta} X_i)^T$ است، اما در زیر نشان می دهیم که این تخمین زنده اریب دارد. در یک مدل رگرسیون بدون جمله ثابت، می دانیم

$$\hat{\beta} = \beta + \frac{\sum X_i U_i}{\sum X_i^T}$$

بنابراین اگر از $(\hat{\beta} X_i)^T$ امید ریاضی بگیریم، خواهیم داشت

$$\begin{aligned} E(\hat{\beta} X_i)^T &= X_i^T E(\hat{\beta})^T = X_i^T E\left(\beta + \frac{\sum X_i U_i}{\sum X_i^T}\right)^T, \\ &= X_i^T E\left[\beta^T + \frac{(\sum X_i U_i)^T}{(\sum X_i^T)^T} + \gamma \beta \frac{\sum X_i U_i}{\sum X_i^T}\right]. \end{aligned}$$

با بسط طرف راست و گرفتن امید ریاضی و با توجه به $E(U_i U_j) = 0$ ، نتیجه می گیریم که

$$E(\hat{\beta} X_i)^T = X_i^T \left(\beta^T + \frac{\sigma_u^2}{\sum X_i^T}\right) \quad (3)$$

خلاصه راه حل تا اینجا این است که می خواستیم Y_i^* را چنان تخمین بزنیم که معادله (۱) صادق باشد. دیدیم طرف راست معادله (۱) با معادله (۲) برابر است. بنابراین باید \hat{Y}_i^* را چنان بدست آورد که اگر از آن امید ریاضی بگیریم یا طرف راست معادله (۲) برابر شود. گفتیم $(\hat{\beta} X_i)^T$ جواب مسأله است. اما دیدیم که این تخمین زنده، اریب دارد؛ زیرا امید ریاضی آن، بر طبق معادله (۳)، با سمت راست معادله (۲) برابر نیست. البته به تعریف «اریب» در عبارتهای فوق توجه داریم - که تنها به معنای عدم صدق در رابطه (۱) به کار می رود - در اینجا باید کوشش کنیم مورد مناسب دیگری به عنوان یک پیش بینی کننده از Y_i^* پیدا کنیم که این اریب را نداشته باشد.

با مقایسه معادلات (۲) و (۳) می‌گوییم که پیش‌بینی‌کننده زیر مطلوب ماست،

$$\hat{Y}_t^* = E(\hat{\beta} X_t)^* - \frac{\sigma_u^2 X_t^*}{\sum X_t^*} + \sigma_u^2 \quad (4)$$

برای اثبات کافی است به جای $E(\hat{\beta} X_t)^*$ مقدار آن را از معادله (۳) قرار دهیم؛

$$\begin{aligned} \hat{Y}_t^* &= \left(\beta^* X_t^* + \frac{\sigma_u^2 X_t^*}{\sum X_t^*} \right) - \frac{\sigma_u^2 X_t^*}{\sum X_t^*} + \sigma_u^2 \\ &= \beta^* X_t^* + \sigma_u^2 \end{aligned}$$

مشاهده می‌شود که سمت راست معادله فوق دقیقاً با سمت راست معادله (۲) برابر است. بنابراین می‌توان نوشت:

$$\hat{Y}_t^* = E(Y_t^* | X_t) ,$$

با توجه به معادله (۲) و ثابت بودن مقادیر β ، X_t و σ^2 ، می‌توان نوشت

$$E(\hat{Y}_t^*) = E(Y_t^* | X_t) ,$$

که همان معادله (۱) است.

بدین ترتیب نشان دادیم که پیش‌بینی‌کننده‌ای که در معادله (۴) پیشنهاد شد، جواب مسأله است. در عمل باید به جای σ_u^2 در معادله (۴)، تخمین آن $(\hat{\sigma}_u^2)$ را قرار داد. می‌دانیم تخمین زنده ناریب از σ^2 عبارت است از:

$$\hat{\sigma}^2 = s^2 = \frac{\sum e_t^2}{n - 2} .$$

اما یادآوری می‌کنیم که در مدل‌های رگرسیون بدون جمله ثابت، باید در مخرج کسر به جای $(n - 2)$ مقدار $(n - 1)$ را قرار دهیم؛ زیرا برای به دست آوردن $\sum e_t^2$ فقط به $\hat{\beta}$ نیاز داریم و در نتیجه فقط یک درجه آزادی را از دست می‌دهیم.

۳.۴ به داخل پراتز سمت چپ $E(Y_{T+P}) \pm$ را اضافه می‌کنیم. خواهیم داشت

$$\begin{aligned} E(\hat{Y}_{T+F} - Y_{T+F})^2 &= E \left\{ [\hat{Y}_{T+F} - E(Y_{T+F})] + [E(Y_{T+F}) - (Y_{T+F})] \right\}^2, \\ &= E \left\{ [\hat{Y}_{T+F} - E(Y_{T+F})]^2 + [E(Y_{T+F}) - Y_{T+F}]^2 \right. \\ &\quad \left. + 2[\hat{Y}_{T+F} - E(Y_{T+F})][E(Y_{T+F}) - Y_{T+F}] \right\}. \end{aligned}$$

می‌دانیم که

$$E[\hat{Y}_{T+F} - E(Y_{T+F})]^2 = \text{Var}^*(Y_{T+F}),$$

و منظور از $\text{Var}^*(Y_{T+F})$ این است که واریانس در حول مقدار میانگین Y_{T+F} تعریف شده است. همچنین از رابطه

$$Y_{T+F} = \alpha + \beta X_{T+F} + U_{T+F},$$

می‌دانیم که

$$E(Y_{T+F}) = \alpha + \beta X_{T+F},$$

بنابراین با کم کردن دو رابطه فوق از یکدیگر داریم

$$E(Y_{T+F}) - Y_{T+F} = -U_{T+F}.$$

بدین ترتیب خواهیم داشت

$$\begin{aligned} E(\hat{Y}_{T+F} - Y_{T+F})^2 &= \text{Var}^*(Y_{T+F}) + \sigma_U^2 \\ &\quad + 2E \left\{ [\hat{Y}_{T+F} - E(Y_{T+F})][E(Y_{T+F}) - (Y_{T+F})] \right\}. \end{aligned}$$

نشان می‌دهیم که در معادله فوق، جمله سوم سمت راست، برابر صفر است،

$$E \left\{ [\hat{Y}_{T+F} - E(Y_{T+F})][E(Y_{T+F}) - (Y_{T+F})] \right\} =$$

$$E[\hat{Y}_{T+F}E(Y_{T+F}) - \hat{Y}_{T+F}Y_{T+F} - E(Y_{T+F})E(Y_{T+F}) + E(Y_{T+F})Y_{T+F}]. \quad (1)$$

در معادله ۳-۱۲ دیدیم که

$$E(Y_t) = \alpha + \beta X_t ,$$

و اگر از معادله ۳-۴ نیز امید ریاضی بگیریم، خواهیم داشت

$$E(\hat{Y}_t) = \alpha + \beta X_t .$$

بنابراین داریم

$$E(Y_t) = E(\hat{Y}_t) .$$

به همین ترتیب خواهیم داشت

$$E(Y_{T+F}) = E(\hat{Y}_{T+F}) , \quad (۲)$$

و در نتیجه طرف راست معادله (۱) برابر صفر می‌شود؛ بنابراین می‌توان گفت

$$\text{Var}(\hat{Y}_{T+F}) = \text{Var}^*(\hat{Y}_{T+F}) + \sigma^2 .$$

۳.۵ از رابطه $Y_t^* = p_1 + q_1 Y_t$ داریم $\bar{Y}^* = p_1 + q_1 \bar{Y}$ ، و در نتیجه خواهیم داشت

$$y_t^* = q_1 y_t . \quad (۱)$$

به همین ترتیب می‌توان نشان داد که

$$x_t^* = q_2 x_t . \quad (۲)$$

می‌دانیم در مدل رگرسیون

$$Y_t^* = a + b X_t^* + U_t^* ,$$

ضریب تعیین (r^{*2}) عبارت است از

$$r^{*2} = \frac{(\sum x_t^* y_t^*)^2}{\sum x_t^{*2} \sum y_t^{*2}} .$$

با جایگزینی رابطه (۱) و (۲) در معادله فوق خواهیم داشت

$$r^{*2} = \frac{(\sum q_i x_i q_i y_i)^2}{\sum (q_i x_i)^2 \sum (q_i y_i)^2}$$

$$= \frac{q_i^2 q_i^2 (\sum x_i y_i)^2}{q_i^2 q_i^2 \sum x_i^2 \sum y_i^2} = r^2$$

۳-۶ پاسخ به این سؤال در متن کتاب و در بحث از معادله ۳-۴۵ ارائه شده است. با وجود این، در اینجا به راه حل دیگری اشاره می‌کنیم. می‌دانیم در مدل رگرسیون

$$Y_i = \alpha + \beta X_i + U_i \quad (1)$$

و با توجه به معادله ۱-۳۵ داریم

$$r^2 = \frac{\hat{\beta} \sum x_i y_i}{\sum y_i^2} = \frac{\hat{\beta}^2 \sum x_i^2}{\sum y_i^2}$$

هرگاه $r^2 = 1$ ، آنگاه خواهیم داشت

$$\hat{\beta}^2 = \frac{\sum y_i^2}{\sum x_i^2}$$

یا

$$\sum x_i^2 = \frac{\sum y_i^2}{\hat{\beta}^2} \quad (2)$$

اما برای مدل (۱) داریم

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

که با جایگزینی معادله (۲) خواهیم داشت

$$1 = \hat{\beta} \frac{\sum x_i y_i}{\sum x_i^2}$$

یا

$$\hat{\beta} = \frac{\sum y_i^2}{\sum x_i y_i} \quad (3)$$

اما برای مدل رگرسیون معکوس $X_i = \alpha + \beta' Y_i + V_i$ داریم

$$\hat{\beta}' = \frac{\sum x_i y_i}{\sum x_i^2}$$

$$\frac{1}{\hat{\beta}'} = \frac{\sum y_i^2}{\sum x_i y_i} \quad (۴)$$

با مقایسه معادله‌های (۳) و (۴) نتیجه می‌گیریم که

$$\hat{\beta} = \frac{1}{\hat{\beta}'} \quad ۳.۷$$

$$\begin{aligned} \sum x_i^2 &= \sum X_i^2 - \frac{1}{n} (\sum X_i)^2 \\ &= ۱۲/۱۶ - \frac{1}{۲۰۰} (۱۱/۳۴)^2 = ۱۱/۵۱۷, \end{aligned}$$

$$\begin{aligned} \sum y_i^2 &= \sum Y_i^2 - \frac{1}{n} (\sum Y_i)^2 \\ &= ۸۴/۹۶ - \frac{1}{۲۰۰} (۲۰/۷۲)^2 = ۸۲/۸۱۳, \end{aligned}$$

$$\begin{aligned} \sum x_i y_i &= \sum X_i Y_i - \frac{1}{n} (\sum X_i) (\sum Y_i) \\ &= ۲۲/۱۳ - \frac{1}{۲۰۰} (۱۱/۳۴)(۱۲/۱۶) = ۲۰/۹۵۵. \end{aligned}$$

با محاسبات فوق به راحتی می‌توان α و β را تخمین زد.

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = ۱/۸۱۹, \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} = ۰/۰۰۰۴.$$

پارامترهای رگرسیون معکوس را به ترتیب زیر تخمین می‌زنیم.

$$\hat{\beta}' = \frac{\sum x_i y_i}{\sum y_i^2} = ۰/۲۵۳۰, \quad \hat{\alpha}' = \bar{X} - \hat{\beta}' \bar{Y} = ۰/۰۳۰.$$

۳-۸ می‌دانیم برای مدل رگرسیون مستقیم داریم

$$\hat{\beta} = \frac{\sum X_i Y_i}{\sum X_i^2}$$

برای تخمین β' در مدل رگرسیون معکوس $X_i = \beta' Y_i + V_i$ خواهیم داشت

$$\hat{\beta}' = \frac{\sum X_i Y_i}{\sum Y_i^2}$$

فوق $\hat{\beta}$ و $\hat{\beta}'$ فقط در مخرج کسر است. اگر $\sum X_i^2 > \sum Y_i^2$ ، آنگاه $\hat{\beta} < \hat{\beta}'$ می‌دانیم r^2 در مدل رگرسیون مستقیم برابر است با

$$r^2 = \frac{(\sum X_i Y_i)^2}{\sum X_i^2 \sum Y_i^2}$$

و برای رگرسیون معکوس نیز دقیقاً همین مقدار خواهد بود؛ بنابراین مقادیر r^2 با یکدیگر برابر است. برای تخمین واریانس جملهٔ اختلال یا $\hat{\sigma}^2$ ، باید ابتدا مجموع مربعات پسماند را در هر یک از دو مدل به دست آورد. برای رگرسیون مستقیم داریم،

$$r^2 = 1 - \frac{\sum e_i^2}{\sum Y_i^2}$$

$$\frac{(\sum X_i Y_i)^2}{\sum X_i^2 \sum Y_i^2} = 1 - \frac{\sum e_i^2}{\sum Y_i^2}$$

در نتیجه:

$$\sum e_i^2 = \frac{\sum X_i^2 \sum Y_i^2 - (\sum X_i Y_i)^2}{\sum X_i^2}$$

برای تخمین σ^2 کافی است $\sum e_i^2$ را بر درجات آزادی، یعنی $(n-1)$ تقسیم کنیم،

$$\hat{\sigma}^2 = \frac{\sum X_i^2 \sum Y_i^2 - (\sum X_i Y_i)^2}{(n-1) \sum X_i^2}$$

به روشی کاملاً مشابه می‌توان $\hat{\sigma}^2$ را برای مدل رگرسیون معکوس به دست آورد.

خواهیم داشت

$$\hat{\sigma}^2 = \frac{\sum X_i^2 \sum Y_i^2 - (\sum X_i Y_i)^2}{(n-1) \sum X_i^2}$$

ملاحظه می شود تفاوت $\hat{\sigma}^2$ در مدل های مستقیم و معکوس فقط در مخرج کسر است. اگر داشته باشیم $\sum Y_i^2 > \sum X_i^2$ آنگاه $\hat{\sigma}^2 < \hat{\sigma}^2$ ، یعنی تخمین واریانس جمله اختلال در مدل رگرسیون مستقیم، از مقدار مشابه در مدل معکوس کمتر خواهد بود.

۳-۹ الف) از دو طرف لگاریتم می گیریم،

$$\ln Y_i = \ln \alpha + \beta X_i + \ln U_i ,$$

و با تعریف جدیدی از متغیرها خواهیم داشت

$$Y_i^* = \alpha^* + \beta X_i + U_i^* ,$$

که برحسب پارامتر β و $\ln \alpha$ خطی است. به هر حال هر دو پارامتر را می توان با روش حداقل مربعات معمولی تخمین زد و برای رسیدن به تخمین α کافی است از $\hat{\alpha}^*$ ، آنتی لگاریتم بگیریم.

ب) این مدل را نمی توان خطی کرد زیرا U_i در مدل، خصوصیت جمع پذیری دارد.

ج) از دو طرف لگاریتم می گیریم،

$$\ln Y_i = \alpha + \beta X_i + U_i ,$$

که یک مدل خطی برحسب α و β است.

د) با اینکه در ظاهر به نظر می رسد که این مدل را نمی توان خطی کرد، چون U_i جمع پذیر است، اما با استفاده از روش «ثابت نگه داشتن یک پارامتر» - که در بند ۵ قسمت ۳-۶ ارائه شد - بسادگی می توان مدل مزبور را خطی کرد. ابتدا مجموع مربعات پسماند را تشکیل می دهیم،

$$\sum e_i^2 = \left[Y_i - \frac{\hat{\alpha}}{(\hat{\beta} - X_i)} \right]^2 . \quad (1)$$

با ثابت نگه داشتن $\hat{\beta}$ ، متغیر Z_1 را به صورت زیر تعریف می‌کنیم،

$$Z_1 = \frac{1}{\hat{\beta} - X_1} \quad (۲)$$

اگر در محدوده تغییرات $\hat{\beta}$ ، یک مقدار اختیاری انتخاب کرده و در معادله (۲) قرار دهیم و Z_1 به دست آمده را در معادله (۱) بگذاریم، خواهیم داشت

$$\sum e_i^2 = (Y_i - \hat{\alpha} Z_1)^2$$

به راحتی می‌توان مقداری از $\hat{\alpha}$ که e_i^2 را حداقل می‌کند به دست آورد. با تغییر مقدار $\hat{\beta}$ در معادله (۲) می‌توان به مقادیر جدیدی از Z_1 و e_i^2 رسید. اگر تمام مقادیر ممکن $\hat{\beta}$ را بدین ترتیب آزمایش کنیم، به تمام مقادیر ممکن e_i^2 دست یافته و سپس کمترین آن را مشخص می‌کنیم. $\hat{\alpha}$ و $\hat{\beta}$ متناظر با این مقدار از e_i^2 ، تخمین زنده‌های روش حداقل مربعات معمولی از پارامترهای مدل مفروض خواهد بود. بدیهی است راه حل فوق، مستلزم شناخت دقیقی از محدوده تغییرات β است.

۳-۱۰ الف) از دو طرف معادله لگاریتم می‌گیریم،

$$\ln Y_i = \ln \alpha + \beta \ln X_i$$

متغیرهای زیر را تعریف می‌کنیم،

$$Y_i^* = \ln Y_i \quad , \quad \alpha^* = \ln \alpha \quad , \quad X_i^* = \ln X_i \quad ,$$

بنابراین خواهیم داشت

$$Y_i^* = \alpha^* + \beta X_i^*$$

اگر جمله اختلال، جمع پذیر باشد، روش فوق موفقیت آمیز نبوده و باید از روشهای تخمین غیرخطی استفاده کرد.

ب) از دو طرف معادله، لگاریتم می‌گیریم،

$$\ln Y_i = \ln \alpha + \beta X_i$$

متغیرهای زیر را تعریف می‌کنیم،

$$Y_i^* = \ln Y_i \quad , \quad \alpha^* = \ln \alpha \quad ,$$

بنابراین خواهیم داشت

$$Y_i^* = \alpha^* + \beta X_i^* .$$

اگر U_i جمع پذیر باشد، باید از روشهای تخمین غیرخطی استفاده کرد.
ج) تعریف می‌کنیم $X_i^* = \ln X_i$. با جایگزینی در مدل اولیه داریم،

$$Y_i = \alpha + \beta X_i^* .$$

U_i می‌تواند جمع پذیر باشد.

(د) دو طرف مدل مفروض را معکوس می‌کنیم،

$$\frac{1}{Y_i} = \alpha - \frac{\beta}{X_i} .$$

این متغیرها را تعریف می‌کنیم،

$$Y_i^* = \frac{1}{Y_i} \quad , \quad X_i^* = \frac{1}{X_i} \quad ,$$

در نتیجه خواهیم داشت

$$Y_i^* = \alpha + \beta X_i^* .$$

اگر U_i جمع پذیر باشد آنگاه روش فوق قابل اجرا نبوده و باید از روشهای تخمین غیرخطی استفاده کرد.

(ه) ابتدا $\frac{Y_i}{1-Y_i}$ را تشکیل می‌دهیم،

$$\begin{aligned} \frac{Y_i}{1-Y_i} &= \frac{e^{(\alpha + \beta X_i)} / [1 + e^{(\alpha + \beta X_i)}]}{1 - e^{(\alpha + \beta X_i)} / [1 + e^{(\alpha + \beta X_i)}]} \\ &= e^{(\alpha + \beta X_i)} . \end{aligned}$$

از دو طرف معادله لگاریتم می‌گیریم،

$$\ln \left(\frac{Y_t}{1 - Y_t} \right) = \alpha + \beta X_t .$$

متغیر زیر را تعریف می‌کنیم.

$$Y_t^* = \frac{Y_t}{1 - Y_t} ,$$

در نتیجه خواهیم داشت

$$Y_t^* = \alpha + \beta X_t .$$

اگر U_t جمع‌پذیر باشد، باید از روشهای تخمین غیرخطی استفاده کرد. (و) تعریف می‌کنیم $X_t^* = \sqrt{X_t}$ ، در این صورت خواهیم داشت

$$Y_t = \alpha + \beta X_t^* .$$

جمله اختلال (U_t) می‌تواند جمع‌پذیر باشد.

ز) در تخمین این مدل باید از روشهای تخمین غیرخطی استفاده کرد، مگر آنکه محدوده تغییرات β معلوم باشد. در این صورت، مشابه بند (د) مسأله ۹-۳ عمل می‌کنیم. برای هر مقدار معلوم β ، به راحتی می‌توان این کمیت را محاسبه کرد،

$$X_t^* = e^{\beta X_t} .$$

فرض می‌کنیم U_t جمع‌پذیر باشد. با جایگذاری معادله فوق در مدل مفروض داریم،

$$Y_t = \alpha + X_t^* + U_t ,$$

$$(Y_t - X_t^*) = \alpha + U_t ,$$

در نتیجه

$$Y_t^* = \alpha + U_t , \quad (1)$$

که در آن $Y_i^* = Y_i - X_i^* \alpha$ است. ملاحظه می‌شود که در مدل (۱)، Y_i^* فقط تابعی از α و U_i است. در مسأله ۱۰-۱ دیدیم که در این گونه مدل‌های رگرسیون تخمین α با میانگین متغیر درون‌زا برابر است،

$$\hat{\alpha} = \bar{Y}^* .$$

خلاصه بحث این است که به ازای مقدار معینی از β ، مدل مفروض را به یک مدل خطی تبدیل کرده و $\hat{\alpha}$ را تخمین زده‌ایم. حال باید مقدار e_i^* را به دست آوریم،

$$\begin{aligned} \sum e_i^* &= \sum (Y_i^* - \hat{\alpha})^2 , \\ &= \sum (Y_i^* - \bar{Y}^*)^2 = \sum y_i^{*2} . \end{aligned}$$

به همین ترتیب از مقادیر ممکن دیگر β نیز استفاده کرده، با به دست آوردن تخمین‌هایی از $\hat{\alpha}$ به مقادیر ممکن $\sum e_i^*$ می‌رسیم. کمترین مقدار $\sum e_i^*$ را مشخص کرده، مقادیر $\hat{\alpha}$ و $\hat{\beta}$ متناظر با آن، جواب مسأله خواهد بود.

ک) این مسأله شبیه مثالی است که در بند ۵ قسمت ۶-۳ و با عنوان «مدلهای غیرخطی برحسب پارامترها» مطرح شده است.

مدل رگرسیون خطی با دو متغیر توضیحی

۴-۱ مقدمه

تا کنون مدلی را مطالعه کرده‌ایم که فقط یک متغیر توضیحی داشته است. در مطالعات کاربردی معمولاً باید بیش از یک متغیر توضیحی را در نظر گرفت؛ زیرا بیان تغییرات یک متغیر درون‌زا به کمک چند متغیر توضیحی به مراتب دقیقتر خواهد بود. ورود بیش از یک متغیر توضیحی در مدل‌های رگرسیون، مسائل مختلفی را مطرح می‌کند. برای اینکه بتوان این مسائل را به نحو منظمی بررسی کرد، ابتدا مدلی را در نظر می‌گیریم که فقط دو متغیر توضیحی دارد و سپس در فصل پنجم این مدل را به k متغیر توضیحی تعمیم می‌دهیم. معمولاً به مدلی که بیش از یک متغیر توضیحی دارد، مدل «رگرسیون چند متغیره»^۱ می‌گوییم. بنابراین مدل‌های رگرسیون با دو متغیر توضیحی در واقع ساده‌ترین مدل رگرسیون چند متغیره است. این مدل را می‌توان چنین نوشت،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t .$$

تعمیم این مدل به حالت عمومی مدل‌های رگرسیون چند متغیره که این است که به جای دو متغیر از k متغیر توضیحی استفاده کنیم،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + U_t .$$

تخمین یک مدل رگرسیون چند متغیره با k متغیر توضیحی، به کمک جبر ماتریسی بسیار ساده‌تر انجام می‌شود، بنابراین بررسی حالت عمومی مدل‌های رگرسیون

با استفاده از ماتریسها^۱، موضوع فصل پنجم خواهد بود.

تخمین پارامترهای α ، β_1 و β_2 در مدل رگرسیون با دو متغیر توضیحی موضوع قسمت ۴-۲ است. در این قسمت مسائل ضریب تعیین نیز بررسی خواهد شد. تحلیل‌های آماری تخمین پارامترها شامل واریانس، کوواریانس و آزمون فرضیه‌های مختلف، موضوع قسمت ۴-۳ است. با اینکه جدول تجزیه واریانس و آزمون F ، در تحلیل‌های رگرسیون از اهمیت ویژه‌ای برخوردار است، اما به نظر می‌رسد بهتر است مباحث تجزیه واریانس در مدل‌های رگرسیون با دو متغیر توضیحی را به صورت مورد خاص از تجزیه واریانس در حالت کلی در نظر بگیریم؛ بنابراین در فصل ششم و در مطالعه تجزیه واریانس در رگرسیون چند متغیره، ابعاد این مسأله در مورد رگرسیون‌هایی با دو متغیر توضیحی نیز روشن خواهد شد.

مسأله پیش‌بینی، همواره یکی از مهمترین مباحث تخمین در مدل‌های رگرسیون بوده است. روال کلی بحث در یک مدل رگرسیون با دو متغیر توضیحی دقیقاً شبیه رگرسیون ساده است؛ با این تفاوت که در محاسبه خطای پیش‌بینی، باید از واریانس و کوواریانس تمام پارامترها استفاده شود، در قسمت ۴-۴ به طور خلاصه مسأله پیش‌بینی را بررسی خواهیم کرد.

هنگامی که یک مدل رگرسیون ساده را تعمیم می‌دهیم و به جای یک متغیر توضیحی، از دو یا چند متغیر توضیحی استفاده می‌کنیم، مسأله جدیدی مطرح می‌شود و آن ضرایب همبستگی جزئی و ضرایب تعیین جزئی است که در تحلیل و تفسیر نتایج حاصل از تخمین پارامترهای مدل، اهمیت فراوان دارد. آیا نمی‌توان رگرسیون چند متغیره را به چند رگرسیون ساده تجزیه کرد؟ به عبارت دیگر آیا نمی‌توان یک بار به کمک Y_1 و X_{11} و بار دیگر با Y_2 و X_{21} دو رگرسیون جداگانه ساخت و پارامترهای هر یک را تخمین زد؟ اساساً ورود همزمان X_{11} و X_{21} به عنوان متغیرهای توضیحی، چه نکته‌هایی را می‌تواند روشن کند که دو مدل رگرسیون ساده مجزا از هم نمی‌تواند آن را

تعیین کند؟ اینها نکته‌هایی است که در قسمت ۴-۵ و با عنوان ضرایب همبستگی و ضرایب تعیین بررسی خواهد شد. در این قسمت نشان خواهیم داد که چگونه می‌توان ضریب تعیین کلی یک مدل رگرسیون را بر حسب ضرایب همبستگی جزئی نوشت. همچنین به این نکته نیز توجه خواهیم کرد که چگونه تخمین پارامترهای یک مدل رگرسیون چند متغیره را می‌توان بر حسب تخمین پارامترهای رگرسیون ساده یا بر حسب ضرایب همبستگی ساده بیان کرد.

۴-۲ تخمین مدل

مدل رگرسیون زیر مفروض است،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t \quad (4-1)$$

که در آن $t = 1, 2, \dots, n$ و Y_t متغیر درون‌زا و X_{1t} و X_{2t} متغیرهای برون‌زا است. مانند گذشته فرض بر این است که U_t یک متغیر تصادفی است. و تمام فرضهای کلاسیک را به این شرح داراست.

۱. میانگین آن صفر است؛ یعنی $E(U_t) = 0$.
۲. واریانس همسانی دارد؛ یعنی واریانس آن به ازای تمام مقادیر t ثابت است، $Var(U_t) = \sigma^2$.
۳. خودهمبستگی ندارد؛ یعنی U_t و U_s به ازای تمام مقادیر $t \neq s$ از یکدیگر مستقل است؛ یعنی $Cov(U_t, U_s) = 0$.
۴. به ازای تمام مقادیر t تابع توزیع احتمال U_t نرمال است. فرضهای فوق در مورد جمله اختلال است. با ارائه فرضهای پنجم و ششم در مورد متغیرهای توضیحی، فرضهای کلاسیک مدلهای رگرسیون با بیش از یک متغیر توضیحی کامل می‌شود.
۵. متغیر توضیحی X_t غیر تصادفی است؛ یعنی X_t از جمله اختلال مستقل است. با استفاده از ۵ فرض فوق و دقیقاً مانند مباحث فصل دوم می‌توان نشان داد که $\hat{\alpha}$ ، $\hat{\beta}_1$ و

$\hat{\beta}_4$ نارایب است و حداقل واریانس را در گروه تخمین زنده‌های خطی دارد. فرض ششم که در ذیل مطرح خواهد شد، فقط به مدل‌های رگرسیون چند متغیره اختصاص دارد.

۶. متغیرهای توضیحی «همخطی» ندارند؛ یعنی هیچ نوع همبستگی کامل خطی بین X_1 و X_2 وجود ندارد؛ به عبارت دیگر نمی‌توان هیچگونه رابطه خطی بین متغیرهای برون‌زا برقرار کرد. مسأله همخطی از مباحث بسیار مهم اقتصادسنجی است و در جلد دوم این کتاب بررسی خواهد شد. اما برای تبیین بیشتر به ذکر مثالی اکتفا می‌کنیم.

در مدل رگرسیون ۴-۱، این رابطه خطی بین متغیرهای توضیحی مفروض است:

$$3X_{1i} + X_{2i} = 0.$$

می‌توان X_{2i} را بر حسب X_{1i} نوشت،

$$X_{2i} = 0 - 3X_{1i}.$$

با جایگزینی معادله فوق در مدل رگرسیون ۴-۱ خواهیم داشت

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 (0 - 3X_{1i}) + U_i.$$

یا

$$Y_i = (\alpha + 0\beta_2) + (\beta_1 - 3\beta_2) X_{1i} + U_i.$$

با تخمین مدل فوق آنچه به دست خواهیم آورد، تخمینهای $(\alpha + 0\beta_2)$ و $(\beta_1 - 3\beta_2)$ است؛ در حالی که هدف ما این بود که α ، β_1 و β_2 را تخمین بزنیم. مشاهده می‌شود که در موارد «همخطی کامل»^۱، یعنی وجود یک رابطه معین و مشخص خطی بین متغیرها، تخمین هر یک از پارامترهای مدل ممکن نیست، مگر اینکه بتوان به کمک نظریه‌های اقتصادی، رابطه دیگری نیز بین پارامترها برقرار کرد. در مثال فوق فقط دو تخمین $(\hat{\alpha} + 0\hat{\beta}_2)$ و $(\hat{\beta}_1 - 3\hat{\beta}_2)$ را به دست می‌آوریم؛ در حالی که سه مجهول $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ داریم. باید یک رابطه خطی دیگر به کمک نظریه‌های اقتصادی بین پارامترها به دست آورد تا تخمینهای جداگانه از آنها ممکن شود. بدیهی است در بسیاری موارد چنین

روابطی را نمی‌توان از نظریه‌های اقتصادی استخراج کرد؛ بنابراین مسأله همخطی از مباحث قابل توجه در اقتصادسنجی است که ابعاد مختلف آن در جلد دوم این کتاب بررسی خواهد شد.

در مواردی (مانند مثال فوق) که همبستگی خطی بین متغیرهای توضیحی کامل باشد؛ یعنی «همخطی کامل»، داشته باشیم، ضریب همبستگی بین X_{1t} و X_{2t} برابر ± 1 است. در این فصل فرض بر این است که بین متغیرهای توضیحی همبستگی وجود دارد، اما کامل نیست. حال که با فرض ششم (مختص رگرسیونهای چند متغیره) آشنا شدیم به تخمین مدل ۴-۱ می‌پردازیم. فرض کنید مدل ۴-۱ را به صورت زیر تخمین زده‌ایم،

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta}_1 X_{1t} + \hat{\beta}_2 X_{2t} \quad (4-2)$$

پسماند را به روال گذشته تعریف می‌کنیم،

$$e_t = Y_t - \hat{Y}_t$$

و مجموع مربعات پسماند را می‌نویسیم،

$$\sum_{i=1}^n e_t^2 = \sum_{i=1}^n (Y_t - \hat{\alpha} - \hat{\beta}_1 X_{1t} - \hat{\beta}_2 X_{2t})^2 \quad (4-3)$$

معیار روش حداقل مربعات معمولی این است که $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ را چنان تعیین کنیم که $\sum e_t^2$ در معادله ۴-۳ حداقل شود؛ بنابراین باید مشتق $\sum e_t^2$ را نسبت به متغیرهای $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ به دست آوریم و مساوی صفر قرار دهیم،

$$\frac{\partial \sum e_t^2}{\partial \hat{\alpha}} = \sum 2 (Y_t - \hat{\alpha} - \hat{\beta}_1 X_{1t} - \hat{\beta}_2 X_{2t}) (-1) = 0$$

$$\frac{\partial \sum e_t^2}{\partial \hat{\beta}_1} = \sum 2 (Y_t - \hat{\alpha} - \hat{\beta}_1 X_{1t} - \hat{\beta}_2 X_{2t}) (-X_{1t}) = 0 \quad (4-4)$$

$$\frac{\partial \sum e_t^2}{\partial \hat{\beta}_2} = \sum 2 (Y_t - \hat{\alpha} - \hat{\beta}_1 X_{1t} - \hat{\beta}_2 X_{2t}) (-X_{2t}) = 0$$

به معادله‌های ۴-۴ سیستم معادلات نرمال می‌گویند. این معادله‌ها را می‌توان به شرح زیر ساده نمود،

$$\begin{aligned} \sum Y_i &= n \hat{\alpha} + \hat{\beta}_1 \sum X_{1i} + \hat{\beta}_2 \sum X_{2i} , \\ \sum X_{1i} Y_i &= \hat{\alpha} \sum X_{1i} + \hat{\beta}_1 \sum X_{1i}^2 + \hat{\beta}_2 \sum X_{1i} X_{2i} , \quad (4.5) \\ \sum X_{2i} Y_i &= \hat{\alpha} \sum X_{2i} + \hat{\beta}_1 \sum X_{1i} X_{2i} + \hat{\beta}_2 \sum X_{2i}^2 . \end{aligned}$$

سیستم سه معادله نرمال ۴-۵ را - که بر حسب سه مجهول $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ است - می‌توان به سهولت حل کرد و تخمینهای حداقل مربعات معمولی از پارامترهای α ، β_1 و β_2 را به دست آورد.

یک روش ساده وجود دارد که می‌توان سیستم معادله‌های نرمال ۴-۵ را بدون مشتق‌گیری از $\sum e_i^2$ نوشت. معادله ۴-۱ را یک بار دیگر می‌نویسیم،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i .$$

از دو طرف معادله فوق \sum گرفته و جمله آخر، یعنی $\sum U_i$ را مساوی صفر فرض می‌کنیم^۱. بدین ترتیب معادله اول نرمال در سیستم معادله‌های ۴-۵ به دست می‌آید، با این تفاوت که به جای $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ مقادیر α ، β_1 و β_2 را خواهیم داشت. در ادامه بحث، اشاره خواهیم کرد که این امر مسأله خاصی را به وجود نمی‌آورد. یک بار دیگر دو طرف معادله ۴-۱ را در X_{1i} ضرب کرده و \sum می‌گیریم و با حذف جمله آخر، یعنی $\sum X_{1i} U_i$ ، به دومین معادله نرمال سیستم ۴-۵ می‌رسیم. سرانجام دو طرف معادله ۴-۱ را در X_{2i} ضرب کرده و با گرفتن \sum و حذف جمله آخر به سومین معادله نرمال خواهیم رسید. اگر به جای دو متغیر توضیحی، k متغیر توضیحی نیز داشتیم، به همین ترتیب، متغیرهای توضیحی دیگر را نیز یک به یک در دو طرف معادله ۴-۱ ضرب کرده و به ترتیب فوق معادله‌های نرمال دیگر حاصل می‌شود. سیستم معادلاتی را - که به این ترتیب و بر حسب α ، β_1 و

۱. در مسأله ۱-۱۱ دیدیم که ضروری نیست $\sum U_i$ برابر صفر باشد، هر چند $E(U_i)$ صفر است.

β_2 به دست می‌آید - حل می‌کنیم. جوابهای حاصل $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ خواهد بود. این راه‌حل بیشتر شبیه روشی برای از بر کردن معادله‌های نرمال است و مبنای نظری ندارد؛ با وجود این، کاربردهای دیگر این روش را در جلد دوم با عنوان روش تخمین «متغیرهای ابزاری»^۱ بررسی خواهیم کرد.

سیستم معادله‌های ۴-۵ را می‌توان به صورت ساده‌تری نیز نوشت. با تقسیم دو طرف معادله اول نرمال بر n ، خواهیم داشت:

$$\bar{Y} = \hat{\alpha} + \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2. \quad (4-6)$$

با جایگزینی مقدار $\hat{\alpha}$ از معادله ۴-۶ در معادله دوم نرمال در سیستم معادله‌های ۴-۵ خواهیم داشت:

$$\sum X_{1i} Y_i = n \bar{X}_1 (\bar{Y} - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2) + \hat{\beta}_1 \sum X_{1i}^2 + \hat{\beta}_2 \sum X_{1i} X_{2i}. \quad (4-7)$$

برای ساده‌تر کردن معادله ۴-۷، متغیرهای زیر را تعریف می‌کنیم،

$$\sum x_{1i}^2 = \sum (X_{1i} - \bar{X}_1)^2 = \sum X_{1i}^2 - n \bar{X}_1^2,$$

$$\sum x_{1i} x_{2i} = \sum (X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2) = \sum X_{1i} X_{2i} - n \bar{X}_1 \bar{X}_2,$$

$$\sum x_{2i}^2 = \sum (X_{2i} - \bar{X}_2)^2 = \sum X_{2i}^2 - n \bar{X}_2^2,$$

$$\sum x_{1i} y_i = \sum (X_{1i} - \bar{X}_1)(Y_i - \bar{Y}) = \sum X_{1i} Y_i - n \bar{X}_1 \bar{Y},$$

$$\sum x_{2i} y_i = \sum (X_{2i} - \bar{X}_2)(Y_i - \bar{Y}) = \sum X_{2i} Y_i - n \bar{X}_2 \bar{Y},$$

$$\sum y_i^2 = \sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - n \bar{Y}^2.$$

با استفاده از تعاریف فوق، معادله ۴-۷ را می‌توان به صورت زیر نوشت،

$$\sum x_{1t} y_t = \hat{\beta}_1 \sum x_{1t}^2 + \hat{\beta}_2 \sum x_{1t} x_{2t} \quad (4.8)$$

با روشی کاملاً مشابه، معادله سوم نرمال را نیز به صورت زیر می‌نویسیم،

$$\sum x_{2t} y_t = \hat{\beta}_1 \sum x_{1t} x_{2t} + \hat{\beta}_2 \sum x_{2t}^2 \quad (4.9)$$

معادله‌های ۴.۸ و ۴.۹ را می‌توان برای $\hat{\beta}_1$ و $\hat{\beta}_2$ حل کرد،

$$\hat{\beta}_1 = \frac{\sum x_{2t}^2 \sum x_{1t} y_t - \sum x_{1t} x_{2t} \sum x_{2t} y_t}{\sum x_{1t}^2 \sum x_{2t}^2 - (\sum x_{1t} x_{2t})^2} \quad (4.10)$$

$$\hat{\beta}_2 = \frac{\sum x_{1t}^2 \sum x_{2t} y_t - \sum x_{1t} x_{2t} \sum x_{1t} y_t}{\sum x_{1t}^2 \sum x_{2t}^2 - (\sum x_{1t} x_{2t})^2}$$

با به دست آوردن $\hat{\beta}_1$ و $\hat{\beta}_2$ می‌توان $\hat{\alpha}$ را از معادله ۴.۶ به دست آورد.

به طور خلاصه برای تخمین پارامترهای مدل ۴-۱ به این ترتیب عمل می‌کنیم.

۱. مقادیر میانگینها را حساب می‌کنیم: \bar{Y} ، \bar{X}_1 ، \bar{X}_2 .

۲. مجموع مربعات و مجموع حاصلضرب تمام متغیرها را به دست می‌آوریم:

$$\sum x_{1t}^2, \sum x_{2t}^2, \sum x_{1t} x_{2t}, \dots$$

۳. مقادیر $\sum x_{1t}^2$ ، $\sum x_{2t}^2$ ، $\sum x_{1t} x_{2t}$ ، $\sum x_{1t} y_t$ ، $\sum x_{2t} y_t$ و $\sum y_t^2$ را محاسبه

می‌کنیم.

۴. با استفاده از معادله‌های ۴-۱۰ مقادیر $\hat{\beta}_1$ و $\hat{\beta}_2$ را به دست می‌آوریم.

۵. از معادله ۴-۶ مقدار $\hat{\alpha}$ را حساب می‌کنیم.

بعد از آشنایی با نحوه تخمین پارامترهای مدل ۴-۱ به چند نکته اشاره می‌شود.

الف) با توجه به معادله ۴-۶ می‌توان گفت که صفحه تخمین رگرسیون از نقطه

میانگینها می‌گذرد.

ب) اگر معادله ۴-۲ را برای تمام مقادیر نمونه جمع و بر n تقسیم کنیم، خواهیم داشت:

$$\bar{Y} = \hat{\alpha} + \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2$$

با مقایسه معادله فوق با معادله ۴-۶ می توان نتیجه گرفت که

$$\hat{Y} = \bar{Y} \quad (۴-۱۱)$$

ج) می دانیم $Y_i = \hat{Y}_i + e_i$. دو طرف معادله را برای تمام مقادیر t جمع کرده و بر n تقسیم می کنیم،

$$\bar{Y} = \bar{\hat{Y}} + \bar{e} \quad ,$$

که با توجه به معادله ۴-۱۱ نتیجه می گیریم که

$$\sum e_i = 0 \quad \text{یا} \quad \bar{e} = 0 \quad . \quad (۴-۱۲)$$

معادله های ۴-۱۱ و ۴-۱۲ بر این دلالت می کند که میانگین مقادیر تخمین زده شده \hat{Y}_i با میانگین مقادیر Y_i برابر است؛ یعنی مجموع پسماندها برابر صفر است.

نتیجه حاصل از معادله ۴-۱۲ را می توان به روش دیگری نیز به دست آورد. از معادله اول در سیستم معادله های نرمال ۴-۴ داریم

$$\sum (Y_i - \hat{Y}_i) = 0 \quad , \quad \text{یا}$$

$$\sum_{i=1}^n e_i = 0 \quad \text{و یا} \quad \bar{e} = 0 \quad .$$

د) معادله دوم نرمال در سیستم معادله های ۴-۴ را میتوان چنین نوشت،

$$\frac{\partial \sum e_i^2}{\partial \hat{\beta}_1} = -2 \sum (Y_i - \hat{Y}_i) X_{1i} = 0 \quad ,$$

و در نتیجه

$$\sum X_{1i} e_i = 0 \quad . \quad (۴-۱۳)$$

به همین ترتیب از معادله سوم نرمال در سیستم معادله های ۴-۴ نتیجه می گیریم که

$$\sum X_{2i} e_i = 0 \quad . \quad (۴-۱۴)$$

اگر دو طرف معادله ۴-۲ را در e_t ضرب کنیم، خواهیم داشت

$$e_t \hat{Y}_t = e_t \hat{\alpha} + e_t \hat{\beta}_1 X_{1t} + e_t \hat{\beta}_2 X_{2t} ,$$

و برای مجموع مشاهدات موجود در نمونه، داریم

$$\sum e_t \hat{Y}_t = \hat{\alpha} \sum e_t + \hat{\beta}_1 \sum X_{1t} e_t + \hat{\beta}_2 \sum X_{2t} e_t .$$

با جایگزینی معادله‌های ۴-۱۲، ۴-۱۳ و ۴-۱۴ در معادله فوق نتیجه می‌گیریم که

$$\sum e_t \hat{Y}_t = 0 . \quad (4-15)$$

از معادله‌های ۴-۱۳، ۴-۱۴ و ۴-۱۵ می‌توان این سه نتیجه بسیار مهم را گرفت که در روش حداقل مربعات معمولی، پسماندها، از متغیرهای توضیحی X_{1t} و X_{2t} و تخمین متغیر درون‌زا \hat{Y}_t ، مستقل است.

۴-۳ ضریب تعیین

برای به دست آوردن فرمول ضریب تعیین در مدل رگرسیون با دو متغیر توضیحی به ترتیب زیر عمل می‌کنیم. می‌دانیم

$$Y_t = \hat{Y}_t + e_t .$$

دو طرف معادله فوق را برای تمام مقادیر t جمع کرده، بر n تقسیم می‌کنیم،

$$\bar{Y} = \bar{\hat{Y}} + \bar{e}_t .$$

دو طرف این دو معادله را از هم کم کرده خواهیم داشت

$$y_t = \hat{y}_t + e_t .$$

دو طرف معادله فوق را به توان ۲ می‌رسانیم و برای تمام مقادیر n جمع می‌کنیم،

$$\sum y_t^2 = \sum \hat{y}_t^2 + \sum e_t^2 + 2 \sum \hat{y}_t e_t .$$

جمله آخر صفر است؛ زیرا با توجه به معادله ۴-۱۵ و نیز معادله ۴-۱۲ داریم

$$\sum e_i \hat{Y}_i = \sum e_i (\hat{y}_i + \hat{Y}) = \sum e_i \hat{y}_i + \hat{Y} \sum e_i = 0 .$$

بنابراین می‌توان نوشت

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2 ,$$

یا:

$$TSS = ESS + RSS .$$

ملاحظه می‌شود که معادله فوق دقیقاً همان قضیه بسیار مهمی است که به صورت معادله ۴-۱۰ در رگرسیون ساده مطرح کردیم.

با استفاده از تعریف ضریب تعیین (معادله ۴-۲۹) می‌توان نوشت

$$R^2 = \frac{ESS}{TSS} = \frac{\sum \hat{y}_i^2}{\sum y_i^2} = 1 - \frac{\sum e_i^2}{\sum y_i^2} . \quad (4-16)$$

ساختار فرمول ضریب تعیین در رگرسیون چند متغیره دقیقاً همان فرمولی است که در رگرسیون ساده و به صورت معادله‌های ۴-۳۳ و ۴-۴۳ به دست آمد. البته ضریب تعیین را در رگرسیون چند متغیره با R^2 نشان می‌دهیم که از مقدار مشابه در رگرسیون ساده متمایز شود. از نظر محاسباتی، با توجه به معادله‌های ۴-۴۵ و ۴-۴۶ می‌دانیم در یک رگرسیون

ساده

$$ESS = \sum \hat{y}_i^2 = \hat{\beta} \sum x_i y_i ,$$

$$RSS = \sum e_i^2 = \sum y_i^2 - \hat{\beta} \sum x_i y_i .$$

اگر دقیقاً به همان روش به کار رفته در استخراج معادله‌های ۴-۴۵ و ۴-۴۶ عمل می‌کنیم، می‌توان تغییرات توضیح داده شده (ESS) و مجموع مربعات پسماند (RSS) را برای رگرسیون ۴-۱ به شرح زیر به دست آورد،

$$ESS = \sum \hat{y}_i^2 = \hat{\beta}_1 \sum x_{1i} y_i + \hat{\beta}_2 \sum x_{2i} y_i , \quad (4-17)$$

$$RSS = \sum e_i^2 = \sum y_i^2 - \hat{\beta}_1 \sum x_{1i} y_i - \hat{\beta}_2 \sum x_{2i} y_i . \quad (4-18)$$

در نتیجه، R^2 با استفاده از تعریف ۱-۲۹ به صورت زیر محاسبه می‌شود،

$$R^2 = \frac{\hat{\beta}_1 \sum x_{1t} y_t + \hat{\beta}_2 \sum x_{2t} y_t}{\sum y_t^2} \quad (۴-۱۹)$$

معمولاً از علامت‌گذارهای مختلفی برای R^2 استفاده می‌کنند. در بعضی موارد، R^2 را با ۳ اندیس معرفی می‌کنند. اندیس اول منعکس‌کننده متغیر درون‌زا است. بعد از این اندیس، یک ممیز گذاشته می‌شود و ۲ اندیس بعدی برای متغیرهای توضیحی به کار می‌رود. برای مثال، $R_{y/x_1, x_2}^2$ نشان می‌دهد که ضریب تعیین برای مدلی محاسبه می‌شود که Y متغیر درون‌زا و X_1 و X_2 متغیرهای برون‌زا هستند. برای سهولت، $R_{y/x_1, x_2}^2$ را به صورت $R_{y/x_1, x_2}^2$ یا حتی $R_{y/12}^2$ نیز نشان می‌دهند. منظور از «۱۲» متغیرهای اول و دوم توضیحی است. می‌توان $R_{y/x_1, x_2}^2$ را به صورت $R_{1/23}^2$ نیز نوشت. در اینجا «۱/۲۳» بدین صورت تفسیر می‌شود که «۱» منعکس‌کننده متغیر درون‌زا است و «۲۳» به معنای دومین و سومین متغیر موجود در مدل است که در واقع همان X_1 و X_2 هستند؛ بنابراین $R_{1/23}^2$ یعنی ضریب تعیین یک مدل رگرسیون که اولین متغیر آن (Y) تابعی از دومین و سومین متغیر آن (X_1, X_2) فرض شده است.

در پایان این بحث، به این نکته اشاره می‌شود که R^2 برابر مجذور ضریب همبستگی بین Y_t و \hat{Y}_t است،

$$R_{y/12}^2 = r_{y, \hat{y}}^2 = \frac{(\sum y_t \hat{y}_t)^2}{\sum y_t^2 \sum \hat{y}_t^2} \quad (۴-۲۰)$$

در مسأله ۱-۱۹ بند (۳)، این فرمول را برای رگرسیون ساده ثابت کردیم. برای رگرسیون چند متغیره نیز به ترتیبی مشابه عمل می‌کنیم. با استفاده از رابطه $y_t = \hat{y}_t + e_t$ داریم

$$\begin{aligned} \sum y_t \hat{y}_t &= \sum (\hat{y}_t + e_t) \hat{y}_t \\ &= \sum \hat{y}_t^2 + \sum e_t \hat{y}_t \end{aligned}$$

با توجه به معادله ۴-۱۵ به سهولت مشاهده می‌شود که جمله آخر ($\sum e_t \hat{y}_t$) برابر صفر

می شود، و در نتیجه

$$\sum y_i \hat{y}_i = \sum \hat{y}_i^2 \quad (4.21)$$

معادله ۴-۲۱ را در معادله ۴-۲۰ جایگزین می کنیم،

$$r_{y, \hat{y}}^2 = \frac{(\sum \hat{y}_i^2)^2}{\sum y_i^2 \sum \hat{y}_i^2} = \frac{\sum \hat{y}_i^2}{\sum y_i^2}$$

بر اساس معادله ۱-۲۹ می دانیم

$$R^2 = \frac{\sum \hat{y}_i^2}{\sum y_i^2}$$

در نتیجه $r_{y, \hat{y}}^2$ با R^2 برابر خواهد شد. لازم است توضیح دهیم که هر جا R^2 را بدون اندیس بنویسیم، منظور همان $r_{y, \hat{y}}^2$ یا صورت‌های مشابه آن است.

مثال ۴-۱ یک نمونه ۵ تایی از کارمندان یک شرکت تولیدی را به طور تصادفی انتخاب کرده ایم و داده‌های حقوق ماهانه (Y_i)، تعداد سالهای تحصیلی بعد از اخذ دیپلم (X_{1i}) و سالهای خدمت (X_{2i})، را در جدول ۴-۱ نشان داده‌ایم.

جدول ۴-۱ داده‌های حقوق، سالهای تحصیلی و خدمت

Y_i	X_{1i}	X_{2i}	$Y_i - \bar{Y}$	$X_{1i} - \bar{X}_1$	$X_{2i} - \bar{X}_2$
۳۰	۴	۱۰	۰	-۱	۰
۲۰	۳	۸	-۱۰	-۲	-۲
۳۶	۶	۱۱	۶	۱	۱
۲۴	۴	۹	-۶	-۱	-۱
۴۰	۸	۱۲	۱۰	۳	۲

۱. این مدل رگرسیون را با استفاده از آمار جدول ۴-۱ تخمین بزنید،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i$$

۲. R^2 را به دست آورید.

۳. درباره تخمین پارامترهایی که به دست آورده‌اید، بحث کنید.

۱. برای پاسخ به قسمت یک، ابتدا میانگینها را حساب می‌کنیم.

$$\bar{Y} = 30, \quad \bar{X}_1 = 0, \quad \bar{X}_2 = 10.$$

با توجه به اینکه مشاهدات، بسیار کم است نوشتن تمام محاسبات در جدول ضروری نیست. ابتدا مقادیر هر یک از متغیرها را بر حسب انحراف از میانگین، در ستونهای چهارم، پنجم و ششم جدول می‌نویسیم و سپس مقادیر لازم را محاسبه می‌کنیم.

$$\sum x_{1i}^2 = 16, \quad \sum x_{1i} y_i = 62,$$

$$\sum x_{2i}^2 = 10, \quad \sum x_{2i} y_i = 52,$$

$$\sum y_i^2 = 272, \quad \sum x_{1i} x_{2i} = 12.$$

با استفاده از معادلات ۱۰-۴ مقادیر $\hat{\beta}_1$ و $\hat{\beta}_2$ به ترتیب زیر محاسبه می‌شود،

$$\hat{\beta}_1 = \frac{\sum x_{2i}^2 \sum x_{1i} y_i - \sum x_{1i} x_{2i} \sum x_{2i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2},$$

$$\hat{\beta}_1 = \frac{10(62) - 12(52)}{16(10) - (12)^2} = \frac{-4}{16} = -0.25.$$

و برای $\hat{\beta}_2$ داریم

$$\hat{\beta}_2 = \frac{\sum x_{1i}^2 \sum x_{2i} y_i - \sum x_{1i} x_{2i} \sum x_{1i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2},$$

$$\hat{\beta}_2 = \frac{16(52) - 12(62)}{16(10) - (12)^2} = \frac{88}{16} = 5.5.$$

برای محاسبه $\hat{\alpha}$ ، معادله ۴-۶ را می‌نویسیم،

$$\bar{Y} = \hat{\alpha} + \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2 ,$$

در نتیجه داریم:

$$۳۰ = \hat{\alpha} + (-۰/۲۵)(۵) + (۵/۵)(۱۰) ,$$

$$\hat{\alpha} = -۲۳/۷۵ .$$

بدین ترتیب تخمین مدل رگرسیون مفروض عبارت است از

$$\hat{Y}_i = -۲۳/۷۵ - ۰/۲۵ X_{1i} + ۵/۵ X_{2i}$$

۲. برای محاسبه R^2 ابتدا معادله ۴-۱۹ را می‌نویسیم،

$$R_{y/12}^2 = R^2 = \frac{\hat{\beta}_1 \sum x_{1i} y_i + \hat{\beta}_2 \sum x_{2i} y_i}{\sum y_i^2}$$

که با استفاده از کمیت‌های محاسبه شده، خواهیم داشت

$$R_{y/12}^2 = R^2 = \frac{-۰/۲۵(۶۲) + ۵/۵(۵۲)}{۲۷۲} = \frac{۲۷۰/۵}{۲۷۲} = ۰/۹۹۴$$

۳. با مشاهده تخمین مدل رگرسیون می‌توان دو نکته را استنتاج کرد. اول اینکه

ضریب متغیر سابقه تحصیلی منفی است؛ یعنی افزایش سالهای تحصیلی، بر سطح حقوق تأثیر منفی دارد. نکته دوم این است که ضریب متعلق به سالهای خدمت از ضریب مشابه برای سابقه تحصیلی بسیار بزرگتر است. به عبارت دیگر، با استفاده از تخمین پارامترها می‌توان گفت که اگر سابقه تحصیلی را ثابت بگیریم، به ازای هر سال سابقه خدمت بیشتر، مبلغ a ۵/۵ به سطح حقوق اضافه خواهد شد که در آن a واحدی است که بر اساس آن حقوق ماهانه را اندازه‌گیری کرده‌ایم. از طرف دیگر، اگر دو نفر را با سالهای خدمت یکسان در نظر بگیریم، پیش‌بینی می‌شود فردی که یک سال بیشتر درس خوانده باشد، a ۰/۲۵ کمتر از دیگری حقوق خواهد گرفت. به این

ترتیب به نظر می‌رسد که نتایج تخمین پارامترهای متغیرهای برون‌زا اساساً رضایت‌بخش نیست.

تخمین جمله ثابت مدل ($\hat{\alpha} = -23/75$) را نیز باید ارزیابی کرد. $\hat{\alpha}$ در مدل اصلی، نشان‌دهنده حقوق ماهانه فردی دیپلمه است که هیچ سابقه خدمت نیز ندارد. پدیهی است حقوق نمی‌تواند منفی شود، اما در تخمین مدل، مقدار $\hat{\alpha}$ منفی شده است؛ یعنی بر اساس پیش‌بینی مدل، حقوق یک کارمند جدید منفی است؛ بنابراین تخمین جمله ثابت مدل نیز چندان جالب نیست.

نتیجه می‌گیریم که نمونه ۵ تایی ما، نمونه خوبی از جامعه کارمندان شرکت مفروض نیست. به نظر می‌رسد که این نمونه از گروه خاصی از کارمندان گرفته شده باشد؛ زیرا ملاحظه می‌شود که سالهای خدمت اعضای متعلق به این نمونه، بین ۸ تا ۱۲ سال است. در چنین حالتی قاعدتاً نمی‌توان خارج از این محدوده را برای متغیر درون‌زا پیش‌بینی کرد و تخمین مدل رگرسیون مفروض نمی‌تواند به این سؤال پاسخ دهد که حقوق یک کارمند تازه‌وارد چه مقدار خواهد بود.

برای تحلیل بیشتر در نتایج تخمین، دو رگرسیون ساده نیز برای این رگرسیون چندمتغیره می‌سازیم. ابتدا برای حقوق ماهانه (Y_i) و سالهای تحصیلی (X_{1i}) و سپس برای حقوق ماهانه و سالهای خدمت (X_{2i})، مدل‌های جداگانه‌ای تخمین می‌زنیم. با استفاده از داده‌های موجود در جدول ۱-۴ تخمین این مدل‌ها به شرح زیر خواهد بود:

$$\hat{Y}_i = 10/625 + 3/875 X_{1i}, \quad R^2 = 0/883$$

$$\hat{Y}_i = -22 + 0/2 X_{2i}, \quad R^2 = 0/994$$

تخمین مدل رگرسیون ساده اول نشان می‌دهد که به ازای هر سال تحصیل بیشتر، مبلغ $3/875$ به حقوق ماهانه اضافه خواهد شد، که در آن a واحد اندازه‌گیری حقوق است. اما در مدل رگرسیون چندمتغیره دیدیم که اضافه شدن سالهای تحصیل هیچگونه کمکی به افزایش سطح حقوق ندارد. این نتیجه می‌تواند بر این حقیقت دلالت کند که عدم حضور متغیر توضیحی سالهای خدمت در مدل رگرسیون ساده اول، باعث شده است که

ما به نتایج کاملاً غلطی از تأثیر متغیر «سوابق تحصیلی» بر سطح حقوق ماهانه برسیم. قدم بعدی در تفسیر تخمین پارامترهای مدل مفروض این است که واریانس تخمینها را محاسبه و فرضیه‌های مختلف را نسبت به آنها آزمون کنیم. اما قبلاً باید به بحث آزمون فرضیه در یک مدل رگرسیون با دو متغیر توضیحی پردازیم.

۴-۴ آزمون فرضیه

مدل رگرسیون ۴-۱ و تخمین آن (۴-۲) را یک بار دیگر می‌نویسیم،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t,$$

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta}_1 X_{1t} + \hat{\beta}_2 X_{2t}.$$

می‌دانیم $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ متغیرهای تصادفی هستند. دقیقاً مانند بحث قسمت ۲-۲ (خصوصیات آماری تخمین زنده‌های حداقل مربعات معمولی) می‌توان نشان داد که $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ توابع خطی از U_t هستند. اگر فرض نرمال بودن تابع توزیع احتمال U_t را بپذیریم، نتیجه می‌گیریم که تخمین پارامترهای مدل ۴-۱ نیز توزیع نرمال دارد،

$$\hat{\alpha} \sim N [E(\hat{\alpha}), \text{Var}(\hat{\alpha})],$$

$$\hat{\beta}_1 \sim N [E(\hat{\beta}_1), \text{Var}(\hat{\beta}_1)],$$

$$\hat{\beta}_2 \sim N [E(\hat{\beta}_2), \text{Var}(\hat{\beta}_2)].$$

با توجه به اینکه اولاً، روش محاسبه میانگین و واریانس تخمین زنده‌های $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ دقیقاً مشابه مدل رگرسیون ساده است که در قسمت ۲-۲ ارائه شد و ثانیاً اثبات نحوه استخراج فرمولهای میانگین و واریانس تخمین پارامترها در حالت کلی با k پارامتر در فصل ششم مطرح خواهد شد؛ بنابراین کافی است در این قسمت فقط به ذکر نتایج مدل رگرسیون ۴-۱ می‌پردازیم.

می‌توان نشان داد که تخمین زنده‌های $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ ناورب هستند؛

$$E(\hat{\alpha}) = \alpha ,$$

$$E(\hat{\beta}_1) = \beta_1 , \quad (4.22)$$

$$E(\hat{\beta}_2) = \beta_2 .$$

همچنین اگر r_{12} به معنای ضریب همبستگی بین X_{1i} و X_{2i} باشد، از معادله

۱.۳۲ می‌دانیم که

$$r_{12}^2 = \frac{(\sum x_{1i} x_{2i})^2}{\sum x_{1i}^2 \sum x_{2i}^2} ,$$

می‌توان نشان داد که واریانس و کواریانس $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ به شرح زیر خواهد بود،

$$\text{Var}(\hat{\alpha}) = \frac{\sigma^2}{n} + \bar{X}_1^2 \text{Var}(\hat{\beta}_1) + 2 \bar{X}_1 \bar{X}_2 \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) + \bar{X}_2^2 \text{Var}(\hat{\beta}_2) \quad (4.23)$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{(1 - r_{12}^2) \sum x_{1i}^2} , \quad (4.24)$$

$$\text{Var}(\hat{\beta}_2) = \frac{\sigma^2}{(1 - r_{12}^2) \sum x_{2i}^2} , \quad (4.25)$$

$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = \frac{-\sigma^2 r_{12}}{(1 - r_{12}^2) \sum x_{1i} \sum x_{2i}} , \quad (4.26)$$

$$\text{Cov}(\hat{\alpha}, \hat{\beta}_1) = -[\bar{X}_1 \text{Var}(\hat{\beta}_1) + \bar{X}_2 \text{Cov}(\hat{\beta}_1, \hat{\beta}_2)] \quad (4.27)$$

$$\text{Cov}(\hat{\alpha}, \hat{\beta}_2) = -[\bar{X}_1 \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) + \bar{X}_2 \text{Var}(\hat{\beta}_2)] \quad (4.28)$$

روش آزمون فرضیه و ساختن فاصله‌های اطمینان برای تخمین پارامترهای

مدل ۴-۱ دقیقاً مانند رگرسیون ساده است که در قسمت ۲-۲ بررسی شد. در اینجا

تنها به ذکر چند نکته اکتفا می‌کنیم،

۱. از معادله‌های ۴-۲۳، ۴-۲۴ و ۴-۲۵ مشاهده می‌شود که به ازای افزایش ضریب همبستگی بین X_{11} و X_{21} ، مقادیر واریانس زیاد می‌شود و از دقت تخمینها کاسته خواهد شد. اگر همبستگی بین X_{11} و X_{21} خیلی زیاد شود؛ یعنی r^2 به سمت یک میل کند، عملاً واریانس تخمینها به قدری زیاد می‌شود که دیگر تخمینها قابل اطمینان نخواهند بود. هنگامی که همبستگی بین X_{11} و X_{21} کامل شود؛ یعنی $r^2 = 1$ باشد، آنگاه تخمین پارامترها ممکن نخواهد بود. به همین ترتیب اگر بین X_{11} و X_{21} مطلقاً رابطه‌ای نباشد؛ یعنی $r^2 = 0$ ، در آن صورت مقادیر واریانسها دقیقاً برابر مقادیری است که از رگرسیون ساده Y_1 بر X_{11} و نیز Y_1 بر X_{21} به دست می‌آید. اینها مسائلی است که به طور خلاصه در مقدمه قسمت ۴-۲ و در چهارچوب یک مثال مطرح شد و دوباره در فصل ششم مطرح خواهد شد.

۲. اگر مجموع مربعات پسماند برای مدل ۴-۱ را با RSS نشان دهیم، در آن صورت مطابق آنچه در معادله ۲-۳۸ گفتیم، کمیت $\frac{RSS}{\sigma^2}$ ، دارای توزیع χ^2 با $(n-3)$ درجه آزادی است. توجه داریم که در مدل ۴-۱ سه پارامتر، تخمین زده می‌شود و در محاسبه RSS سه درجه آزادی از دست می‌دهیم؛ بنابراین

$$\frac{RSS}{\sigma^2} \sim \chi^2 (n-3) .$$

دقیقاً مشابه نامسای ۲-۴۴، می‌توان برای σ^2 نیز فاصله اطمینان ساخت؛ برای مثال، یک فاصله اطمینان ۹۵ درصد برای σ^2 عبارت خواهد بود از

$$\frac{(n-3) \hat{\sigma}^2}{\chi^2_{.975}} < \sigma^2 < \frac{(n-3) \hat{\sigma}^2}{\chi^2_{.025}} \quad (4-29)$$

۳. برای به دست آوردن تخمین واریانس U_1 ، مقدار مجموع مربعات پسماند را بر درجات آزادی آن تقسیم می‌کنیم. مشابه استدلالی که برای معادله ۲-۲۸ ارائه شده می‌توان ثابت کرد که بدین ترتیب تخمین نااریبی از واریانس U_1 خواهیم داشت؛

$$\hat{\sigma}_u^2 = \frac{RSS}{(n-3)} \quad (4-30)$$

۴. با جایگزینی $\hat{\sigma}^2$ از معادله ۴-۳۰ در معادله‌های ۴-۲۳ و ۴-۲۶ به تخمین واریانس پارامترها می‌رسیم. اگر جذر واریانسها را حساب کنیم، انحراف تخمین پارامترها را به دست خواهیم آورد. با داشتن میانگین و انحراف معیار پارامترها، به راحتی می‌توان آنها را استاندارد کرد. پارامترهای استاندارد زیر هر یک توزیع t با $(n-3)$ درجه آزادی دارند:

$$t = \frac{\hat{\alpha} - \alpha}{SE(\hat{\alpha})}, \quad (4.31)$$

$$t = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)}, \quad (4.32)$$

$$t = \frac{\hat{\beta}_2 - \beta_2}{SE(\hat{\beta}_2)}, \quad (4.33)$$

مثال ۴-۲ تابع تولید زیر را در نظر می‌گیریم:

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i,$$

که در آن $U_i \sim IN(0, \sigma^2)$ و

Y_i = لگاریتم تولید ،

X_{1i} = لگاریتم کار ،

X_{2i} = لگاریتم سرمایه ،

X_{1i} و X_{2i} غیر تصادفی است. این محاسبات بر اساس یک نمونه شامل ۲۳ شرکت مختلف تولیدی به دست آمده است.

$$\bar{X}_1 = 10, \quad \sum x_{1i}^2 = 12, \quad \sum x_{1i} y_i = 10,$$

$$\bar{X}_2 = 0, \quad \sum x_{1i} x_{2i} = 8, \quad \sum x_{2i} y_i = 8,$$

$$\bar{Y} = 12, \quad \sum x_{1i}^2 = 12, \quad \sum y_i^2 = 10.$$

الف) $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ و مقادیر انحراف معیار آنها را محاسبه کنید و تخمین مدل را بنویسید.

ب) فاصله اطمینان ۹۵ درصد برای α ، β_1 و β_2 را به دست آورید و فرضیه‌های $\beta_1 = 1$ و $\beta_2 = 0$ را به صورت جداگانه در سطح معنی‌دار ۵ درصد آزمون کنید.

الف) از معادله‌های ۴-۱۰ استفاده می‌کنیم و مقادیر $\hat{\beta}_1$ و $\hat{\beta}_2$ را به دست می‌آوریم،

$$\hat{\beta}_1 = \frac{\sum x_{2i}^2 \sum x_{1i} y_i - \sum x_{1i} x_{2i} \sum x_{2i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2}$$

$$= \frac{12(10) - 8(8)}{12(12) - (8)^2} = \frac{56}{80} = 0.7$$

$$\hat{\beta}_2 = \frac{\sum x_{1i}^2 \sum x_{2i} y_i - \sum x_{1i} x_{2i} \sum x_{1i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2}$$

$$= \frac{12(8) - 8(10)}{12(12) - (8)^2} = \frac{16}{80} = 0.2$$

برای محاسبه $\hat{\alpha}$ ، معادله ۴-۶ را می‌نویسیم،

$$\bar{Y} = \hat{\alpha} + \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2$$

در نتیجه داریم:

$$12 = \hat{\alpha} + 0.7(10) + 0.2(5)$$

و بدین ترتیب

$$\hat{\alpha} = 4$$

برای تخمین مقادیر انحراف معیار، ابتدا $\hat{\sigma}^2$ را تخمین می‌زنیم. برای این کار ضروری است که ابتدا R^2 را به دست آوریم و سپس به کمک آن RSS (مجموع مربعات پسماند) را محاسبه کنیم تا $\hat{\sigma}^2$ به دست آید. از فرمول ۴-۱۹ استفاده می‌کنیم،

$$R_{y/12}^2 = \frac{\hat{\beta}_1 \sum x_{1i} y_i + \hat{\beta}_2 \sum x_{2i} y_i}{\sum y_i^2} ,$$

$$= \frac{0.7(10) + 0.2(8)}{10} = 0.86 .$$

برای محاسبه مجموع مربعات پسماند از معادله ۴-۱۸ استفاده می‌کنیم،

$$RSS = \sum e_i^2 = \sum y_i^2 - \hat{\beta} \sum x_{1i} y_i - \hat{\beta}_2 \sum x_{2i} y_i ,$$

$$= 10 - 0.7(10) - 0.2(8) = 1/4 .$$

با استفاده از معادله ۴-۳۰ داریم

$$\hat{\sigma}_u^2 = \frac{RSS}{(n-3)} = \frac{1/4}{20-3} = 0.07 .$$

با به دست آوردن تخمین واریانس U_i ابتدا باید از معادله‌های ۴-۲۴ و ۴-۲۵ استفاده کرد تا واریانس $\hat{\beta}_1$ و $\hat{\beta}_2$ به دست آید. اما لازمه این کار محاسبه مجذور ضریب همبستگی بین X_{1i} و X_{2i} یعنی r_{12}^2 است. از معادله ۴-۳۲ می‌دانیم که

$$r_{12}^2 = \frac{(\sum x_{1i} x_{2i})^2}{\sum x_{1i}^2 \sum x_{2i}^2} ,$$

$$= \frac{(8)^2}{12(12)} = \frac{64}{144} ,$$

بنابراین

$$\text{Var}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum x_{1i}^2 (1 - r_{12}^2)} ,$$

$$= \frac{0.07}{12 \left(1 - \frac{74}{144}\right)} = \frac{0.07}{\frac{80}{12}} = \frac{7}{100} \left(\frac{3}{20}\right) = \frac{21}{2000}$$

$$\text{Var}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{\sum x_{1i}^2 (1 - r_{12}^2)}$$

$$= \frac{0.07}{12 \left(1 - \frac{74}{100}\right)} = \frac{21}{2000}$$

برای تخمین واریانس $\hat{\alpha}$ به تخمین $\text{Cov}(\hat{\beta}_1, \hat{\beta}_2)$ نیاز داریم. از معادله ۴-۲۶ استفاده می‌کنیم،

$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = \frac{-\hat{\sigma}^2 r_{12}^2}{\sum x_{1i} x_{2i} (1 - r_{12}^2)}$$

$$= \frac{-0.07 \left(\frac{74}{144}\right)}{8 \left(\frac{80}{144}\right)} = \frac{-0.07(74)}{64} = -0.007$$

و با توجه به $\bar{X}_1 = 10$ و $\bar{X}_2 = 0$ ، واریانس $\hat{\alpha}$ به صورت زیر محاسبه می‌شود،

$$\text{Var}(\hat{\alpha}) = \frac{\hat{\sigma}^2}{n} + \bar{X}_1^2 \text{Var}(\hat{\beta}_1) + 2 \bar{X}_1 \bar{X}_2 \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) + \bar{X}_2^2 \text{Var}(\hat{\beta}_2)$$

$$= \frac{0.07}{23} + 100(0.07) \left(\frac{3}{20}\right) + 2(0)(10)(0.07) \left(\frac{1}{10}\right)$$

$$+ 20(0.07) \left(\frac{3}{20}\right)$$

$$= 0.07 \left[\frac{1}{23} + \frac{300}{20} + 20 + \frac{70}{20} \right] = 0.07(17.935)$$

$$= 0.125545$$

برای محاسبه مقادیر انحراف معیار $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ از واریانس آنها جذر می‌گیریم،

$$SE(\hat{\alpha}) = \sqrt{0/615545} = 0/78 ,$$

$$SE(\hat{\beta}_1) = \sqrt{\frac{21}{2000}} = 0/102 ,$$

$$SE(\hat{\beta}_2) = 0/102 .$$

تخمین مدل رگرسیون مفروض به شرح زیر خواهد بود،

$$\hat{Y}_i = \varepsilon + 0/7 X_{1i} + 0/2 X_{2i} ,$$

(0/78) (0/102) (0/102)

$$R^2 = 0/86 .$$

ب) مقدار t از جدول t با $23 - 3 = 20$ درجه آزادی و در سطح معنی‌دار ۵ درصد برابر است با $t = \pm 2/086$. فاصله‌های اطمینان ۹۵ درصد برای α ، β_1 و β_2 به شرح زیر است،

$$\hat{\alpha} - t_{\alpha} SE(\hat{\alpha}) < \alpha < \hat{\alpha} + t_{\alpha} SE(\hat{\alpha}) ,$$

$$\varepsilon - 2/086 (0/78) < \alpha < \varepsilon + 2/086 (0/78) ,$$

$$2/37 < \alpha < 0/63 .$$

برای β_1 داریم:

$$\hat{\beta}_1 - t_{\alpha} SE(\hat{\beta}_1) < \beta_1 < \hat{\beta}_1 + t_{\alpha} SE(\hat{\beta}_1) ,$$

$$0/7 - 2/086 (0/102) < \beta_1 < 0/7 + 2/086 (0/102) ,$$

$$0/49 < \beta_1 < 0/91 .$$

برای β_2 خواهیم داشت

$$\hat{\beta}_2 - t_{\alpha} SE(\hat{\beta}_2) < \beta_2 < \hat{\beta}_2 + t_{\alpha} SE(\hat{\beta}_2) ,$$

$$0/2 - 2/0.86 (0/1.02) < \beta_1 < 0/2 + 2/0.86 (0/1.02) ,$$

$$-0/01 < \beta_1 < 0/41 .$$

فاصله اطمینان برای σ^2 از توزیع χ^2 به دست می آید. از معادله ۴-۲۹ داریم:

$$\frac{(n-3) \hat{\sigma}^2}{\chi^2_{0/95}} < \sigma^2 < \frac{(n-3) \hat{\sigma}^2}{\chi^2_{0/05}} .$$

مقادیر به دست آمده از جدول χ^2 با $(n-3) = 23 - 3 = 20$ درجه آزادی عبارت است از

$$\chi^2_{0/05} = 9/59 , \quad \chi^2_{0/95} = 34/2$$

در نتیجه داریم:

$$\Pr \left[\frac{20 (0/07)}{34/2} < \sigma^2 < \frac{20 (0/07)}{9/59} \right] = 0/90 .$$

به عبارت دیگر با ۹۰ درصد احتمال، σ^2 بین دو مقدار زیر قرار می گیرد،

$$\frac{20 (0/07)}{34/2} < \sigma^2 < \frac{20 (0/07)}{9/59} ,$$

یا

$$0/041 < \sigma^2 < 0/145 .$$

فرضیه $\beta_1 = 1$ در سطح معنی دار ۵ درصد رد می شود؛ زیرا β_1 در فاصله اطمینان ۹۰ درصد قرار نمی گیرد. اما فرضیه $\beta_1 = 0$ تقطه ای است که در فاصله اطمینان ۹۰ درصد برای β_1 قرار دارد.

۴-۵ پیش بینی

مسئله پیش بینی در مدل های رگرسیون با دو متغیر توضیحی شبیه پیش بینی در مدل های رگرسیون ساده است؛ با این تفاوت که برای محاسبه خطای پیش بینی

باید از واریانس و کوواریانس تمام پارامترها استفاده کرد.
مدل رگرسیون ۴-۱ و تخمین آن ۴-۲ را یک بار دیگر می‌نویسیم،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t ,$$

$$\hat{Y}_t = \hat{\alpha} + \hat{\beta}_1 X_{1t} + \hat{\beta}_2 X_{2t} .$$

فرض کنید می‌خواهیم به ازای مقادیر معینی از متغیرهای برون‌زا، مقدار متغیر درون‌زا را پیش‌بینی کنیم. دوره‌ای را که می‌خواهیم پیش‌بینی کنیم، $t = \bar{t}$ می‌نامیم. سؤال این است که به ازای $X_{1\bar{t}} = X_{1\bar{t}}$ و $X_{2\bar{t}} = X_{2\bar{t}}$ اولاً مقدار $\hat{Y}_{\bar{t}}$ چقدر است و ثانیاً چگونه می‌توان یک فاصله اطمینان ۹۵ درصد برای آن ساخت؟
مدل ۴-۱ برای دوره \bar{t} عبارت است از

$$Y_{\bar{t}} = \alpha + \beta_1 X_{1\bar{t}} + \beta_2 X_{2\bar{t}} + U_{\bar{t}} ,$$

که تخمین آن را می‌توان چنین نوشت،

$$\hat{Y}_{\bar{t}} = \hat{\alpha} + \hat{\beta}_1 X_{1\bar{t}} + \hat{\beta}_2 X_{2\bar{t}} . \quad (4.34)$$

مانند معادله ۳-۵ خطای تخمین را برای معادله فوق تعریف می‌کنیم،

$$e_t = Y_t - \hat{Y}_t , \quad (4.35)$$

خواهیم داشت

$$e_{\bar{t}} = U_{\bar{t}} - (\hat{\alpha} - \alpha) - (\hat{\beta}_1 - \beta_1) X_{1\bar{t}} - (\hat{\beta}_2 - \beta_2) X_{2\bar{t}} . \quad (4.36)$$

می‌دانیم توابع توزیع احتمال $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ و $U_{\bar{t}}$ نرمال هستند و چون $e_{\bar{t}}$ یک تابع خطی از آنهاست توزیع احتمال نرمال خواهد داشت. برای پیدا کردن میانگین $e_{\bar{t}}$ باید از معادله ۴-۳۶ امید ریاضی بگیریم،

$$E(e_{\bar{t}}) = E(U_{\bar{t}}) - E(\hat{\alpha} - \alpha) - X_{1\bar{t}} E(\hat{\beta}_1 - \beta_1) - X_{2\bar{t}} E(\hat{\beta}_2 - \beta_2) .$$

با توجه به اینکه $\hat{\alpha}$ ، $\hat{\beta}_1$ و $\hat{\beta}_2$ تخمینهای ناریب هستند و نیز با توجه به $E(U_i) = 0$ خواهیم داشت

$$E(e_i) = 0 \quad (4.37)$$

با گرفتن امید ریاضی از معادله ۴-۳۵ و با توجه به معادله ۴-۳۷ خواهیم داشت

$$E(\hat{Y}_i) = E(Y_i) \quad (4.38)$$

بر اساس معادله ۴-۳۸ می توان گفت که \hat{Y}_i یک پیش بینی کننده ناریب است. باید دقت کنیم که در اینجا ناریبی را به صورت $E(\hat{Y}_i)$ برابر با امید ریاضی Y_i گرفته ایم، Y_i یک متغیر تصادفی است؛ در حالی که امید ریاضی Y_i دیگر تصادفی نیست، بلکه مقدار معنی دارد.

برای محاسبه واریانس e_i ، به ترتیبی کاملاً مشابه روش استخراج معادله ۳-۸ عمل می کنیم. خواهیم داشت

$$\begin{aligned} \text{Var}(e_i) = \sigma_u^2 \left(1 + \frac{1}{n}\right) + (X_{1i} - \bar{X}_1)^2 \text{Var}(\hat{\beta}_1) + 2(X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2) \\ \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) + (X_{2i} - \bar{X}_2)^2 \text{Var}(\hat{\beta}_2) \end{aligned} \quad (4.39)$$

با داشتن میانگین و واریانس e_i ، به راحتی می توان آن را استاندارد کرد. اگر به جای σ^2 در معادله ۴-۳۹ مقدار تخمین آن را از معادله ۴-۳۰ قرار دهیم،

$$\hat{\sigma}_u^2 = \frac{\text{RSS}}{(n-3)}$$

به تخمین واریانس e_i می رسمیم. بدین ترتیب مقدار استاندارد شده e_i دارای توزیع t با $(n-3)$ درجه آزادی خواهد بود،

$$\frac{e_i - E(e_i)}{\text{SE}(e_i)} \sim t(n-3)$$

با توجه به معادله های ۴-۳۵ و ۴-۳۷ خواهیم داشت

$$\frac{Y_i - \hat{Y}_i}{\text{SE}(e_i)} \sim t(n-3)$$

در نتیجه فاصله اطمینان ۹۵ درصد برای Y_f عبارت است از

$$\hat{Y}_f - t_{\alpha/2} SE(e_f) < Y_f < \hat{Y}_f + t_{\alpha/2} SE(e_f) , \quad (۴-۴۰)$$

که از نظر ساختار ریاضی دقیقاً مانند نامساوی ۳-۱۱ در مدل‌های رگرسیون ساده است. مشاهده می‌شود که تفاوت معادله‌های ۴-۴۰ و ۳-۱۱ فقط در فرمول $SE(e_f)$ است.

به ترتیبی کاملاً مشابه با رگرسیون ساده و معادله ۳-۱۸، می‌توان برای $E(Y_f)$ نیز در مدل رگرسیون با دو متغیر توضیحی، فاصله اطمینان ساخت. با توجه به اینکه حالت کلی این مسأله را در فصل نهم مطرح خواهیم کرد، به نظر می‌رسد، تبیین بیشتر آن در این قسمت ضروری نیست.

مثال ۴-۳ مدل رگرسیون، موضوع مثال ۴-۲ را یک بار دیگر در نظر می‌گیریم:

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i .$$

می‌دانیم تخمین این مدل به صورت زیر به دست آمده است:

$$\hat{Y}_i = 4 + 0.7 X_{1i} + 0.2 X_{2i} .$$

همچنین نتایج محاسباتی زیر را از مثال ۴-۲ داریم،

$$\bar{X}_1 = 10 , \quad \text{Var}(\hat{\beta}_1) = \frac{3}{20} (0.07) ,$$

$$\bar{X}_2 = 5 , \quad \text{Var}(\hat{\beta}_2) = \frac{3}{20} (0.07) ,$$

$$\hat{\sigma}_u^2 = 0.07 , \quad \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = -\frac{1}{10} (0.07) .$$

اولاً، مقدار Y_i را به ازای $X_{1i}=12$ و $X_{2i}=7$ پیش‌بینی کنید.

ثانیاً، برای مقدار پیش‌بینی شده Y_i یک فاصله اطمینان ۹۵ درصد بسازید.

برای قسمت اول تخمین مدل را برای دوره t می‌نویسیم:

$$\hat{Y}_t = 4 + 0.7 X_{1t} + 0.2 X_{2t} ,$$

که به ازای $X_{1f}=12$ و $X_{2f}=7$ خواهیم داشت

$$\hat{Y}_f = 4 + 0.7(12) + 0.2(7) = 13/8 .$$

برای قسمت دوم، نامساوی $4-4.0$ را برای فاصله اطمینان Y_f می نویسیم،

$$\hat{Y}_f - t_{\alpha/2} SE(e_f) < Y_f < \hat{Y}_f + t_{\alpha/2} SE(e_f) .$$

بنابراین باید ابتدا $Var(e_f)$ را حساب کنیم. با توجه به

$$X_{1f} - \bar{X} = 12 - 10 = 2 ,$$

$$X_{2f} - \bar{X} = 7 - 0 = 7 ,$$

و با استفاده از معادله $4-39$ خواهیم داشت

$$\begin{aligned} Var(e_f) &= 0.7 \left(1 + \frac{1}{23} \right) + (4) \left(\frac{3}{20} \right) 0.7 - 2(2)(7) \left(\frac{1}{10} \right) 0.7 + 4 \left(\frac{3}{20} \right) 0.7 , \\ &= 0.7 \left(1 + \frac{1}{23} \right) + 4(0.7) \left(\frac{3}{20} - \frac{2}{10} + \frac{3}{20} \right) = 0.101 . \end{aligned}$$

با گرفتن جذر از $Var(e_f)$ ، به انحراف معیار e_f خواهیم رسید،

$$SE(e_f) = \sqrt{0.101} = 0.318 .$$

بدین ترتیب نامساوی $4-4.0$ را می توان با توجه به $t = \pm 2.086$ چنین نوشت،

$$13/8 - 2.086(0.318) < Y_f < 13/8 + 2.086(0.318) ,$$

$$13/8 - 0.66 < Y_f < 13/8 + 0.66 ,$$

یا

$$13/14 < Y_f < 14/66 .$$

دقت پیش‌بینی

فاصله اطمینان برای Y_f در مدل رگرسیون ساده، یعنی رابطه ۳-۱۰ را یک بار دیگر می‌نویسیم:

$$\hat{Y}_f - t_{\alpha/2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}} < Y_f < \hat{Y}_f + t_{\alpha/2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(X_f - \bar{X})^2}{\sum x_i^2}}.$$

در بحث مربوط به نمودار ۳-۱ دیدیم که به ازای $X_f = \bar{X}$ ، فاصله اطمینان برای Y_f در کمترین حد ممکن است؛ زیرا به ازای این مقدار، جمله $\frac{(X_f - \bar{X})^2}{\sum x_i^2}$ حذف خواهد شد؛ به عبارت دیگر، هر چه X_f از \bar{X} فاصله بگیرد، واریانس و انحراف معیار e_f بیشتر می‌شود و از دقت پیش‌بینی فاصله اطمینان برای Y_f کاسته خواهد شد.

نتیجه فوق در مورد مدل‌هایی که بیش از یک متغیر توضیحی دارد، معمولاً صادق نیست. علت آن، وجود جمله حاوی کوواریانس $\hat{\beta}_1$ و $\hat{\beta}_2$ در فرمول محاسبه واریانس e_f است. در معادله ۴-۳۹ می‌بینیم که واریانس e_f نه فقط تابعی از واریانس U_f ، واریانس $\hat{\beta}_1$ و واریانس $\hat{\beta}_2$ است، بلکه از کوواریانس $\hat{\beta}_1$ و $\hat{\beta}_2$ نیز تبعیت می‌کند. معمولاً علامت کوواریانس، که با توجه به فرمول ۴-۲۶ منفی است، و ضریب این کوواریانس، که حاصل ضرب $(X_{1f} - \bar{X}_1)$ و $(X_{2f} - \bar{X}_2)$ است، عوامل مهمی در نادرستی نتیجه حاصل از رگرسیون ساده به رگرسیون چند متغیره است.

برای تبیین بیشتر این نکته، به معادله ۴-۳۹ مراجعه می‌کنیم. اگر X_{1f} از \bar{X}_1 فاصله بگیرد، چون اختلاف آن از \bar{X}_1 به توان ۲ می‌رسد، در نهایت، نتیجه یک عدد مثبت خواهد بود. بنابراین فرقی نمی‌کند که جهت انحراف X_{1f} از \bar{X}_1 چه باشد. برای هر گونه انحراف X_{1f} از \bar{X}_1 نیز این بحث دقیقاً معتبر خواهد بود. اما برای ضریب کوواریانس $\hat{\beta}_1$ و $\hat{\beta}_2$ این گونه نیست. فرض می‌کنیم $\text{Cov}(\hat{\beta}_1, \hat{\beta}_2)$ یک عدد منفی است. اگر X_{1f} به \bar{X}_1 چنان نزدیک شود که نتیجه $(X_{1f} - \bar{X}_1)$ یک عدد منفی شود، حاصل ضرب $(X_{1f} - \bar{X}_1)$ و $\text{Cov}(\hat{\beta}_1, \hat{\beta}_2)$ یک عدد مثبت شده و واریانس e_f بشدت افزایش خواهد یافت؛ در حالی که با نزدیک شدن X_{1f} به میانگین \bar{X}_1 انتظار ما این بود که واریانس e_f کاهش یابد. بدین ترتیب نتیجه می‌گیریم که در مدل‌های رگرسیون چند متغیره نمی‌توان

به همان نتیجه ساده و قابل توجهی برسیم که در مبحث دقت پیش‌بینی در رگرسیون‌های ساده حاصل شد. در این مورد مثالی می‌آوریم.

مثال ۴-۴ در مثال ۴-۳ دیدیم که \bar{X}_Y برابر ۵ بود و به ازای $X_{Yf} = 7$ مقدار واریانس e_f را برابر $0/101$ تخمین زدیم. با توجه به معادله ۴-۳۹ مشاهده می‌شود که $(X_{Yf} - \bar{X}_Y)$ در مثال ۴-۳ با $2 = (7.5)$ برابر است و هنگامی که به صورت ضریب واریانس $\hat{\beta}_1$ به توان ۲ برسد، برابر ۴ خواهد شد.

حال فرض کنید که می‌خواهیم به ازای $X_{Yf} = 3$ پیش‌بینی کنیم که مقدار $(X_{Yf} - \bar{X}_Y)$ در معادله ۴-۳۹ برابر $-2 = (3.5)$ می‌شود که مربع آن یا ضریب واریانس $\hat{\beta}_1$ همان مقدار ۴ را خواهد داشت. اما تأثیر مقدار جدید X_{Yf} در ضریب کوواریانس، کاملاً با مثال قبلی متفاوت است. مقدار این ضریب در مثال ۴-۳ عبارت بود از

$$2 (X_{Yf} - \bar{X}_Y) (X_{Yf} - \bar{X}_Y) = 2 (12 - 10) (7 - 5) = 2 (4) ,$$

در حالی که در این مثال برابر است با

$$2 (X_{Yf} - \bar{X}_Y) (X_{Yf} - \bar{X}_Y) = 2 (12 - 10) (3 - 5) = -2 (4) .$$

با توجه به مقدار کوواریانس $\hat{\beta}_1$ و $\hat{\beta}_2$ خواهیم داشت

$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = 8 (-0/007) = -0/056 , \quad \text{در مثال ۴-۳}$$

$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = -8 (-0/007) = 0/056 , \quad \text{در مثال ۴-۴}$$

و بدیهی است که این دو مقدار، کاملاً با یکدیگر تفاوت دارند و تأثیر چشمگیری بر $\text{Var}(e_f)$ خواهند داشت؛ به عبارت دقیقتر، برای این مثال داریم

$$\text{Var}(e_f) = 0/07 \left(1 + \frac{1}{23}\right) + 4 (0/07) \left(\frac{3}{20}\right) + 0/056 + 4 \left(\frac{3}{20}\right) 0/07 = 0/213 ,$$

که از مقدار قبلی، یعنی $0/101$ ، بزرگتر است؛ زیرا به جای یک جمله منفی شامل یک جمله مثبت است.

۴-۶ ضرایب همبستگی و ضرایب تعیین*
مدل رگرسیون ۴-۱ را یک بار دیگر مشاهده کنید،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t .$$

ضریب همبستگی Y_t و X_{1t} را با ρ_{y_1} ، ضریب همبستگی Y_t و X_{2t} را با ρ_{y_2} و ضریب همبستگی X_{1t} و X_{2t} را با ρ_{12} نشان می‌دهیم. به ρ_{y_1} ، ρ_{y_2} و ρ_{12} «ضرایب همبستگی ساده»^۱ می‌گوییم. همچنین سه تخمین زیر را تعریف می‌کنیم،

$$\hat{\beta}_{y_1} = \text{تخمین شیب رگرسیون } Y_t \text{ بر } X_{1t} ,$$

$$\hat{\beta}_{y_2} = \text{تخمین شیب رگرسیون } Y_t \text{ بر } X_{2t} ,$$

$$\hat{\beta}_{12} = \text{تخمین شیب رگرسیون } X_{1t} \text{ بر } X_{2t} .$$

به همین ترتیب می‌توان $\hat{\beta}_{21}$ و $\hat{\beta}_{1y}$ ، $\hat{\beta}_{2y}$ را تعریف کرد. برای مثال

$$\hat{\beta}_{12} = \frac{\sum x_{1t} x_{2t}}{\sum x_{2t}^2} , \quad \hat{\beta}_{y_1} = \frac{\sum x_{1t} y_t}{\sum x_{1t}^2} , \quad \hat{\beta}_{1y} = \frac{\sum x_{1t} y_t}{\sum y_t^2} .$$

به شش پارامتر فوق «پارامترهای مدل جزئی» می‌گوییم. مدل جزئی، هر زیر مدلی است که بتوان از مدل رگرسیون مرکب ۴-۱ به دست آورد.

ابتدا این سؤال را مطرح می‌کنیم که آیا می‌توان رابطه بین تخمین پارامترهای مدل و پارامترهای مدل جزئی را به دست آورد؟ فرمول $\hat{\beta}_1$ را از معادله‌های ۴-۱۰ می‌نویسیم،

$$\hat{\beta}_1 = \frac{\sum x_{1t}^2 \sum x_{1t} y_t - \sum x_{1t} x_{2t} \sum x_{2t} y_t}{\sum x_{1t}^2 \sum x_{2t}^2 - (\sum x_{1t} x_{2t})^2} .$$

صورت و مخرج کسر فوق را بر $\sum x_{1t}^2 \sum x_{2t}^2$ تقسیم می‌کنیم،

$$\hat{\beta}_1 = \frac{\hat{\beta}_{y_1} - \hat{\beta}_{21} \hat{\beta}_{y_2}}{1 - \hat{\beta}_{12} \hat{\beta}_{21}} . \quad (4-41)$$

با استفاده از فرمول $\hat{\beta}_y$ و به صورتی کاملاً مشابه خواهیم داشت

$$\hat{\beta}_y = \frac{\hat{\beta}_{y_2} - \hat{\beta}_{1_2} \hat{\beta}_{y_1}}{1 - \hat{\beta}_{1_2} \hat{\beta}_{y_1}} \quad (4-42)$$

مفید است به طور خلاصه تحلیلی از معادله‌های 4-41 و 4-42 ارائه دهیم. معادله 4-41 را به صورت زیر می‌نویسیم،

$$\hat{\beta}_1 = \hat{\beta}_{y_1} - \frac{\hat{\beta}_{y_2} \hat{\beta}_{y_1}}{1 - \hat{\beta}_{1_2} \hat{\beta}_{y_1}}$$

یا:

$$\hat{\beta}_1 = \hat{\beta}_{y_1} - P$$

که P را می‌توان «جمله تعدیل» نامید؛ به عبارت دیگر، رابطه فوق نشان می‌دهد که اگر $\hat{\beta}_{y_1}$ را تعدیل کنیم، $\hat{\beta}_1$ به دست می‌آید. $\hat{\beta}_1$ ضریب X_{1i} در مدل رگرسیون دو متغیره 4-1 و $\hat{\beta}_{y_1}$ ضریب X_{1i} در مدل جزئی زیر است

$$Y_i = \alpha + \beta_1 X_{1i} + V_i$$

اگر X_{1i} و X_{2i} در نمونه با یکدیگر همبستگی نداشته باشند، یعنی $\sum x_{1i} \sum x_{2i} = 0$ ، آنگاه $\hat{\beta}_{y_1} = \hat{\beta}_{1_2} = 0$ زیرا

$$\hat{\beta}_{y_1} = \frac{\sum x_{2i} x_{1i}}{\sum x_{1i}^2} = 0$$

و در نتیجه جمله تعدیل در معادله‌های 4-41 و 4-42 برابر صفر می‌شود؛ بنابراین

$$\hat{\beta}_1 = \hat{\beta}_{y_1}, \quad \hat{\beta}_2 = \hat{\beta}_{y_2}$$

در چنین حالتی می‌گوییم که متغیرهای توضیحی، «متعامد» هستند و بدین ترتیب تخمین پارامترها در رگرسیون 4-1 دقیقاً برابر تخمین پارامترها در مدل‌های رگرسیون

ساده جزئی خواهد بود. این مسأله را در جلد دوم این کتاب با عنوان «مسأله همخطی» مطالعه خواهیم کرد.

سؤال دوم این است که آیا می توان $\hat{\beta}_1$ و $\hat{\beta}_2$ را برحسب ضرایب همبستگی ساده نوشت؟ پاسخ مثبت است. ابتدا نشان می دهیم که چگونه می توان ضریب همبستگی ساده بین X_1 و Y_1 را برحسب β نوشت.

مدل رگرسیون ساده زیر را مشاهده کنید،

$$Y_i = a + \beta X_i + U_i .$$

در معادله ۱-۳۲ دیدیم که ضریب همبستگی ساده بین دو متغیر X_1 و Y_1 عبارت است از

$$\begin{aligned} r_{yx} &= \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \sum y_i^2}} , \\ &= \frac{\sum x_i y_i}{\sum x_i^2} \cdot \frac{\sqrt{\sum x_i^2}}{\sqrt{\sum y_i^2}} . \end{aligned}$$

واریانس X_1 در نمونه برابر است با

$$\text{Var}(X_1) = S^2 = \frac{\sum (X_i - \bar{X})^2}{n} = \frac{\sum x_i^2}{n} ,$$

و انحراف معیار X_1 در نمونه عبارت است از

$$S = \sqrt{\frac{\sum x_i^2}{n}} . \quad (۴-۴۳)$$

بنابراین با استفاده از معادله ۴-۴۳ خواهیم داشت

$$r_{yx} = \hat{\beta}_{yx} \frac{S_x}{S_y} , \quad (۴-۴۴)$$

یا

$$\hat{\beta}_{yx} = r_{yx} \frac{S_y}{S_x} . \quad (۴-۴۵)$$

با استفاده از معادله ۴-۴۵ می توان معادله ۴-۴۱ را به صورت زیر نوشت،

$$\hat{\beta}_1 = \frac{r_{y1} \frac{S_y}{S_1} - r_{11} \frac{S_r}{S_1} r_{y2} \frac{S_y}{S_r}}{1 - r_{12} \frac{S_1}{S_r} r_{21} \frac{S_r}{S_1}} = \frac{r_{y1} - r_{11} r_{y2}}{1 - r_{12} r_{21}} \cdot \frac{S_y}{S_1}$$

یا:

$$\hat{\beta}_1 = \frac{r_{y1} - r_{11} r_{y2}}{1 - r_{12}^2} \cdot \frac{S_y}{S_1} \quad (4-46)$$

به ترتیبی کاملاً مشابه خواهیم داشت

$$\hat{\beta}_2 = \frac{r_{y2} - r_{12} r_{y1}}{1 - r_{12}^2} \cdot \frac{S_y}{S_r} \quad (4-47)$$

به تعریف دو اصطلاح می پردازیم: به تخمین پارامترهایی، مانند $\hat{\beta}_{y1}$ یا $\hat{\beta}_{y2}$ «ضرایب درجه اول»^۱ می گویند. به تخمین پارامترهایی مانند $\hat{\beta}_{y1/2}$ «ضرایب درجه اول»^۲ گفته می شود. منظور از $\hat{\beta}_{y1/2}$ ، تخمین ضریب اولین متغیر توضیحی در مدلی است که متغیر درون زای آن Y_1 است. عدد ۲ بعد از ممیز نشان می دهد که این مدل، یک متغیر توضیحی دیگر نیز دارد؛ بنابراین، معادله های ۴-۴۱ و ۴-۴۲، ضرایب درجه اول را برحسب ضرایب درجه صفر در یک مدل رگرسیون چند متغیره مشخص می سازد. سؤالی که مطرح می شود، این است که آیا در ضرایب همبستگی نیز چنین مراتبی وجود دارد؟

۱. ضرایب همبستگی درجه صفر و یک

ابتدا می گویم ضرایب همبستگی ساده (در مقدمه قسمت ۴-۵) را می توان «ضرایب همبستگی درجه صفر»^۲ نامید. مدل رگرسیون ۴-۱ را یک بار دیگر می نویسیم،

$$Y_i = \alpha + \beta_1 X_{i1} + \beta_2 X_{i2} + U_i$$

1. Zero-order Coefficients

2. First-order Coefficients

3. Zero-order Correlation Coefficients

می‌دانیم X_{1t} بر متغیر Y_t تأثیر می‌گذارد. از مدل فوق ملاحظه می‌شود که متغیر دیگری نیز به نام X_{2t} وجود دارد که Y_t از آن متأثر است. ضریب همبستگی ساده بین Y_t و X_{1t} عبارت است از

$$r_{y1} = \frac{\sum y_t x_{1t}}{\sqrt{\sum y_t^2 \sum x_{1t}^2}}$$

برای محاسبه r_{y1} ، تغییرات ملاحظه شده در Y_t را با تغییراتی که در X_{1t} وجود داشته است مرتبط کرده‌ایم. اما باید توجه داشت که تمام تغییرات Y_t به علت تغییر در X_{1t} نبوده، بلکه قسمتی از آن بر اثر تغییرات X_{2t} حاصل شده است؛ بنابراین r_{y1} واقعاً ضریب همبستگی تغییرات X_{1t} و تغییرات Y_t حاصل از آن را اندازه‌گیری نمی‌کند، زیرا Y_t از X_{2t} نیز متأثر است.

اگر مدل رگرسیون Y_t روی X_{1t} را بسازیم، خواهیم داشت

$$Y_t = \alpha + \beta_{y2} X_{2t} + V_t$$

مقادیر پسماند این مدل عبارت است از

$$e_t = Y_t - \hat{\beta}_{y2} X_{2t} \quad (4-48)$$

مقادیر پسماند را می‌توان این گونه تفسیر کرد که اگر تأثیر خطی تغییرات X_{2t} بر Y_t را کنار بگذاریم، آنچه باقی می‌ماند، تغییرات پسماند است، یعنی آن قسمتی از تغییرات Y_t که به تغییرات X_{2t} مربوط نیست؛ بنابراین باید ضریب همبستگی این قسمت از تغییرات Y_t را با X_{1t} حساب کنیم؛ به عبارت دیگر، باید ضریب همبستگی X_{1t} و $(Y_t - \hat{\beta}_{y2} X_{2t})$ را به دست آوریم. این ضریب همبستگی را با r_{y1}^* نشان می‌دهیم.

$$r_{y1}^* = \frac{\sum (y_t - \hat{\beta}_{y2} x_{2t}) x_{1t}}{\sqrt{\sum (y_t - \hat{\beta}_{y2} x_{2t})^2 \sum x_{1t}^2}} \quad (4-49)$$

نکته مهمی که باید در نظر داشت این است که X_{1t} و X_{2t} نیز در مدل (4-1) از یکدیگر مستقل نیستند. همین امر باعث می‌شود که نتوانیم ضریب همبستگی X_{1t} و Y_t را دقیقاً حساب کنیم؛ زیرا وقتی تغییرات X_{1t} را با Y_t یا با $(Y_t - \hat{\beta}_{y2} X_{2t})$ ، مرتبط

می‌کنیم، این «تغییرات خالص» X_{1t} نیست که با Y_t مرتبط است، بلکه قسمتی از تغییرات X_{1t} در واقع ناشی از تغییرات X_{2t} است. اکنون باید به همین ترتیب بتوان تغییرات X_{1t} را که به علت X_{2t} بوده است حذف کرد. کافی است یک مدل رگرسیون بسازیم که X_{1t} را تابعی از X_{2t} فرض کند. داریم

$$X_{1t} = \alpha + \beta_{12} X_{2t} + \varepsilon_t .$$

مقادیر پسماند عبارت است از

$$e_t = X_{1t} - \hat{\beta}_{12} X_{2t} . \quad (4.50)$$

یعنی اگر تأثیر خطی تغییرات X_{2t} بر X_{1t} را کنار بگذاریم، آنچه باقی می‌ماند، تغییرات پسماند است، یعنی آن قسمتی از تغییرات X_{1t} که به تغییرات X_{2t} مربوط نیست؛ بنابراین باید در معادله ۴-۴۹ به جای X_{1t} ، مقدار پسماند ۴-۵۰ را قرار داد. خواهیم داشت

$$r_{y1/2} = \frac{\sum (y_t - \hat{\beta}_{y2} x_{2t}) (x_{1t} - \hat{\beta}_{12} x_{2t})}{\sqrt{\sum (y_t - \hat{\beta}_{y2} x_{2t})^2 \sum (x_{1t} - \hat{\beta}_{12} x_{2t})^2}} . \quad (4.51)$$

معادله ۴-۵۱ فرمول «ضریب همبستگی درجه اول»^۱ است. $r_{y1/2}$ نشان می‌دهد که می‌خواهیم ضریب همبستگی بین Y_t و اولین متغیر توضیحی (X_{1t}) را در مدلی حساب کنیم که یک متغیر توضیحی دیگر نیز به نام X_{2t} دارد؛ بنابراین باید توجه کرد که عدد ۲ بعد از ممیز بدین معنی نیست که متغیر توضیحی دیگر «ثابت» فرض شده است، بلکه می‌گوییم X_{2t} به منزله یک متغیر توضیحی وجود دارد، اما سعی می‌کنیم تأثیر تغییرات آن را در Y_t و نیز در X_{1t} به نحوی حذف کنیم که بتوانیم خالص تغییرات X_{1t} و Y_t را با یکدیگر مرتبط ساخته، ضریب همبستگی آنها را محاسبه کنیم. به $r_{y1/2}$ ، یعنی ضریب همبستگی درجه اول، معمولاً «ضریب همبستگی جزئی»^۲ می‌گویند.

باید توجه داشت که اگر بین X_{1t} و X_{2t} در نمونه، همبستگی وجود نداشته باشد،

1. First-order Correlation Coefficient
2. Partial Correlation Coefficient

در معادله ۴-۴۹ جمله $\sum x_{1t} x_{2t}$ برابر صفر شده و ضریب همبستگی جزئی عبارت خواهد بود از

$$r_{y_1/r}^* = \frac{\sum (y_t x_{1t})}{\sqrt{\sum (y_t - \hat{\beta}_{y_1} x_{1t})^2} \sqrt{\sum x_{1t}^2}}$$

۲. ضریب همبستگی جزئی بر حسب ضرایب همبستگی ساده

در معادله ۴-۵۱ توانستیم ضریب همبستگی جزئی بین X_{1t} و Y_t را در مدلی به دست آوریم که شامل X_{2t} است. در این قسمت می‌خواهیم این ضریب همبستگی جزئی را بر حسب ضرایب همبستگی ساده بنویسیم. ابتدا صورت کسر ۴-۵۱ را به شرح زیر می‌نویسیم:

$$r_{y_1/r} = \frac{\sum y_t x_{1t} - \hat{\beta}_{12} \sum y_t x_{2t} - \hat{\beta}_{21} \sum x_{1t} x_{2t} + \hat{\beta}_{12} \hat{\beta}_{21} \sum x_{2t}^2}{\sqrt{\sum (y_t - \hat{\beta}_{12} x_{2t} - \hat{\beta}_{21} x_{1t})^2} \sqrt{\sum x_{1t}^2}}$$

جمله اول $(\sum y_t x_{1t})$ را ملاحظه می‌کنیم. از معادله ۱-۳۲ می‌دانیم

$$r_{y_1} = \frac{\sum y_t x_{1t}}{\sqrt{\sum y_t^2} \sqrt{\sum x_{1t}^2}}$$

در نتیجه

$$\sum y_t x_{1t} = n r_{y_1} \sqrt{\frac{\sum y_t^2}{n}} \sqrt{\frac{\sum x_{1t}^2}{n}}$$

همچنین از معادله ۴-۴۳ می‌دانیم

$$s_y = \sqrt{\frac{\sum y_t^2}{n}}, \quad s_1 = \sqrt{\frac{\sum x_{1t}^2}{n}}$$

بنابراین جمله اول صورت کسر را می‌توان به صورت زیر نوشت،

$$\sum y_t x_{1t} = n r_{y_1} s_y s_1 \quad (۴-۵۲)$$

حال به جمله دوم توجه می‌کنیم. از معادله ۴-۴۵ می‌دانیم

$$\hat{\beta}_{12} = r_{12} \frac{s_1}{s_2}$$

دقیقاً مشابه معادله ۴-۵۲ خواهیم داشت

$$\sum y_t x_{rt} = n r_{yt} S_y S_r ,$$

در نتیجه جمله دوم رابه صورت زیر می نویسیم،

$$\hat{\beta}_{1r} = \sum y_t x_{rt} = r_{1r} \frac{S_1}{S_r} n r_{yt} S_y S_r .$$

به همین ترتیب خواهیم داشت

$$\hat{\beta}_{yr} \sum x_{1t} x_{rt} = r_{yr} \frac{S_y}{S_r} n r_{1r} S_1 S_r ,$$

$$\hat{\beta}_{yr} \hat{\beta}_{1r} \sum x_{rt}^2 = r_{yr} r_{1r} \frac{S_y S_1}{S_r^2} n S_r^2 .$$

عبارتهای فوق را در صورت کسر جایگزین می کنیم،

$$\text{صورت کسر} = n r_{1r} S_y S_1 - r_{1r} \frac{S_1}{S_r} n r_{yt} S_y S_r - r_{yr} \frac{S_y}{S_r} n r_{1r} S_1 S_r$$

$$+ r_{yr} r_{1r} \frac{S_y S_1}{S_r^2} n S_r^2 .$$

و بعد از ساده کردن داریم

$$\text{صورت کسر} = n S_y S_1 (r_{1r} - r_{yr} r_{1r}) .$$

مخرج کسر ۴-۵۱ را نیز می توان به راحتی ساده کرد. یک بار دیگر مخرج کسر را

می نویسیم،

$$\text{مخرج کسر} = \sqrt{\sum (y_t - \hat{\beta}_{yr} x_{rt})^2 \sum (x_{1t} - \hat{\beta}_{1r} x_{rt})^2} .$$

جمله $\sum (y_t - \hat{\beta}_{yr} x_{rt})^2$ ، در واقع مجموع مربعات پسماند برای مدل رگرسیون زیر است،

$$Y_t = \alpha + \beta_{yr} X_{rt} + V_t .$$

با استفاده از معادله ۴-۴۳، برای معادله فوق داریم

$$r_{y1}^2 = 1 - \frac{RSS}{\sum y_t^2}$$

یا

$$RSS = \sum y_t^2 (1 - r_{y1}^2)$$

با توجه به معادله ۴۳-۴ می‌دانیم

$$\sum y_t^2 = \frac{n \sum y_t^2}{n} = n S_y^2$$

بنابراین

$$RSS = n S_y^2 (1 - r_{y1}^2)$$

جمله $\sum (y_t - \hat{\beta}_{y1} x_{1t})^2$ (مجموع مربعات پسماند برای مدل جزئی فوق) را می‌توان به صورت زیر نوشت،

$$\sum (y_t - \hat{\beta}_{y1} x_{1t})^2 = n S_y^2 (1 - r_{y1}^2)$$

به همین ترتیب می‌توان گفت که جمله $\sum (x_{1t} - \hat{\beta}_{11} x_{1t})^2$ یعنی مجموع مربعات پسماند برای مدل

$$x_{1t} = \alpha + \beta_{11} x_{1t} + \varepsilon_t$$

عبارت خواهد بود از

$$\sum (x_{1t} - \hat{\beta}_{11} x_{1t})^2 = n S_1^2 (1 - r_{11}^2)$$

اکنون مخرج کسر ۴۵-۱ را به صورت زیر می‌نویسیم،

$$\text{مخرج کسر} = n S_y S_1 \sqrt{(1 - r_{y1}^2)(1 - r_{11}^2)}$$

با توجه به نتایج فوق برای صورت و مخرج، می‌توان کسر ۴۵-۱ را به شرح زیر

نوشت،

$$r_{y1/1} = \frac{r_{y1} - r_{y1} r_{11}}{\sqrt{(1 - r_{y1}^2)(1 - r_{11}^2)}} \quad (4.53)$$

معادله ۴-۵۳ مقدار ضریب همبستگی جزئی را بین متغیر درون‌زای Y_1 و متغیر برون‌زای X_{12} مشخص می‌کند. به همین ترتیب می‌توان ضریب همبستگی جزئی را بین متغیر درون‌زای Y_1 و متغیر برون‌زای X_{21} تعریف کرد. به همین ترتیب و با توجه به آنچه که تاکنون داشتیم، می‌توان گفت که

$$r_{y_1/1} = \frac{r_{y_2} - r_{y_1} r_{12}}{\sqrt{(1 - r_{y_1}^2)(1 - r_{12}^2)}} \quad (4-54)$$

یادآوری این نکته مفید است که $r_{y_1/1}$ ، یعنی ضریب همبستگی جزئی بین Y_1 و X_{12} ، در حقیقت یک ضریب همبستگی ساده بین دو سری پسماند به شرح زیر است.

$$1. (y_1 - \hat{\beta}_{y_1} x_{12}) \text{ که مقدار } e_1 \text{ برای مدل زیر است،}$$

$$Y_1 = \alpha + \beta_{y_1} X_{12} + V_1 .$$

$$2. (x_{21} - \hat{\beta}_{21} x_{12}) \text{ که مقدار } e_1 \text{ برای مدل زیر خواهد بود،}$$

$$X_{21} = \alpha + \beta_{21} x_{12} + e_1 .$$

نتیجه کلی بحث را یک بار دیگر ذکر می‌کنیم: ضریب همبستگی جزئی بین هر دو متغیر در حقیقت همبستگی بین پسماندهای هر یک از این متغیرهاست وقتی تأثیرات خطی تمام متغیرهای توضیحی دیگر بر روی آنها را حذف کرده باشیم.

در پایان این قسمت به معرفی یک اصطلاح می‌پردازیم. به مجذور ضریب همبستگی جزئی اصطلاحاً «ضریب تعیین جزئی»^۱ می‌گویند.

۳. ضرایب تعیین و پسماندها

در این قسمت می‌خواهیم اهمیت پسماندها و ارتباط آن را با ضرایب تعیین بیشتر تبیین کنیم. اولین نکته این است که می‌توان به کمک پسماندهای به دست آمده از مدل‌های جزئی رگرسیون، به تخمین پارامترهای یک رگرسیون چند متغیره رسید.

مدل رگرسیون ۴-۱ را یک بار دیگر می نویسیم،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i .$$

مدل رگرسیون جزئی که Y_i را بر حسب X_{2i} بیان می کند، در نظر می گیریم.

$$Y_i = \alpha + \beta_{y2} X_{2i} + V_i .$$

پسماند این مدل جزئی عبارت است از

$$(y_i - \hat{\beta}_{y2} x_{2i}) . \quad (4.55)$$

به همین ترتیب یک مدل رگرسیون جزئی دیگر می سازیم که در آن X_{1i} بر حسب X_{2i} بیان شده باشد.

$$X_{1i} = \alpha + \beta_{12} X_{2i} + \varepsilon_i . \quad (4.56)$$

پسماند این مدل به شرح زیر است،

$$(x_{1i} - \hat{\beta}_{12} x_{2i}) . \quad (4.57)$$

یک مدل جزئی می سازیم که متغیرهای ۴-۵۶ و ۴-۵۷ به ترتیب متغیرهای درون‌زا و برون‌زای آن باشد. تخمین شیب این مدل خطی ساده را می نویسیم،

$$\begin{aligned} \text{شیب مدل جزئی} &= \frac{\sum (y_i - \hat{\beta}_{y2} x_{2i}) (x_{1i} - \hat{\beta}_{12} x_{2i})}{\sum (x_{1i} - \hat{\beta}_{12} x_{2i})^2} , \\ &= \frac{\sum y_i (x_{1i} - \hat{\beta}_{12} x_{2i}) - \hat{\beta}_{y2} \sum x_{2i} (x_{1i} - \hat{\beta}_{12} x_{2i})}{\sum (x_{1i} - \hat{\beta}_{12} x_{2i})^2} . \end{aligned}$$

در مدل رگرسیون ۴-۵۶ می دانیم X_{2i} متغیر توضیحی و $(x_{1i} - \hat{\beta}_{12} x_{2i})$ همان پسماند یا ε_i است. از طرف دیگر از معادله ۱-۲۴ می دانیم که در تخمینهای حداقل مربعات معمولی، متغیر توضیحی از جمله پسماند مستقل است،

$$\sum x_{2i} (x_{1i} - \hat{\beta}_{12} x_{2i}) = 0 ,$$

بدین ترتیب شیب مدل جزئی برابر است با

$$\text{شیب مدل جزئی} = \frac{\sum y_i x_{1i} - \hat{\beta}_{12} \sum y_i x_{2i}}{\sum (x_{1i} - \hat{\beta}_{12} x_{2i})^2}$$

با جایگزینی $\hat{\beta}_{12} = \frac{\sum x_{1i} x_{2i}}{\sum x_{2i}^2}$ در رابطه مذکور، خواهیم داشت

$$\text{شیب مدل جزئی} = \frac{\sum x_{1i}^2 \sum x_{1i} y_i - \sum x_{1i} x_{2i} \sum x_{2i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2} = \beta_1$$

با توجه به معادله ۴-۱۰ معلوم می‌شود که طرف راست معادله فوق دقیقاً همان $\hat{\beta}_1$ است. بنابراین توانستیم به کمک پسماندهای حاصل از مدل‌های جزئی رگرسیون، به تخمین پارامتر β_1 در مدل رگرسیون چند متغیره برسیم.

به روشی کاملاً مشابه، می‌توان به کمک پسماندهای حاصل از مدل رگرسیون

جزئی

$$Y_i = \alpha + \beta_{y1} X_{1i} + V_i$$

یعنی $(y_i - \hat{\beta}_{y1} x_{1i})$ و نیز پسماندهای حاصل از مدل رگرسیون جزئی

$$X_{2i} = \alpha + \beta_{x21} X_{1i} + V_i$$

یعنی $(x_{2i} - \hat{\beta}_{x21} x_{1i})$ ، یک مدل رگرسیون ساخت و نشان داد که تخمین شیب این مدل دقیقاً برابر $\hat{\beta}_2$ است.

۴. ضریب تعیین برحسب ضرایب همبستگی جزئی

فرض کنید ضرایب همبستگی جزئی در یک مدل رگرسیون چند متغیره را می‌دانیم. آیا می‌توان ضریب تعیین این مدل را برحسب ضرایب همبستگی جزئی نوشت؟ معادله ۴-۱۹ را یک بار دیگر می‌نویسیم،

$$R^2 = R_{y/12}^2 = \frac{\hat{\beta}_1 \sum x_{1i} y_i + \hat{\beta}_2 \sum x_{2i} y_i}{\sum y_i^2}$$

قیلاً در معادله‌های ۴-۴۶ و ۴-۴۷ توانستیم تخمین پارامترهای مدل رگرسیون چند متغیره، یعنی $\hat{\beta}_1$ ، $\hat{\beta}_2$ را بر حسب ضرایب همبستگی ساده بنویسیم. این مقادیر را در معادله فوق قرار می‌دهیم. در صورت کسر، دو جمله وجود دارد. ابتدا هر یک از این دو جمله را جداگانه محاسبه می‌کنیم،

$$\begin{aligned} \hat{\beta}_1 \sum x_{1t} y_t &= \frac{(r_{y1} - r_{r1} r_{y2}) S_y}{(1 - r_{12}^2) S_1} \cdot R \frac{\sum x_{1t} y_t}{1} \\ &= \frac{(r_{y1} - r_{r1} r_{y2})}{(1 - r_{12}^2)} \cdot \frac{\sqrt{\sum y_t^2}}{\sqrt{\sum x_{1t}^2}} \cdot \frac{\sqrt{\sum x_{1t} y_t}}{1} \cdot \frac{\sqrt{\sum x_{1t} y_t}}{\sqrt{\sum y_t^2}} \cdot \frac{\sqrt{\sum y_t^2}}{1} \\ &= \frac{(r_{y1} - r_{r1} r_{y2})}{(1 - r_{12}^2)} \cdot \sqrt{\hat{\beta}_{y1}} \cdot \hat{\beta}_{1y} \cdot \sum y_t^2 \end{aligned}$$

با توجه به معادله ۴-۴۵ می‌توان چنین نوشت،

$$\hat{\beta}_{y1} = r_{y1} \frac{S_y}{S_1} \quad , \quad \hat{\beta}_{1y} = r_{1y} \frac{S_1}{S_y}$$

در نتیجه خواهیم داشت

$$\hat{\beta}_1 \sum x_{1t} y_t = \frac{(r_{y1} - r_{r1} r_{y2})}{(1 - r_{12}^2)} \cdot \sqrt{r_{y1} \frac{S_y}{S_1} r_{1y} \frac{S_1}{S_y}} \cdot \sum y_t^2$$

با توجه به اینکه $r_{y1} = r_{1y}$ ، نتیجه می‌شود که

$$\hat{\beta}_1 \sum x_{1t} y_t = \frac{(r_{y1}^2 - r_{r1} r_{y2} r_{y1})}{(1 - r_{12}^2)} \cdot \sum y_t^2 \quad (4-58)$$

برای جمله دوم در صورت کسر ۴-۱۹ داریم

$$\hat{\beta}_2 \sum x_{2t} y_t = \frac{(r_{y2}^2 - r_{r1} r_{y2} r_{y1})}{(1 - r_{12}^2)} \cdot \sum y_t^2 \quad (4-59)$$

معادله‌های ۴-۵۸ و ۴-۵۹ را در معادله ۴-۱۹ جایگزین می‌کنیم و بعد از ساده کردن داریم

$$R_{y12}^2 = \frac{r_{y1}^2 + r_{y2}^2 - 2 r_{y1} r_{y2} r_{r1}}{(1 - r_{12}^2)} \quad (4-60)$$

و به این ترتیب با داشتن ضرایب همبستگی ساده می‌توان به ضریب تعیین رسید.

۵. تجزیه تغییرات توضیح داده شده به کمک ضرایب تعیین

در یک مدل رگرسیون می‌توان تغییرات توضیح داده شده (ESS) را به کمک ضرایب تعیین ساده و جزئی تجزیه کرد و مکانیسم شکل‌گیری آنها را بررسی نمود. بدیهی است این مسأله در تحلیل آنالیز واریانس برای رگرسیون چند متغیره که در فصل ششم بررسی خواهد شد، اهمیت فراوان دارد.

فرض کنید یک متغیر درون‌زا (Y_t) و دو متغیر توضیحی (X_{1t} , X_{2t}) داریم. ابتدا مدل رگرسیون Y_t بر روی X_{1t} را می‌سازیم،

$$Y_t = \alpha + \beta_{y1} X_{1t} + V_t .$$

با استفاده از تعریف ضریب تعیین، معادله ۱-۲۹ داریم

$$r_{y1}^2 = \frac{ESS}{\sum y_t^2} ,$$

یا

$$ESS = r_{y1}^2 \sum y_t^2 . \quad (۴-۶۱)$$

با توجه به معادله ۱-۴۳ می‌توان مجموع مربعات پسماند را به صورت زیر نوشت،

$$RSS = (1 - r_{y1}^2) \sum y_t^2 . \quad (۴-۶۲)$$

بدیهی است که این پسماندها، یعنی $(y_t - \hat{\beta}_{y1} x_{1t})$ باید توضیح داده شوند.

X_{2t} را وارد مدل می‌کنیم. انتظار این است که بتوانیم پسماندهای Y_t را با X_{2t} توضیح دهیم. اما X_{2t} از X_{1t} مستقل نیست؛ بنابراین قسمتی از تغییرات X_{2t} با X_{1t} و با استفاده از مدل زیر توضیح داده می‌شود.

$$X_{2t} = \alpha + \beta_{21} X_{1t} + \varepsilon_t ,$$

در نتیجه پسماندهای X_{2t} عبارت است از $(x_{2t} - \hat{\beta}_{21} x_{1t})$. باید بتوانیم به کمک پسماندهای X_{2t} ، پسماندهای Y_t را توضیح دهیم. یک رگرسیون می‌سازیم که پسماندهای Y_t را بر حسب پسماندهای X_{2t} بیان کند.

ضریب تعیین این رگرسیون در واقع یک ضریب تعیین جزئی است که می‌توان آن را به صورت $r_{y_2/1}^2$ نوشت، که دلالت بر مدلی می‌کند که می‌خواهد تأثیر X_{2t} را بر Y_t اندازه‌گیری کند، وقتی X_{1t} تأثیر خود را گذاشته است. با توجه به معادله ۱-۲۹ می‌توان تغییرات توضیح داده شده را در این مدل حساب کرد. می‌دانیم برای به دست آوردن تغییرات توضیح داده شده (ESS) باید ضریب تعیین را در کل تغییرات (TSS) ضرب کنیم. اما کل تغییراتی که در این مدل باید توضیح داده شود چیزی نیست جز همان مجموع مربعات پسماند Y_t ، یعنی معادله ۴-۶۲، که از مرحله قبل باقی مانده است. بنابراین برای به دست آوردن تغییرات توضیح داده شده کافی است که $r_{y_2/1}^2$ را در معادله ۴-۶۱ ضرب کنیم،

$$ESS = r_{y_2/1}^2 \left[(1 - r_{y_1}^2) \sum y_t^2 \right]. \quad (۴-۶۳)$$

روی پسماندهای X_{2t} در مدل رگرسیون پسماندهای Y_t

به سهولت می‌توان مجموع مربعات پسماند در مدل فوق را نیز به دست آورد. با توجه به معادله ۱-۴۳ داریم،

$$RSS = (1 - r_{y_2/1}^2) \left[(1 - r_{y_1}^2) \sum y_t^2 \right]. \quad (۴-۶۴)$$

روی پسماندهای X_{2t} در مدل رگرسیون پسماندهای Y_t

بدین ترتیب نتیجه می‌گیریم که تغییرات توضیح داده شده Y_t ، در دو مرحله ساخته شده است. مرحله اول که در واقع معادله ۴-۶۱ است، و آنچه که در مرحله دوم به دست آمده و در معادله ۴-۶۳ منعکس است. این دو را با هم جمع می‌کنیم تا کل تغییرات توضیح داده شده Y_t در مدل رگرسیون مفروض به دست آید،

$$ESS = \sum y_t^2 [r_{y_1}^2 + r_{y_2/1}^2 (1 - r_{y_1}^2)]. \quad (۴-۶۵)$$

به عبارت دیگر کل تغییرات توضیح داده شده Y_i شامل دو قسمت است: تغییراتی که از رگرسیون ساده Y_i روی X_{1i} به دست می‌آید و تغییراتی که به رگرسیون ساده پسماندهای Y_i روی پسماندهای X_{2i} مربوط است؛ بنابراین با ورود X_{2i} به مدل، تغییرات توضیح داده شده به اندازه $[(1 - r_{y_1}^2) \sum y_i^2]$ $r_{y_2/1}^2$ اضافه می‌شود، که در آن $\sum y_i^2 (1 - r_{y_1}^2)$ مجموع مربعات پسماندهایی است که حاصل رگرسیون Y_i روی X_{1i} است.

یک بار دیگر به تعریف $r_{y_2/1}^2$ اشاره می‌کنیم. در یک رگرسیون ساده Y_i روی X_{1i} ، تغییراتی که توضیح داده نشده است، (پسماندها) را در نظر می‌گیریم. $r_{y_2/1}^2$ در واقع درصدی از این تغییرات توضیح داده نشده را حساب می‌کند که با ورود X_{2i} به مدل می‌تواند توضیح داده شود.

معادله ۴-۶۵ را دوباره ملاحظه کنید. وقتی متغیر X_{2i} به X_{1i} اضافه شود، ESS کل تغییرات توضیح داده نشده در Y_i است؛ به عبارت دیگر، کل تغییرات توضیح داده شده چیزی نیست جز تغییرات توضیح داده شده در مدل زیر

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i .$$

اما از معادله ۴-۱۶ می‌دانیم که

$$R_{y/12}^2 = \frac{ESS}{\sum y_i^2} ,$$

یا

$$ESS = R_{y/12}^2 \sum y_i^2 . \quad (4-66)$$

با جایگزینی معادله ۴-۶۶ در ۴-۶۵ خواهیم داشت

$$R_{y/12}^2 = r_{y_1}^2 + r_{y_2/1}^2 (1 - r_{y_1}^2) . \quad (4-67)$$

معادله فوق رابطه بین ضریب تعیین کل، ضریب تعیین جزئی و ضریب تعیین ساده در یک مدل رگرسیون، شامل دو متغیر توضیحی را نشان می‌دهد.

مثال ۴-۵ برای تخمین تغییرات دستمزدها مدل رگرسیون زیر مفروض است،

$$W_i = \alpha + \beta_1 \frac{1}{U_i} + \beta_2 P_{i-1} + \varepsilon_i ,$$

که در آن W_t درصد تغییر دستمزد و U_t نرخ بیکاری به صورت درصدی از نیروی کار و P_t درصد تغییر در شاخص قیمت کالاهای مصرفی است. آمار این متغیرها به صورت فصلی و برای ۲۶ دوره جمع آوری و محاسبات زیر انجام شده است،

$$r_{w1} = 0/6508 \quad , \quad r_{w2} = 0/3567 \quad , \quad r_{12} = 0/0726 \quad .$$

الف) ضرایب تعیین ساده r_{w1}^T و r_{w2}^T را حساب کنید.

ب) ضرایب همبستگی جزئی و ضرایب تعیین جزئی بین W_t و $\frac{1}{U_t}$ و نیز بین W_t و P_{t-1} را به دست آورید.

ج) ضریب تعیین مدل رگرسیون مفروض، یعنی $R_{w/12}^2$ ، را محاسبه کنید.

د) نتایج به دست آمده را توضیح دهید.

الف) باید ضرایب همبستگی ساده را مجدور کنیم،

$$r_{w1}^T = (r_{w1})^T = (0/6508)^T = 0/4235 \quad ,$$

$$r_{w2}^T = (r_{w2})^T = (0/3567)^T = 0/1272 \quad .$$

ب) برای محاسبه ضریب همبستگی جزئی بین W_t و $\frac{1}{U_t}$ از معادله ۴-۵۳ استفاده

می کنیم،

$$\begin{aligned} r_{w/12} &= \frac{r_{w1} - r_{w2} r_{12}}{\sqrt{(1 - r_{w2}^T)(1 - r_{12}^T)}} \quad , \\ &= \frac{0/6508 - 0/3567(0/0726)}{\sqrt{(1 - 0/1272)(1 - 0/0726^2)}} = 0/6707 \quad . \end{aligned}$$

برای به دست آوردن ضریب همبستگی جزئی بین W_t و P_{t-1} کافی است از معادله

۴-۵۴ استفاده شود،

$$\begin{aligned} r_{w2/1} &= \frac{r_{w2} - r_{w1} r_{12}}{\sqrt{(1 - r_{w1}^T)(1 - r_{12}^T)}} \quad , \\ &= \frac{0/3567 - 0/6508(0/0726)}{\sqrt{(1 - 0/4235)(1 - 0/0726^2)}} = 0/4087 \quad . \end{aligned}$$

برای به دست آوردن ضریب تعیین جزئی $r_{w1/2}^2$ ، باید $r_{w1/2}$ را مجذور کنیم،

$$r_{w1/2}^2 = (r_{w1/2})^2 = (0/6707)^2 = 0/4498 ,$$

و به همین ترتیب

$$r_{w2/1}^2 = (r_{w2/1})^2 = (0/4087)^2 = 0/1670 .$$

ملاحظه می شود که در این مثال چون همبستگی بین متغیرهای توضیحی بسیار ضعیف، یعنی $r_{12} = 0/0726$ ؛ تفاوت بسیار مختصری بین ضرایب همبستگی مرتبه صفر و یک وجود دارد.

ج) از معادله ۴-۶۰ می توان برای به دست آوردن $R_{w1/2}^2$ استفاده کرد؛

$$R_{w1/2}^2 = \frac{r_{w1}^2 + r_{w2}^2 - 2 r_{w1} r_{w2} r_{12}}{(1 - r_{12}^2)}$$

بنابراین خواهیم داشت

$$R_{w1/2}^2 = \frac{0/4235 + 0/1272 - 2(0/6008)(0/3567)(0/0726)}{1 - (0/0726)^2} ,$$

$$= \frac{0/517}{0/9948} = 0/5197 .$$

د) با ملاحظه $r_{w1}^2 = 0/4235$ می توان گفت که متغیر $\frac{1}{U_1}$ به تنهایی می تواند بیشتر از ۴۲ درصد تغییرات دستمزدها را توضیح دهد و با توجه به $r_{w2}^2 = 0/1272$ به این نتیجه می رسیم که حدود ۱۳ درصد از تغییر در دستمزدها می تواند به تنهایی به تغییر در قیمت‌ها نسبت داده شود. اگر هر دو متغیر توضیحی، یعنی بیکاری و قیمت، با هم عمل کنند حدود ۵۲ درصد تغییرات دستمزدها توضیح داده خواهد شد و این حقیقتی است که با ملاحظه $R_{w1/2}^2 = 0/5197$ به دست می آید. وقتی قیمت تنها متغیر توضیحی باشد، نمی تواند همه تغییرات دستمزدها را توضیح دهد و قسمتی از این تغییرات به صورت پسماند باقی می ماند. با توجه به ضریب تعیین جزئی $r_{w1/2}^2 = 0/4498$ ، می توان گفت

که بیکاری می‌تواند حدود ۴۵ درصد از این تغییرات توضیح داده نشده را توضیح دهد. همچنین وقتی بیکاری به عنوان تنها متغیر توضیحی عمل کند، قطعاً پسماندهایی در تغییرات دستمزد خواهد بود که متغیر قیمت، با توجه به ضریب تعیین جزئی $\tau_{w/p}^1 = 0/1670$ ، می‌تواند حدود ۱۷ درصد تغییرات آن را تبیین کند.

مسائل فصل چهارم

۴-۱ تخمین معادلات رگرسیون زیر مفروض است،

$$C_t = \hat{\alpha}_1 + 0.92 Y_t + e_{1t} ,$$

$$C_t = \hat{\alpha}_2 + 0.84 C_{t-1} + e_{2t} ,$$

$$C_{t-1} = \hat{\alpha}_3 + 0.78 Y_t + e_{3t} ,$$

$$Y_t = \hat{\alpha}_4 + 0.55 C_{t-1} + e_{4t} ,$$

که در آن $\hat{\alpha}_i$ تخمین جمله ثابت و e_{it} مقادیر پسماند است. با استفاده از اطلاعات مذکور، پارامترهای β_1 و β_2 در مدل رگرسیون زیر را تخمین بزنید،

$$C_t = \alpha + \beta_1 Y_t + \beta_2 C_{t-1} + U_t .$$

۴-۲ فرض کنید دو نمونه با مشخصات زیر موجود است.

نمونه اول	نمونه دوم
$n_1 = 20$	$n_2 = 25$
$\bar{X}_1 = 20$	$\bar{X}_2 = 23$
$\bar{Y}_1 = 25$	$\bar{Y}_2 = 28$
$\sum x_{1t}^2 = 80$	$\sum x_{2t}^2 = 100$
$\sum x_{1t} y_{1t} = 120$	$\sum x_{2t} y_{2t} = 150$
$\sum y_{1t}^2 = 200$	$\sum y_{2t}^2 = 250$

الف) برای هر نمونه یک معادله رگرسیون خطی تخمین بزنید.

ب) دو نمونه را با یکدیگر تلفیق کرده، یک معادله رگرسیون خطی تخمین بزنید.

ج) برای اینکه تخمین پارامترها در «رگرسیون تلفیقی»^۱ معتبر باشد، چه فرضهایی را باید پذیرفت.

۴-۳ مدل رگرسیون زیر مفروض است،

$$Y_i = \alpha + \beta X_i + U_i .$$

این مدل را تخمین زده، سپس متغیر توضیحی دیگری، مانند Z_i به آن اضافه می‌کنیم:

$$Y_i = \alpha' + \beta' X_i + \gamma' Z_i + U'$$

در چه شرایطی موارد زیر صحیح است.

الف) $\hat{\beta} = \hat{\beta}'$

ب) $\sum e_i' \geq \sum e_i$

ج) در سطح معنی‌دار ۵ درصد، $\hat{\beta}$ معنی‌دار است، اما $\hat{\beta}'$ معنی‌دار نیست.

د) در سطح معنی‌دار ۵ درصد، $\hat{\beta}'$ معنی‌دار است، اما $\hat{\beta}$ معنی‌دار نیست.

۴-۴ ابتدا مدل شماره (۱) را تخمین می‌زنیم،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i , \quad (1)$$

آنگاه مدل زیر را تخمین می‌زنیم

$$X_{1i} = \gamma + \lambda X_{2i} + V_i , \quad (2)$$

و پسماندهای رگرسیون مدل آخر، یعنی \hat{V}_i را حساب می‌کنیم. سپس مدل شماره (۳) را تخمین می‌زنیم،

$$Y_i = \alpha' + \beta_1' \hat{V}_i + \beta_2' X_{2i} + \varepsilon_i \quad (3)$$

نشان دهید که، $\hat{\beta}_1 = \tilde{\beta}_1$. این نتیجه را چگونه توضیح می‌دهید؟
 ۴-۵. در مدل رگرسیون

$$y_t = \beta x_{1t} + \gamma x_{2t} + U_t,$$

که $t = 1, 2, \dots, n$ ، تمام متغیرها بر حسب انحراف از میانگین نوشته شده است. دو روش مختلف زیر را برای تخمین β در نظر می‌گیریم:

(الف) مدل فوق را به همین صورت موجود تخمین می‌زنیم و $\hat{\beta}$ و $\hat{\gamma}$ را حساب می‌کنیم.

(ب) ابتدا رگرسیون y_t بر روی x_{2t} را تخمین زده، پسماندهای رگرسیون، یعنی \hat{e}_t را به دست می‌آوریم. آنگاه رگرسیون x_{1t} روی x_{2t} را تخمین می‌زنیم و پسماندهای این رگرسیون، یعنی \hat{e}_t^* را محاسبه می‌کنیم. سپس رگرسیون \hat{e}_t بر روی x_{1t}^* را تخمین می‌زنیم و شیب این تخمین را $\hat{\rho}$ می‌نامیم.

$$\text{اولاً، نشان دهید } \hat{\rho} = \hat{\beta}.$$

ثانیاً، نشان دهید پسماندهای رگرسیون در دو روش فوق با یکدیگر برابر است.

۴-۶. به جای اینکه پارامترهای β_1 و β_2 را از مدل رگرسیون زیر تخمین بزنیم،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t,$$

از این مدل استفاده می‌شود:

$$Y_t = \alpha' + \beta_1 X_{1t}^* + \beta_2 X_{2t} + V_t,$$

که در آن X_{1t}^* ، پسماند رگرسیون X_{1t} روی X_{2t} ، و V_t جمله اختلال است.

(الف) نشان دهید تخمین β_2 از مدل دوم دقیقاً برابر است با تخمین شیب مدلی که از رگرسیون Y_t روی X_{2t} به دست می‌آید.

(ب) مقدار اریب این تخمین را حساب کنید.

۴-۷ یک پژوهشگر اقتصادسنجی می‌خواهد مدل مصرف زیر را تخمین بزند،

$$C_t = \alpha + \beta Y_t + \gamma S_t + U_t ,$$

که در آن C_t مصرف، Y_t درآمد و S_t پس‌انداز در سطح کلان است. می‌دانیم

$$Y_t = C_t + S_t .$$

به نظر شما چه مسائلی در تخمین این مدل مصرف وجود دارد؟

۴-۸ نتایج تخمین دو مدل مختلف رگرسیون را ارائه می‌کنیم. چگونه می‌توان گفت

که در کدام موارد حتماً اشتباهاتی صورت گرفته است.

$$R_{y/12}^2 = 0/89 \quad , \quad R_{y/123}^2 = 0/86 \quad (\text{الف})$$

$$R_{y/12}^2 = 0/70 \quad , \quad r_{y/12}^2 = 0/13 \quad , \quad r_{y/12}^2 = 0/23 \quad (\text{ب})$$

۴-۹ نتایج تخمین دو مدل مختلف رگرسیون را بیان می‌کنیم. به نظر شما در چه

مواردی قطعاً اشتباهاتی در محاسبات صورت گرفته است.

$$R_{y/12}^2 = 0/701 \quad , \quad r_{y2}^2 = 0/126 \quad , \quad r_{y1}^2 = 0/227 \quad (\text{الف})$$

$$(\sum x_i^2) (\sum y_i^2) - (\sum x_i y_i)^2 = -1732/86 \quad (\text{ب})$$

۴-۱۰ مدل زیر را در نظر می‌گیریم که در آن U_t ، تمام فرضهای کلاسیک را شامل است،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t .$$

نشان دهید که

$$R^2 = \frac{\hat{\beta}_1^2 \sum x_{1t}^2 + \hat{\beta}_2^2 \sum x_{2t}^2 + 2 \hat{\beta}_1 \hat{\beta}_2 \sum x_{1t} x_{2t}}{\sum y_t^2} . \quad (\text{الف})$$

$$R^2 = \frac{\hat{\beta}_1 \sum x_{1t} y_t + \hat{\beta}_2 \sum x_{2t} y_t}{\sum y_t^2} . \quad (\text{ب})$$

$$R^2 = \frac{r_{y1}^2 + r_{y2}^2 - 2 r_{y1} r_{y2} r_{12}}{1 - r_{12}^2} \quad (\text{ج})$$

$$r_{y1/2}^2 = \frac{(r_{y1} - r_{y2} r_{12})^2}{(1 - r_{12}^2)(1 - r_{y2}^2)} \quad (\text{د})$$

۴-۱۱ مدل رگرسیون زیر را در نظر می‌گیریم،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t .$$

این مدل را با روش حداقل مربعات معمولی تخمین می‌زنیم. اگر ضریب همبستگی بین $\hat{\beta}_1$ و $\hat{\beta}_2$ را یا $r_{\hat{\beta}_1 \hat{\beta}_2}$ تعریف کنیم، نشان دهید که مجذور ضریب همبستگی بین $\hat{\beta}_1$ و $\hat{\beta}_2$ برابر با مجذور ضریب همبستگی بین X_{1t} و X_{2t} است،

$$r_{\hat{\beta}_1 \hat{\beta}_2} = r_{12} .$$

۴-۱۲ در مدل رگرسیون

$$Y_t = \beta_1 X_{1t} + \beta_2 X_{2t} + U_t ,$$

داده‌های متغیرهای Y_t ، X_{1t} و X_{2t} را تغییر مقیاس داده‌ایم، به طوری که:

$$Y_t^* = \lambda Y_t ,$$

$$X_{1t}^* = \mu_1 X_{1t} ,$$

$$X_{2t}^* = \mu_2 X_{2t} .$$

الف) به نظر شما چه رابطه‌ای بین تخمین پارامترها در مدل اولیه، یعنی $\hat{\beta}_1$ و $\hat{\beta}_2$ یا

تخمین پارامترها در مدل جدید، یعنی $\hat{\beta}_1^*$ و $\hat{\beta}_2^*$ وجود دارد؟

ب) تغییر مقیاس در اندازه‌گیری متغیرها، چه تغییری در R^2 و SEE (خطای معیار

تخمین) ایجاد می‌کند؟

حل مسائل فصل چهارم

۴-۱ با توجه به معادله ۴-۴ می‌دانیم در مدل رگرسیون $Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t$

داریم

$$\hat{\beta}_1 = \frac{\hat{\beta}_{y1} - \hat{\beta}_{y2} \hat{\beta}_{21}}{1 - \hat{\beta}_{12} \hat{\beta}_{21}}$$

یا توجه به تخمین معادلات مفروض داریم

$$\hat{\beta}_{y1} = 0/92 ,$$

$$\hat{\beta}_{y2} = 0/84 ,$$

$$\hat{\beta}_{21} = 0/78 ,$$

$$\hat{\beta}_{12} = 0/55 ,$$

در نتیجه خواهیم داشت

$$\hat{\beta}_1 = \frac{0/92 - 0/78(0/84)}{1 - 0/55(0/78)} = \frac{0/2668}{0/571} = 0/4668 .$$

برای $\hat{\beta}_2$ نیز از معادله ۴-۴ داریم

$$\begin{aligned} \hat{\beta}_2 &= \frac{\hat{\beta}_{y2} - \hat{\beta}_{12} \hat{\beta}_{y1}}{1 - \hat{\beta}_{12} \hat{\beta}_{21}} \\ &= \frac{0/84 - 0/55(0/92)}{1 - 0/55(0/78)} = \frac{0/334}{0/571} = 0/585 \end{aligned}$$

۴-۲ برای نمونه‌های اول و دوم، به ترتیب اندیسه‌های یک و دو می‌گذاریم. در محاسبات زیر اگر متغیری بدون اندیس باشد، بدین معنی است که برای نمونه تلفیقی محاسبه شده است.

$$n = n_1 + n_2 ,$$

$$\bar{X} = \frac{n_1 \bar{X}_1 + n_2 \bar{X}_2}{n_1 + n_2} , \quad \bar{Y} = \frac{n_1 \bar{Y}_1 + n_2 \bar{Y}_2}{n_1 + n_2} .$$

برای محاسبه $\sum x_i^2$ ، ابتدا به این نکته توجه می‌کنیم که

$$\sum_1^{n_1+n_2} x_i^2 = \sum_1^{n_1+n_2} X_i^2 - n \bar{X}^2 . \quad (1)$$

اما می‌دانیم که

$$\sum_1^{n_1+n_2} X_i^2 = \sum_1^{n_1} X_{1i}^2 + \sum_1^{n_2} X_{2i}^2 .$$

باید هر یک از دو جمله سمت راست را محاسبه کنیم،

$$\sum_1^{n_1} X_{1i}^2 = \sum_1^{n_1} x_{1i}^2 + n_1 \bar{X}_1^2 , \quad \sum_1^{n_2} X_{2i}^2 = \sum_1^{n_2} x_{2i}^2 + n_2 \bar{X}_2^2 ,$$

بنابراین

$$\sum X_i^2 = \left(\sum_1^{n_1} x_{1i}^2 + \sum_1^{n_2} x_{2i}^2 \right) + (n_1 \bar{X}_1^2 + n_2 \bar{X}_2^2) .$$

رابطه فوق را در معادله (۱) قرار می‌دهیم،

$$\sum_1^{n_1+n_2} x_i^2 = \left(\sum_1^{n_1} x_{1i}^2 + \sum_1^{n_2} x_{2i}^2 \right) + (n_1 \bar{X}_1^2 + n_2 \bar{X}_2^2 - n \bar{X}^2) . \quad (2)$$

جمله دوم سمت راست عبارت فوق را چنین می‌نویسیم،

$$\begin{aligned} n_1 \bar{X}_1^2 + n_2 \bar{X}_2^2 - n \bar{X}^2 &= n_1 \bar{X}_1^2 + n_2 \bar{X}_2^2 - (n_1 + n_2) \left[\frac{(n_1 \bar{X}_1 + n_2 \bar{X}_2)^2}{(n_1 + n_2)^2} \right] , \\ &= n_1 \bar{X}_1^2 + n_2 \bar{X}_2^2 - \frac{(n_1 \bar{X}_1 + n_2 \bar{X}_2)^2}{(n_1 + n_2)} , \\ &= \frac{n_1 n_2}{n_1 + n_2} (\bar{X}_1^2 - \bar{X}_2^2) . \end{aligned}$$

رابطه فوق را در معادله (۲) جایگزین می‌کنیم،

$$\sum_1^{n_1+n_2} x_{it}^* = \sum_1^{n_1} x_{1t}^* + \sum_1^{n_2} x_{2t}^* + \frac{n_1 n_2}{n_1 + n_2} (\bar{X}_1^* - \bar{X}_2^*) .$$

بدین ترتیب می‌توان نوشت،

$$\sum_1^{n_1+n_2} x_{it}^* y_{it} = \sum_1^{n_1} x_{1t}^* y_{1t} + \sum_1^{n_2} x_{2t}^* y_{2t} + \frac{n_1 n_2}{n_1 + n_2} (\bar{X}_1^* - \bar{X}_2^*) (\bar{Y}_1 - \bar{Y}_2) ,$$

$$\sum_1^{n_1+n_2} y_{it}^* = \sum_1^{n_1} y_{1t}^* + \sum_1^{n_2} y_{2t}^* + \frac{n_1 n_2}{n_1 + n_2} (\bar{Y}_1^* - \bar{Y}_2^*) .$$

با اطلاعات فوق به محاسبات زیر می‌پردازیم.

الف) تخمین معادله رگرسیون برای نمونه اول به صورت زیر محاسبه می‌شود،

$$\hat{Y}_{1t} = -0 + 1/0 X_{1t} ,$$

(۲/۳۷) (۰/۱۲)

$$\sum e_{1t}^* = \text{RSS}_1 = 20 , \quad r_1^* = 0/9 , \quad \text{d.f.} = 18 .$$

برای مدل دوم داریم

$$\hat{Y}_{2t} = -6/0 + 1/0 X_{2t} ,$$

(۲/۱۲) (۰/۱۲)

$$\sum e_{2t}^* = \text{RSS}_2 = 20 , \quad r_2^* = 0/9 , \quad \text{d.f.} = 23 .$$

ب) برای نمونه تلفیقی داریم:

$$n = 40 , \quad \bar{X} = \frac{70}{3} , \quad \bar{Y} = \frac{10}{3} ,$$

$$\sum x_{it}^* = 280 , \quad \sum x_{it}^* y_{it} = 370 , \quad \sum y_{it}^* = 500 ,$$

تخمین مدل رگرسیون تلفیقی، عبارت است از

$$\hat{Y}_t = -1/16 + 1/32 X_t , \quad \text{RSS} = 61/1 , \quad r^* = 0/89 .$$

ج) برای اینکه تخمین پارامترها در مدل رگرسیون تلفیقی معتبر باشد، باید فرض کنیم که مدل‌های رگرسیون نمونه‌های اول و دوم در واقع یکی است،

$$\alpha_1 = \alpha_2, \beta_1 = \beta_2, \sigma_1^2 = \sigma_2^2,$$

در فصل هشتم، قسمت ۵-۸ با عنوان آزمون تغییر ساختاری، خواهیم دید که چگونه این فرضها را می‌توان آزمون کرد. بعلاوه در مورد این مسأله خاص می‌توان نشان داد که فرضیه $\sigma_1^2 = \sigma_2^2$ را نمی‌شود رد کرد، اما فرضیه $\beta_1 = \beta_2$ و $\alpha_1 = \alpha_2$ در سطح معنی‌دار یک درصد، رد می‌شود، چنین نتیجه‌ای بدین معنی است که نباید این دو نمونه را تلفیق کرد. لازم است یادآوری شود که این نتیجه از اهمیت خاصی برخوردار است؛ زیرا با ملاحظه تخمین مدل‌های رگرسیون در هر یک از دو نمونه ملاحظه می‌شود که $\hat{\beta}_1 = 1/5$ و $\hat{\beta}_2 = 1/5$ و تخمین‌های $\hat{\alpha}_1$ و $\hat{\alpha}_2$ نیز تفاوت بسیار مختصری با هم دارند. بنابراین اگر بخواهیم تنها با توجه به صورت تخمین معادله‌ها دربارهٔ مجوز تلفیق نظر بدهیم، قطعاً جواب مثبت است؛ زیرا این دو تخمین بسیار نزدیک به هم بوده و بر این دلالت می‌کند که از یک جامعه آماری اخذ شده است. اما در آزمون فرضیه مشاهده می‌شود که با وجود $\hat{\beta}_1 = \hat{\beta}_2 = 1/5$ ، حتی نمی‌توان فرضیه $\beta_1 = \beta_2$ را قبول کرد؛ بنابراین نتایج حاصل از تخمینها به تنهایی نمی‌تواند ما را به استنتاج آماری برساند و باید حتماً فرضیه‌های مختلف را آزمون کرد. ادامهٔ مباحث مربوط به این مسأله را می‌توان در مسألهٔ ۸-۸ ملاحظه کرد.

۴-۳ الف) هرگاه Z_1 از X_1 مستقل باشد، اضافه کردن Z_1 به مدل رگرسیون اول، در تخمین β تغییری نمی‌دهد؛ بنابراین موقعی $\hat{\beta} = \hat{\beta}'$ که Z_1 و X_1 همبستگی نداشته باشد.

ب) می‌دانیم اگر تعداد متغیرهای توضیحی اضافه شود، مجموع مربعات پسماند کاهش می‌یابد یا حداقل تغییری نمی‌کند؛ بنابراین رابطهٔ $\sum e_i^2 \geq \sum e_i'^2$ همواره صادق است.

ج) وقتی متغیر جدیدی که اضافه می‌کنیم (مثلاً Z_1) همبستگی خیلی زیاد با X_1

داشته باشد، ممکن است معنی دار بودن تخمین ضریب X_1 در مدل جدید، یعنی $\hat{\beta}'$ ، را تغییر دهد.

د) این مسأله، حالتی است که یک متغیر جدید، می تواند ضریبی را که معنی دار نبوده است، معنی دار کند؛ یعنی می تواند ضریبی را که مثلاً در سطح معنی دار ۵ درصد رد شده است قابل قبول کند. این حالت در شرایطی ایجاد می شود که متغیر جدید، یعنی Z_1 ، اولاً با متغیر قدیمی یعنی X_1 همبستگی ضعیفی داشته باشد و ثانیاً از قدرت توضیحی بسیار قوی برخوردار باشد. برای تبیین این نکته، واریانس $\hat{\beta}$ و $\hat{\beta}'$ را می نویسیم.

$$\text{Var}(\hat{\beta}) = \frac{\sum e_i^2}{(n-2) \sum x_i^2}, \quad \text{Var}(\hat{\beta}') = \frac{\sum e_i^2}{(n-3) \sum x_i^2 (1 - r_{x,z}^2)}$$

فرض می کنیم $\text{Var}(\hat{\beta})$ بقدری زیاد است که در آزمونهای آماری، فرضیه $\beta = 0$ رد نشده است؛ یعنی $\hat{\beta}$ معنی دار نیست. اما چون ورود متغیر Z_1 که قدرت توضیحی بسیاری هم دارد - باعث می شود که مجموع مربعات پسماند در مدل جدید ($\sum e_i^2$) کمتر از $\sum e_i^2$ بشود؛ در نتیجه قطعاً $\text{Var}(\hat{\beta}')$ از $\text{Var}(\hat{\beta})$ کمتر خواهد شد. بنابراین چه بسا در آزمونهای آماری $\hat{\beta}'$ معنی دار شود در حالی که قبلاً معنی دار نبود. البته این امر به شرطی واقع می شود که ضریب همبستگی بین X_1 و Z_1 بسیار کم باشد، این $r_{x,z}^2$ بسیار کوچک باشد؛ زیرا در غیر این صورت، با اینکه صورت کسر ($\text{Var}(\hat{\beta}')$) کم می شود، اما مخرج کسر نیز کم شده، هیچ دلیلی در دست نخواهد بود که $\text{Var}(\hat{\beta}')$ کاهش یابد.

۴-۴ با توجه به معادله ۴-۱۰ می دانیم اگر β_1 را از مدل اول تخمین بزنیم، خواهیم داشت

$$\hat{\beta}_1 = \frac{\sum x_{1i}^2 \sum x_{1i} y_i - \sum x_{1i} x_{2i} \sum x_{2i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2}$$

نشان می دهیم که نتیجه تخمین $\hat{\beta}_1$ از مدل سوم نیز دقیقاً به همین نتیجه خواهد رسید. برای تخمین $\hat{\beta}_1$ از مدل سوم، می توان از مدل زیر استفاده کرد،

$$Y_i = \alpha'' + \beta_1' \hat{V}_i + \epsilon_i'$$

زیرا می‌دانیم \hat{V}_i پسماند رگرسیون X_{2i} بر X_{1i} است و با توجه به معادله‌های نرمال می‌توان گفت که \hat{V}_i از X_{2i} مستقل است؛ بنابراین وقتی در معادله سوم، متغیر X_{2i} از \hat{V}_i مستقل باشد، می‌توان به جای تخمین یک مدل رگرسیون دو متغیره، β_1 را با استفاده از مدل رگرسیون ساده فوق تخمین زد:

$$\hat{\beta}_1 = \frac{\sum y_i \hat{v}_i}{\sum \hat{v}_i^2}$$

می‌دانیم $\hat{v}_i = x_{2i} - \hat{\lambda} x_{1i}$ با جایگزینی در معادله فوق داریم

$$\hat{\beta}_1 = \frac{\sum y_i (x_{2i} - \hat{\lambda} x_{1i})}{\sum (x_{2i} - \hat{\lambda} x_{1i})^2}$$

مقدار $\hat{\lambda} = \frac{\sum x_{1i} x_{2i}}{\sum x_{1i}^2}$ را در معادله فوق قرار می‌دهیم،

$$\hat{\beta}_1 = \frac{\sum x_{2i}^2 \sum y_i x_{1i} - \sum x_{1i} x_{2i} \sum x_{2i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2}$$

که دقیقاً همان فرمول $\hat{\beta}_1$ است؛ بنابراین $\hat{\beta}_1 = \beta_1$. توضیح این نتیجه در متن کتاب آمده است.

۴.۵ الف) تخمین β از مدل اولیه به صورت زیر است،

$$\hat{\beta} = \frac{\sum x_{2i}^2 \sum x_{1i} y_i - \sum x_{1i} x_{2i} \sum x_{2i} y_i}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2} \quad (1)$$

رگرسیون y_i روی x_{2i} عبارت است از

$$y_i = \alpha x_{2i} + \varepsilon_i$$

و پسماندهای آن برابر است با

$$y_i^e = y_i - \hat{\alpha} x_{2i}$$

رگرسیون x_{1i} روی x_{2i} عبارت است از

$$x_{1i} = \lambda x_{2i} + v_i \quad (2)$$

و پسماندهای آن برابر است با

$$x_{1t}^e = x_{1t} - \hat{\lambda} x_{2t} .$$

رگرسیون y_t^e روی x_{1t}^e عبارت است از

$$y_t^e = \rho x_{1t}^e + w_t , \quad (3)$$

و پسماندهای آن برابر است با

$$y_t^{e0} = y_t^e - \hat{\rho} x_{1t}^e . \quad (4)$$

باید ثابت کنیم که $\hat{\rho} = \hat{\beta}$. برای این منظور، به جای x_{1t}^e و y_t^e در معادله (۳) مقادیرشان را قرار می‌دهیم،

$$(y_t - \hat{\alpha} x_{2t}) = \rho (x_{1t} - \hat{\lambda} x_{2t}) + w_t .$$

$\hat{\rho}$ از این معادله برابر است با

$$\hat{\rho} = \frac{\sum (y_t - \hat{\alpha} x_{2t}) (x_{1t} - \hat{\lambda} x_{2t})}{\sum (x_{1t} - \hat{\lambda} x_{2t})^2} . \quad (5)$$

می‌دانیم $(x_{1t} - \hat{\lambda} x_{2t})$ در واقع پسماند رگرسیون در معادله (۲) است که بنا بر معادلات نرمال از متغیر توضیحی یعنی x_{2t} مستقل است؛ بنابراین

$$\hat{\alpha} \sum x_{2t} (x_{1t} - \hat{\lambda} x_{2t}) = 0 .$$

با جایگزینی رابطه فوق در معادله (۵)، خواهیم داشت

$$\hat{\rho} = \frac{\sum y_t (x_{1t} - \hat{\lambda} x_{2t})}{\sum (x_{1t} - \hat{\lambda} x_{2t})^2} .$$

اما از معادله (۲) می‌دانیم که

$$\hat{\lambda} = \frac{\sum x_{1t} x_{2t}}{\sum x_{2t}^2} ,$$

در نتیجه خواهیم داشت

$$\hat{\rho} = \frac{\sum x_{1t}^2 x_{1t} y_t - \sum x_{1t} x_{2t} \sum x_{2t} y_t}{\sum x_{1t}^2 \sum x_{2t}^2 - (\sum x_{1t} x_{2t})^2}$$

که دقیقاً همان معادله (۱) است؛ یعنی $\hat{\beta} = \hat{\rho}$.

ب) پسماندهای مدل مفروض عبارت است از $y_t - \hat{\beta} x_{1t} - \hat{\gamma} x_{2t}$. با توجه به معادله (۴)، باید ثابت کنیم که

$$y_t - \hat{\rho} x_{1t} = y_t - \hat{\beta} x_{1t} - \hat{\gamma} x_{2t} \quad (۶)$$

برای این منظور ابتدا $(y_t - \rho x_{1t})$ را بسط می‌دهیم،

$$y_t - \hat{\rho} x_{1t} = (y_t - \hat{\alpha} x_{2t}) - \hat{\rho} (x_{1t} - \hat{\lambda} x_{2t})$$

در قسمت اول ثابت کردیم که $\hat{\beta} = \hat{\rho}$. در نتیجه داریم

$$y_t - \hat{\rho} x_{1t} = y_t - \hat{\beta} x_{1t} - (\hat{\alpha} - \hat{\beta} \hat{\lambda}) x_{2t}$$

بنابراین برای اثبات معادله (۶) باید ثابت کنیم که

$$(\hat{\alpha} - \hat{\beta} \hat{\lambda}) = \hat{\gamma} \quad (۷)$$

برای اثبات رابطه (۷) کافی است از معادله‌های ۴-۴۱ و ۴-۴۲ استفاده کنیم. بر اساس این معادله‌ها خواهیم داشت

$$\hat{\beta} = \frac{\hat{\beta}_{yx_1} - \hat{\beta}_{x_2x_1} \hat{\alpha}}{1 - \hat{\lambda} \hat{\beta}_{x_2x_1}}, \quad \hat{\gamma} = \frac{\hat{\alpha} - \hat{\lambda} \hat{\beta}_{yx_1}}{1 - \hat{\lambda} \hat{\beta}_{x_2x_1}} \quad (۸)$$

بنابراین تعداد $(\hat{\alpha} - \hat{\beta} \hat{\lambda})$ برابر است با

$$(\hat{\alpha} - \hat{\beta} \hat{\lambda}) = \hat{\alpha} - \frac{(\hat{\lambda} \hat{\beta}_{yx_1} - \hat{\beta}_{x_2x_1} \hat{\alpha})}{1 - \hat{\lambda} \hat{\beta}_{x_2x_1}}$$

$$= \frac{\hat{\alpha} - \hat{\alpha} \hat{\lambda} \hat{\beta}_{x_1 x_1} - \hat{\lambda} \hat{\beta}_{y x_1} + \hat{\alpha} \hat{\lambda} \hat{\beta}_{x_1 x_1}}{1 - \hat{\lambda} \hat{\beta}_{x_1 x_1}}$$

$$= \frac{\hat{\alpha} - \hat{\lambda} \hat{\beta}_{y x_1}}{1 - \hat{\lambda} \hat{\beta}_{x_1 x_1}}$$

با توجه به معادله (۸)، خواهیم داشت

$$(\hat{\alpha} - \hat{\beta} \hat{\lambda}) = \hat{\lambda}$$

که دقیقاً اثبات رابطه (۷) است.

۴-۶ الف) می‌دانیم اگر یک مدل رگرسیون را با روش حداقل مربعات معمولی تخمین بزنیم، مقادیر پسماند از متغیرهای توضیحی مستقل است؛ بنابراین X_{1t}^* - که پسماند مدل رگرسیون مفروض است - از متغیر توضیحی X_{1t} مستقل خواهد بود. در چنین حالتی، رگرسیون Y_t بر روی X_{1t}^* و X_{1t} دقیقاً به همان تخمینی از β_1 می‌رسد که رگرسیون Y_t بر روی X_{1t} نتیجه خواهد داد.

ب) برای تعیین مقدار اریب $\hat{\beta}_1$ ، ابتدا فرمول آن را می‌نویسیم،

$$\hat{\beta}_1 = \frac{\sum x_{1t} y_t}{\sum x_{1t}^2}$$

مقدار y_t را از مدل اولیه (مدل صحیح) در فرمول فوق قرار می‌دهیم،

$$\hat{\beta}_1 = \frac{\sum x_{1t} (\beta_1 x_{1t} + \beta_2 x_{2t} + U_t)}{\sum x_{1t}^2}$$

$$= \beta_1 + \beta_2 \frac{\sum x_{1t} x_{2t}}{\sum x_{1t}^2} + \frac{\sum x_{1t} U_t}{\sum x_{1t}^2}$$

از دو طرف رابطه بالا امید ریاضی می‌گیریم. یا توجه به $E(\sum x_{1t} U_t) = 0$ ، خواهیم داشت

$$E(\hat{\beta}_\gamma) = \beta_\gamma + \beta_1 \frac{\sum x_{1t} x_{\gamma t}}{\sum x_{\gamma t}^2}$$

اگر یک مدل رگرسیون بسازیم که در آن X_{1t} بر حسب $X_{\gamma t}$ بیان شده باشد،

$$X_{1t} = \gamma + \beta_{1\gamma} X_{\gamma t} + \varepsilon_t$$

آنگاه

$$\hat{\beta}_{1\gamma} = \frac{\sum x_{1t} x_{\gamma t}}{\sum x_{\gamma t}^2}$$

نتیجه می‌گیریم که

$$\text{مقدار اریب} = E(\hat{\beta}_\gamma) - \beta_\gamma = \beta_1 \hat{\beta}_{1\gamma}$$

۴-۷ با توجه به $C_t = Y_t - S_t$ ، به وضوح ملاحظه می‌شود که در تخمین تابع مصرف خواهیم داشت: $\hat{\alpha} = 0$ ، $\hat{\beta} = 1$ و $\hat{\gamma} = -1$. ضریب تعیین (R^2) نیز برابر یک خواهد بود. همچنین مجموع پسماندها نیز صفر می‌شود؛ بنابراین نتیجه می‌گیریم که مدل مفروض در واقع تابع تغییرات ساختاری مصرف را بیان نمی‌کند، بلکه مبتنی بر یک رابطهٔ حسابداری است. در این گونه موارد باید در مدل اصلی تجدیدنظر کرد به گونه‌ای که یکی از متغیرهای توضیحی حذف شود. استفاده از نظریه‌های اقتصادی، باید محور اصلی این تجدیدنظرها باشد.

۴-۸ الف) $R_{y/1234}^2 = 0/86$ بر این دلالت می‌کند که ضریب تعیین در یک مدل رگرسیون که Y_t را بر حسب چهار متغیر توضیحی X_{1t} ، X_{2t} ، X_{3t} و X_{4t} بیان می‌کند، برابر است با $0/86$. اما $R_{y/123}^2 = 0/89$ نشان می‌دهد که Y_t بر حسب سه متغیر توضیحی X_{1t} ، X_{2t} و X_{3t} بیان شده است؛ با وجود این، ضریب تعیین آن $0/89$ است. می‌دانیم که اگر تعداد متغیرهای توضیحی اضافه شود، قدرت توضیحی مدل بیشتر شده یا حداقل ثابت باقی می‌ماند؛ بنابراین $R_{y/1234}^2$ نمی‌تواند کمتر از $R_{y/123}^2$ شود. نتیجه می‌گیریم که در مورد الف) قطعاً اشتباهی واقع شده است.

ب) $R_{y/123}^2$ نشان می‌دهد که در مدل رگرسیون مفروض سه متغیر توضیحی وجود

دارد. اطلاعات موجود فقط در مورد دو ضریب تعیین جزئی ارائه شده است؛ بنابراین هیچ راهی وجود ندارد که بتوان صحت رابطه $R_{y/12}^2 = 0.70$ را با توجه به دو اطلاع داده شده ارزیابی کرد.

۴-۹ الف) با توجه به معادله ۴-۶۰ داریم

$$R_{y/12}^2 = \frac{r_{y1}^2 + r_{y2}^2 - 2 r_{y1} r_{y2} r_{11}}{(1 - r_{12}^2)}$$

بنابراین چهار متغیر در ارتباط با یکدیگر قرار دارند. در قسمت (الف) این مسأله سه کمیت در مورد سه متغیر ارائه شده است: r_{y1}^2 ، r_{y2}^2 و $R_{y/12}^2$. اما مقدار r_{12} که در واقع همان r_{21} است - معین نشده است و می تواند مقادیر متعددی داشته باشد که با توجه به معادله فوق قابل محاسبه است. کافی است این مقادیر را در معادله فوق قرار دهیم. خواهیم داشت

$$0.701 = \frac{0.227 + 0.126 - 2(\sqrt{0.227(0.126)})r_{11}}{(1 - r_{12}^2)}$$

بعد از ساده کردن داریم

$$0.701 r_{12}^2 - 0.3282 r_{12} - 0.348 = 0$$

ریشه های این معادله برابر است با

$$r_{12} = 0.987 \quad , \quad r_{12} = -0.504$$

مشاهده می شود که هر دو مقدار r_{12} ، مقادیر قابل قبولی از ضریب همبستگی است؛ بنابراین نتیجه می گیریم که داده های مسأله می تواند صحیح باشد.

ب) داده های مسأله در این قسمت غلط است؛ زیرا می دانیم

$$r_{x,y}^2 = \frac{(\sum x_i y_i)^2}{\sum x_i^2 \sum y_i^2}$$

از طرف دیگر می‌دانیم

$$0 \leq r_{x,y}^2 \leq 1,$$

باید داشته باشیم:

$$\sum x_t^2 \sum y_t^2 \geq (\sum x_t y_t)^2,$$

یا

$$\sum x_t^2 \sum y_t^2 - (\sum x_t y_t)^2 \geq 0.$$

اما مقدار داده شده برای عبارت فوق برابر $1732/86$ است؛ پس نتیجه می‌گیریم که قطعاً اشتباهاتی در محاسبات واقع شده است.

۴-۱۰ الف) از تعریف R^2 شروع می‌کنیم،

$$R^2 = \frac{ESS}{TSS} = \frac{\sum \hat{y}_t^2}{\sum y_t^2}.$$

به جای \hat{y}_t مقدار آن را از رابطه زیر قرار می‌دهیم،

$$\hat{y}_t = \hat{\beta}_1 x_{1t} + \hat{\beta}_2 x_{2t}.$$

نتیجه می‌گیریم که

$$R^2 = \frac{\hat{\beta}_1^2 \sum x_{1t}^2 + \hat{\beta}_2^2 \sum x_{2t}^2 + 2 \hat{\beta}_1 \hat{\beta}_2 \sum x_{1t} x_{2t}}{\sum y_t^2}.$$

ب) دقیقاً همان فرمول ۴-۱۹ در متن کتاب است.

ج) دقیقاً همان فرمول ۴-۶۰ در متن کتاب است.

د) اگر فرمول ۴-۵۳ را مجذور کنیم، این رابطه به اثبات می‌رسد.

۴-۱۱ با استفاده از تعریف ضریب همبستگی بین دو متغیر، می‌توان چنین نوشت

$$r_{\hat{\beta}_1, \hat{\beta}_2} = \frac{E[(\hat{\beta}_1 - \beta_1)(\hat{\beta}_2 - \beta_2)]}{\sqrt{E(\hat{\beta}_1 - \beta_1)^2 E(\hat{\beta}_2 - \beta_2)^2}},$$

یا

$$r_{\hat{\beta}_1, \hat{\beta}_2}^2 = \frac{\{E[(\hat{\beta}_1 - \beta_1)(\hat{\beta}_2 - \beta_2)]\}^2}{E(\hat{\beta}_1 - \beta_1)^2 E(\hat{\beta}_2 - \beta_2)^2}$$

$$= \frac{[\text{Cov}(\hat{\beta}_1, \hat{\beta}_2)]^2}{\text{Var}(\hat{\beta}_1) \text{Var}(\hat{\beta}_2)}$$

از طرف دیگر، با توجه به معادله ۴-۲۶ می‌دانیم

$$\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = \frac{-\sigma^2 r_{12}^2}{\sum x_{1t} x_{2t} (1 - r_{12}^2)} = \frac{-\sigma^2 (\sum x_{1t} x_{2t})^2}{\sum x_{1t}^2 \sum x_{2t}^2 \sum x_{1t} x_{2t} (1 - r_{12}^2)}$$

$$= \frac{-\sigma^2 \sum x_{1t} x_{2t}}{\sum x_{1t}^2 \sum x_{2t}^2 (1 - r_{12}^2)}$$

با استفاده از فرمول فوق و نیز فرمولهای ۴-۲۴ و ۴-۲۵ و جایگزینی در فرمول $r_{\hat{\beta}_1, \hat{\beta}_2}^2$ خواهیم داشت

$$r_{\hat{\beta}_1, \hat{\beta}_2}^2 = \frac{\sigma^4 (\sum x_{1t} x_{2t})^2 (\sum x_{1t}^2) (1 - r_{12}^2) (\sum x_{2t}^2) (1 - r_{12}^2)}{(\sum x_{1t}^2)^2 (\sum x_{2t}^2)^2 (1 - r_{12}^2)^2 \sigma^4 \sigma^4}$$

$$= \frac{(\sum x_{1t} x_{2t})^2}{\sum x_{1t}^2 \sum x_{2t}^2} = r_{12}^2$$

۴-۱۲ الف) مدل جدید را که در آن مقیاس اندازه‌گیری متغیرها تغییر یافته است، به

صورت زیر می‌نویسیم،

$$\frac{Y_t^*}{\lambda} = \frac{\beta_1 X_{1t}^*}{\mu_1} + \frac{\beta_2 X_{2t}^*}{\mu_2} + U_t$$

یا

$$Y_t^* = \frac{\lambda \beta_1 X_{1t}^*}{\mu_1} + \frac{\lambda \beta_2 X_{2t}^*}{\mu_2} + \lambda U_t$$

$$= \beta_1^* X_{1t}^* + \beta_2^* X_{2t}^* + U_t^* \quad (1)$$

تخمین β_1^* از مدل جدید عبارت است از

$$\hat{\beta}_1^* = \left(\frac{\lambda \beta_1}{\mu_1} \right) = \frac{\lambda}{\mu_1} \hat{\beta}_1$$

در نتیجه

$$\hat{\beta}_1^* = \frac{\lambda}{\mu_1} \hat{\beta}_1$$

به همین ترتیب می توان نشان داد که

$$\hat{\beta}_2^* = \frac{\lambda}{\mu_2} \hat{\beta}_2$$

ب) اگر ضریب تعیین در مدل جدید را با R^{*2} نشان دهیم، با توجه به معادله ۴-۲۰ خواهیم داشت

$$R^{*2} = \frac{\sum (Y_t^* \hat{Y}_t^*)^2}{\sum Y_t^{*2} \sum \hat{Y}_t^{*2}}$$

عبارت $Y_t^* = \lambda Y_t$ و $\hat{Y}_t^* = \lambda \hat{Y}_t$ را در فرمول فوق جایگزین می کنیم،

$$R^{*2} = \frac{\sum (\lambda Y_t \lambda \hat{Y}_t)^2}{\sum \lambda^2 Y_t^2 \sum \lambda^2 \hat{Y}_t^2} = \frac{\sum (Y_t \hat{Y}_t)^2}{\sum Y_t^2 \sum \hat{Y}_t^2} = R^2$$

نتیجه می گیریم که تغییر در مقیاس اندازه گیری متغیرها در ضریب تعیین تأثیری ندارد. همچنین اگر انحراف معیار رگرسیون در مدل جدید را با SEE^* نشان دهیم،

با توجه به معادله ۱-۴۵ داریم

$$SEE^* = \sqrt{\frac{\sum e_t^{*2}}{n-2}} = \sqrt{\frac{\sum \hat{U}_t^{*2}}{n-2}}$$

با توجه به معادله (۱)، می‌دانیم: $\hat{U}_t^* = \lambda \hat{U}_t^*$

$$SEE^* = \sqrt{\frac{\sum (\lambda \hat{U}_t^{*2})}{(n-2)}} = \lambda \sqrt{\frac{\sum \hat{U}_t^{*2}}{n-2}} = \lambda \sqrt{\frac{\sum e_t^2}{n-2}}$$

در نتیجه خواهیم داشت

$$SEE^* = \lambda SEE$$

یعنی تغییر مقیاس در متغیر درون‌زا، انحراف معیار رگرسیون را به اندازه λ تغییر می‌دهد.

مفاهیم و تخمین مدل رگرسیون خطی چند متغیره

۵-۱ مقدمه

مدلهایی که تا به حال مطالعه کرده ایم، شامل یک متغیر توضیحی و به صورت

$$Y_t = \alpha + \beta X_t + U_t,$$

یا دو متغیر توضیحی و از نوع:

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + U_t,$$

بوده اند. اما می توان فرض کرد که در یک مدل رگرسیون، k متغیر توضیحی وجود داشته باشد،

$$Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_k X_{kt}$$

برای سهولت در محاسبات، تخمین این گونه مدلها با جبر ماتریسی انجام می شود؛ بنابراین باید ابتدا مدل رگرسیون چند متغیره را به زبان ماتریسی نوشت، سپس آن را با روش حداقل مربعات معمولی تخمین زد. یک پیوست ماتریسی نیز در بخش پیوستها در نظر گرفته شده است^۱ تا بتوان تعاریف و عملیات اصلی ماتریسها را در حد ضرورت مرور کرد. بنابراین آشنایی با جبر ماتریسی پیش نیاز این فصل محسوب می شود. با وجود این، سعی شده است بسیاری از مطالب اصلی به گونه ای ارائه شود که آشنایی با مفاهیم و عملیات مقدماتی ماتریسها کافی باشد. در قسمت ۵-۲ یک مدل رگرسیون چند متغیره و نیز فرضهای کلاسیک آن به زبان ماتریسی ارائه شده است. رسیدن به فرمولهای تخمین

پارامترها در رگرسیون چند متغیره، با استفاده از مشتق‌گیری ماتریسی بسیار ساده است. باید این نکته را نیز در نظر گرفت که محاسبات تخمین پارامترها می‌تواند برحسب مشاهدات اصلی یا انحراف از میانگین باشد. توجه به خصوصیات ماتریسها در هریک از این دو حالت اهمیت فراوان دارد. این نکات موضوع قسمت ۳-۵ است.

ضریب تعیین در مدل‌های رگرسیون چند متغیره، موضوع قسمت ۴-۵ خواهد بود. در این قسمت فرمول‌های مختلفی برای محاسبه R^2 به زبان ماتریسی ارائه شده است. در فصل اول دیدیم که اگر فرمول عمومی ضریب تعیین را برای مدلی به کار بریم که فاقد جمله ثابت است، چه بسا ممکن است، مقدار R^2 منفی شود. در این قسمت، این نکته را برای مدل‌های رگرسیون چند متغیره تعمیم داده‌ایم و همچنین شرایطی را که R^2 می‌تواند منفی شود، ارزیابی کرده‌ایم.

یکی از مسائل مهمی که در رگرسیون‌های چند متغیره مطرح می‌شود، افزایش یا حداقل عدم کاهش R^2 به ازای افزایش متغیرهای توضیحی است. این مسأله را -که ریشه در تعریف R^2 دارد- می‌توان با وارد کردن درجات آزادی در تعریف R^2 تا حدی تعدیل کرد. به همین دلیل در مواردی باید ضرورتاً از معیار ضریب تعیین تعدیل شده یا \bar{R}^2 استفاده کرد. مباحث \bar{R}^2 و رابطه آن با R^2 در قسمت ۵-۵ مطرح شده است. همچنین در این قسمت نشان داده‌ایم که آن دسته از متغیرهای توضیحی که واریانس جمله اختلال را حداقل می‌کنند، ضریب تعیین تعدیل شده را حداکثر خواهند کرد. بدیهی است این امر می‌تواند به منزله یکی از معیارهای مفید در حفظ یا حذف متغیرهای توضیحی مورد استفاده قرار گیرد.

در قسمت ۶-۵ نکاتی پیرامون R^2 و \bar{R}^2 مطرح شده است. معمول این است که در تخمین یک مدل رگرسیون، هرگاه مقدار R^2 زیاد باشد، این امر را نشانه‌ای از موفقیت یا خوبی تخمین می‌دانند. در این قسمت سعی بر این است تا به این سؤال پاسخ داده شود که آیا چنین استنتاجی صحیح است؟ مسأله حداکثرسازی R^2 یا استفاده از R^2 در مقایسه تخمین‌های مختلف از یک مدل یا مقایسه تخمین مدل‌های مختلف، یکی از مباحث بحث‌انگیز اقتصادسنجی است، بنابراین سعی شده است که در این قسمت ابعاد مختلف

آن روشتر شود. سرانجام در قسمت ۵.۷ مسائل بیان ماتریسی R^1 بر حسب مقادیر اصلی متغیرها به طور خلاصه بررسی خواهد شد.

۵.۲ بیان ماتریسی رگرسیون چندمتغیره
 مدل رگرسیون چندمتغیره زیر را ملاحظه کنید،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i$$

که در آن $t = 1, 2, \dots, n$. اگر بخواهیم این مدل را به زبان ماتریسی بیان کنیم، پارامترهای مدل باید با یک بردار مشخص شود. انتظار این است که این بردار را β بنامیم. در این صورت عناصر بردار β عبارت است از:

$$\beta' = [\alpha, \beta_1, \beta_2, \dots, \beta_k] \cdot$$

بیان فوق از بردار β چندان مرسوم نیست؛ زیرا اولین عنصر آن α است؛ بنابراین باید نوع نوشتن مدل رگرسیون چندمتغیره را تغییر دهیم. یک راه حل این است که مدل را به صورت زیر بنویسیم،

$$Y_i = \beta_1 + \beta_2 X_{1i} + \beta_3 X_{2i} + \dots + \beta_{k+1} X_{ki} + U_i \cdot$$

در این حالت، بردار β عبارت خواهد بود از:

$$\beta' = [\beta_1, \beta_2, \beta_3, \dots, \beta_k] \cdot$$

ملاحظه می شود که مشکل فوق برای β رفع شده است؛ زیرا دیگر پارامتر α در بین عناصر بردار β وجود ندارد. اما مسأله دیگری ایجاد شده و آن عدم تطابق اندیس پارامترهای β با اندیس متغیرهای توضیحی است؛ به عبارت دیگر، معمول این است که مثلاً β_2 ضریب X_{1i} باشد، اما در اینجا β_2 ضریب X_{1i} فرض شده است. برای رفع این مشکل بهترین راه حل این است که مدل رگرسیون را به صورت زیر بنویسیم،

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i \cdot$$

که در آن y بردار مقادیر متغیر درون‌زا و $n \times 1$ است $(y \rightarrow n \times 1)$ ، و X ماتریس مقادیر متغیرهای برون‌زا یا توضیحی و $n \times k$ است، $(X \rightarrow n \times k)$. پارامترهای مدل با بردار β مشخص شده است؛ بنابراین $k \times 1$ خواهد بود $(\beta \rightarrow k \times 1)$ ، سرانجام بردار u شامل مقادیر جمله‌های اختلال بوده، $n \times 1$ می‌باشد $(u \rightarrow n \times 1)$. لازم است یادآوری شود که عناصر ماتریس X برخلاف روال معمول در ماتریسها نوشته شده است. می‌دانیم در ماتریسها، عنصر a_{ij} بیان‌کننده مقداری است که در ردیف i و ستون j قرار دارد؛ مثلاً X_{22} باید عنصری باشد که در ردیف سوم و ستون دوم قرار می‌گیرد. اما در معادله ۵-۴ ملاحظه می‌شود که مثلاً عنصر X_{22} در ماتریس X ، در ردیف دوم و ستون سوم واقع شده است. این جابجایی به طور اساسی مسأله خاصی را ایجاد نمی‌کند و بهتر است به همین شکل حفظ شود؛ زیرا در غیر این صورت علامتگذاری را در معادلات پیچیده‌تر خواهد کرد.

معادله ۵-۳ را می‌توان به صورت زیر نیز نوشت،

$$y = \beta_1 + \beta_2 X_1 + \beta_3 X_2 + \dots + \beta_k X_k + u \quad (5.5)$$

که در آن $X_k \rightarrow n \times 1$ معادله ۵-۵ بر این دلالت می‌کند که بردار مشاهدات متغیر درون‌زا برابر با مجموع بردار جمله‌های اختلال u و یک ترکیب خطی از ستونهای ماتریس X است.

بعد از آشنایی با بیان ماتریسی یک مدل رگرسیون چندمتغیره، به بررسی فرضهای کلاسیک در مدل‌های رگرسیون چندمتغیره می‌پردازیم.

۱. فرضهای کلاسیک مدل‌های رگرسیون چندمتغیره

در این قسمت فرضهای کلاسیک جمله‌های اختلال در بردار u و نیز فرضهای کلاسیک متغیرهای توضیحی را بررسی می‌کنیم. در فصل چهارم دیدیم که باید فرضهای میانگین، واریانس، کوواریانس و سرانجام تابع توزیع احتمال U را بررسی کنیم. در مورد متغیرهای توضیحی نیز دو فرض را مطرح کردیم. روال بحث دقیقاً مانند مطالب مطرح

شده در قسمت ۲-۴ است، با این تفاوت که در این قسمت از بیان ماتریسی استفاده خواهد شد.

ابتدا می‌گوییم که در معادله ۵-۵، مقادیر $\beta_1, \beta_2, \dots, \beta_k$ به عنوان پارامترهای جامعه، برای ما مجهول است و باید آنها را تخمین بزنیم. حتی اگر این پارامترها معلوم باشد، ترکیب خطی $(\beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k)$ نمی‌تواند بردار y را به طور دقیق برای ما مشخص کند؛ زیرا قطعاً متغیرهای دیگری هم وجود دارد که می‌تواند تغییرات متغیر درون‌زا را توضیح دهد، اما در مدل ما وارد نشده است. معادله‌ای که می‌خواهد رفتار یک متغیر اقتصادی را توضیح دهد، معمولاً دقیق و معین نیست بلکه تصادفی است. نتیجه می‌گیریم که بردار u در مدل ۵-۵ شامل متغیرهای تصادفی است و در نتیجه یک بردار تصادفی خواهد بود. سؤال این است که میانگین، واریانس، کوواریانس و تابع توزیع احتمال آن چیست؟ لازم است توضیح داده شود که اطلاعات ما از خصوصیات u در نهایت روشن‌کننده این نکته خواهد بود که بردار y چگونه ایجاد شده است.

اولین فرض، به میانگین یا امید ریاضی u مربوط است. مانند آنچه در فصل اول و در نمودار ۱-۱ و معادله ۱-۹ مطرح شد، می‌توان گفت که در معادله‌های ۲-۵، برای هریک از مقادیر $U_t, t = 1, 2, \dots, n$ ، فرض بر این است که $E(U_t) = 0$ ؛ بنابراین

$$\begin{bmatrix} E(U_1) \\ E(U_2) \\ \vdots \\ E(U_n) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$

یا

$$E(u) = 0 \quad (5.6)$$

که در آن $n \times 1 \rightarrow 0$. از دو طرف معادله ۵-۴، امید ریاضی می‌گیریم. داریم

$$E(y) = E(X\beta) + E(u) .$$

با توجه به معادله ۵-۶ و نیز در نظر گرفتن این نکته که X در آزمایشهای فرضی تکراری ثابت نگه داشته می شود، خواهیم داشت:

$$E(y) = X\beta \quad (5.7)$$

معادله ۵-۷ نتیجه مستقیم اولین فرض u است.

دومین و سومین فرض u به واریانس و کوواریانس آن مربوط است. ابتدا ماتریسی می سازیم که عناصر قطری و غیرقطری آن به ترتیب واریانس و کوواریانس عناصر بردار u باشد. این ماتریس را معمولاً با Σ_u نشان می دهند. برای ساختن این ماتریس باید ابتدا ترانهاد بردار u را حساب کرده، سپس ماتریس uu' را تشکیل دهیم. چون $u \rightarrow n \times 1$ ، بنابراین $u' \rightarrow 1 \times n$ ، و لذا $uu' \rightarrow n \times n$. اگر از ماتریس uu' امید ریاضی بگیریم، ماتریس واریانس - کوواریانس مقادیر مختلف جمله های اختلال به ترتیب زیر به دست می آید،

$$E(uu') = E \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_n \end{bmatrix} [U_1 \ U_2 \ \dots \ U_n] = E \begin{bmatrix} U_1' & U_1 U_2 & \dots & U_1 U_n \\ U_2 U_1 & U_2' & \dots & U_2 U_n \\ \vdots & \vdots & \ddots & \vdots \\ U_n U_1 & U_n U_2 & \dots & U_n' \end{bmatrix} \quad \text{یا}$$

$$E(uu') = \begin{bmatrix} E(U_1') & E(U_1 U_2) & \dots & E(U_1 U_n) \\ E(U_2 U_1) & E(U_2') & \dots & E(U_2 U_n) \\ \vdots & \vdots & \ddots & \vdots \\ E(U_n U_1) & E(U_n U_2) & \dots & E(U_n') \end{bmatrix}$$

می دانیم واریانس U_1 برابر است با

$$\begin{aligned} \text{Var}(U_1) &= E[U_1 - E(U_1)]^2 = E(U_1 - \bar{U}_1)^2 \\ &= E(U_1') \end{aligned}$$

عناصر قطری ماتریس $E(uu')$ ، در واقع واریانس مقادیر U_1, U_2, \dots, U_n است. همچنین می‌دانیم کوواریانس U_i و U_j برابر است با:

$$\begin{aligned} \text{Cov}(U_i, U_j) &= E[U_i - E(U_i)][U_j - E(U_j)], \\ &= E(U_i - \bar{U})(U_j - \bar{U}), \\ &= E(U_i U_j). \end{aligned}$$

بنابراین کوواریانس U_i و U_j به ازای تمام مقادیر i و j استثنای مقادیر $i=j$ ، برابر عناصر غیرقطری ماتریس $E(uu')$ است. بدین ترتیب می‌توان نوشت،

$$\Sigma_u = E(uu') = \begin{bmatrix} \text{Var}(U_1) & \text{Cov}(U_1, U_2) & \dots & \text{Cov}(U_1, U_n) \\ \text{Cov}(U_2, U_1) & \text{Var}(U_2) & \dots & \text{Cov}(U_2, U_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(U_n, U_1) & \text{Cov}(U_n, U_2) & \dots & \text{Var}(U_n) \end{bmatrix}, \quad (5.8)$$

می‌دانیم یکی از فرضهای کلاسیک U_i ، فرض واریانس همسانی است. براساس این فرض، معادله ۵-۸ را می‌توان چنین نوشت،

$$\Sigma_u = E(uu') = \begin{bmatrix} \sigma^2 & \text{Cov}(U_1, U_2) & \dots & \text{Cov}(U_1, U_n) \\ \text{Cov}(U_2, U_1) & \sigma^2 & \dots & \text{Cov}(U_2, U_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(U_n, U_1) & \text{Cov}(U_n, U_2) & \dots & \sigma^2 \end{bmatrix}, \quad (5.9)$$

که در آن σ^2 برابر واریانس جمله اختلال است.

عدم خودهمبستگی، یعنی صفر بودن کوواریانس هر دو مقدار U_i و U_j ، یکی دیگر از فرضهای کلاسیک جمله‌های اختلال است؛ یعنی $\text{Cov}(U_i, U_j) = 0$. براساس

این فرض، تمام جمله‌های غیرقطری ماتریس Σ_u برابر صفر می‌شود؛ یعنی Σ_u یک ماتریس قطری خواهد بود. با استفاده از معادلهٔ ۵.۸ داریم:

$$\Sigma_u = E(uu') = \begin{bmatrix} \text{Var}(U_1) & \cdot & \dots & \cdot \\ \cdot & \text{Var}(U_2) & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \dots & \text{Var}(U_n) \end{bmatrix} \quad (5-10)$$

اگر هر دو فرض واریانس همسانی و عدم خود همبستگی را با هم در نظر بگیریم، در معادلهٔ ۵.۸، تمام عناصر قطری برابر σ^2 و عناصر غیرقطری برابر صفر می‌شود و ماتریس واریانس - کوواریانس Σ_u را می‌توان به صورت زیر نوشت،

$$\Sigma_u = E(uu') = \begin{bmatrix} \sigma^2 & \cdot & \dots & \cdot \\ \cdot & \sigma^2 & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \dots & \sigma^2 \end{bmatrix}$$

از σ^2 فاکتور گرفته خواهیم داشت:

$$\Sigma_u = E(uu') = \sigma^2 \begin{bmatrix} 1 & \cdot & \dots & \cdot \\ \cdot & 1 & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \dots & 1 \end{bmatrix}$$

می‌دانیم ماتریسی که هر یک از عناصر قطری آن یک و همهٔ عناصر غیرقطری آن صفر باشد، یک ماتریس «یکه»^۱ است که معمولاً آن را با I نشان می‌دهند. بدین ترتیب، ماتریس واریانس - کوواریانس جمله‌های اختلال مدل رگرسیون، با فرضهای واریانس

ناهمسانی و عدم خودهمبستگی برابر است با

$$\Sigma_u = E(uu') = \sigma^2 \mathbf{I}. \quad (5-11)$$

چهارمین فرض u به تابع توزیع احتمال آن مربوط است. فرض می شود که تابع توزیع احتمال u نرمال است. با توجه به میانگین 0 و واریانس $\sigma^2 \mathbf{I}$ ، می توان چنین نوشت،

$$u \sim N(0, \sigma^2 \mathbf{I}) \quad (5-12)$$

بعد از بیان فرضهای u ، باید به بررسی فرضهای ماتریس متغیرهای توضیحی، یعنی X بپردازیم. اولین فرض این است که X یک ماتریس غیر تصادفی است. این فرض بدین معنی است که اگر یک نمونه n تایی دیگر از مشاهدات داشته باشیم، ماتریس متغیرهای توضیحی، یعنی X بدون تغییر باقی می ماند. در ابتدا به نظر می رسد، قبول چنین فرضی بسیار دشوار باشد؛ زیرا انتظار این است که مقادیر متغیرهای توضیحی از یک نمونه به نمونه دیگر فرق کند. اما همان گونه که در فصل اول اشاره کردیم - منظور از فرض اخیر این است که متغیرهای توضیحی را در آزمایشهای تکراری ثابت نگه می داریم تا به این ترتیب بتوانیم تأثیر تغییرات U_1 بر Y_1 را مشاهده کنیم. ممکن است این نکته مطرح شود که در علوم اجتماعی، ثابت نگه داشتن متغیرهای توضیحی و تکرار آزمایش معمولاً امکانپذیر نیست. دقیقاً به همین علت است که باید از آزمایشهای فرضی تکراری صحبت کنیم. بنابراین فرض غیر تصادفی بودن X ، معمولاً در اقتصاد نمی تواند مابه ازای عینی داشته باشد و تنها به این سبب مطرح می شود که بتواند خصوصیات تأثیرگذاری متغیر تصادفی U_1 را بر تغییرات Y_1 بیان کند. در واقع درک استنتاجات آماری به شرط ثبات X ، می تواند بسیار روشن کننده باشد. همان گونه که در جلد دوم این کتاب خواهیم دید، می توان فرض غیر تصادفی بودن X را کنار گذاشت و تأثیر آن را در استنتاجات قبلی مشاهده کرد.

دومین فرض مربوط به ماتریس X این است که بگوییم متغیرهای توضیحی، همخطی کامل ندارد؛ یعنی در یک رابطه خطی کامل با یکدیگر قرار نمی گیرد. همان گونه که در قسمت ۲-۴ بیان شد - اگر بین متغیرهای توضیحی همبستگی کامل

خطی برقرار باشد، تخمین مستقل هریک از پارامترها، عملاً غیرممکن خواهد بود. بررسی مسأله همخطی، موضوع جلد دوم این کتاب خواهد بود. با استفاده از تعریف (رتبه ماتریس)^۱ می توان این فرض را به صورت زیر بیان کرد،

$$\rho(X) = k$$

که در آن $\rho(X)$ ، در واقع مرتبه ماتریس X و k تعداد متغیرهای توضیحی است.^۱ اکنون که با فرضهای کلاسیک مدلهای رگرسیون چندمتغیره آشنا شدیم، می توان به تخمین مدل به زبان ماتریسی پرداخت. اما بهتر است قبل از آن، به بیان ماتریسی جمله های پسماند نیز اشاره کنیم تا مبحث بیان ماتریسی رگرسیون چندمتغیره کاملتر شود.

۲. بیان ماتریسی جمله های پسماند

مدل رگرسیون چندمتغیره (۵-۱) را یک بار دیگر می نویسیم؛

$$Y_i = \beta_1 + \beta_2 X_{i1} + \dots + \beta_k X_{ik} + U_i.$$

اگر به جای پارامترهای مدل، یعنی $\beta_1, \beta_2, \dots, \beta_k$ ، مقادیر تخمین آنها را قرار دهیم، تخمین مدل رگرسیون به دست می آید؛ بنابراین:

$$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_{i1} + \dots + \hat{\beta}_k X_{ik}. \quad (5-12)$$

تفاوت بین مقدار \hat{Y}_i - که از معادله ۵-۱۳ حاصل می شود - با مقدار مشاهده شده Y_i برابر با مقدار پسماند است؛ بنابراین با توجه به تعریف

$$\text{پسماند} = \text{تخمین} - \text{مشاهده}$$

1. Rank of a Matrix

به پیوست «الف» مراجعه شود.

۲. می توان فرض سومی نیز برای متغیرهای توضیحی مطرح کرد: تعداد مشاهدات بیشتر از تعداد پارامترهایی است که باید تخمین زده شود؛ یعنی $n > k$.

خواهیم داشت:

$$e_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_k X_{ik}) \quad (5.14)$$

که در آن e_i پسماند و $i = 1, 2, \dots, n$.

برای بیان ماتریسی پسماندها، باید به جای مقادیر β در سیستم معادله‌های 5.3، مقادیر $\hat{\beta}$ را قرار دهیم. اگر بردار $\hat{\beta}$ را به صورت زیر تعریف کنیم،

$$\hat{\beta} = [\hat{\beta}_0 \hat{\beta}_1 \dots \hat{\beta}_k],$$

با توجه به معادله 5.4 خواهیم داشت:

$$\hat{y} = X\hat{\beta}, \quad (5.15)$$

که در آن $n \times 1 \rightarrow \hat{y}$. با استفاده از معادله 5.14، می‌توان بردار پسماند، یعنی e را چنین نوشت،

$$e = y - \hat{y}, \quad (5.16)$$

که در آن $n \times 1 \rightarrow \hat{e}$. معادله 5.16 را می‌توان به صورت زیر نوشت،

$$e = y - X\hat{\beta}. \quad (5.17)$$

5.3 تخمین مدل رگرسیون خطی چندمتغیره: روش حداقل مربعات معمولی
می‌دانیم معیار روش حداقل مربعات معمولی این است که باید پارامترها را چنان تخمین زد که مجموع مربعات پسماند حداقل شود؛ بنابراین ابتدا e_i^2 را به زبان ماتریسی می‌نویسیم و سپس از آن نسبت به پارامترها مشتق می‌گیریم. مقادیری از پارامترها که این مشتق را صفر کند، جواب مسأله است. ابتدا نشان می‌دهیم که $e'e$ ، بیان ماتریسی مجموع مربعات پسماند، یعنی $\sum e_i^2$ است. برای این منظور چنین می‌نویسیم،

$$e'e = [e_1 \ e_2 \ \dots \ e_n] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix},$$

$$= e_1^2 + e_2^2 + \dots + e_n^2 = \sum_{i=1}^n e_i^2.$$

با جایگزینی معادله ۵-۱۷ در $e'e$ خواهیم داشت:

$$\begin{aligned} e'e &= (y - X\hat{\beta})' (y - X\hat{\beta}), \\ &= (y' - \hat{\beta}' X') (y - X\hat{\beta}), \\ &= y'y - y' X\hat{\beta} - \hat{\beta}' X' y + \hat{\beta}' X' X\hat{\beta}. \quad (5-18) \end{aligned}$$

به ابعاد $y' X\hat{\beta}$ و $y' X\hat{\beta}$ توجه می‌کنیم.

$$y' X\hat{\beta} \rightarrow (1 \times n) (n \times k) (k \times 1) \rightarrow (1 \times 1),$$

$$\hat{\beta}' X' y \rightarrow (1 \times k) (k \times n) (n \times 1) \rightarrow (1 \times 1).$$

بنابراین $y' X\hat{\beta}$ و $\hat{\beta}' X' y$ هر دو اسکالر هستند و چون یکی ترانهاد دیگری است، پس با یکدیگر برابرند؛ یعنی:

$$y' X\hat{\beta} = \hat{\beta}' X' y.$$

با جایگزینی رابطه فوق در معادله ۵-۱۸، خواهیم داشت:

$$e'e = y'y - 2y' X\hat{\beta} + \hat{\beta}' X' X\hat{\beta}. \quad (5-19)$$

باید از $e'e$ در معادله فوق نسبت به $\hat{\beta}$ مشتق بگیریم. در پیوست (۵-الف) دیدیم که در

معادله $y = ax$ ، داریم:

$$\frac{\partial y}{\partial x} = a'$$

همچنین می‌دانیم که اگر داشته باشیم: $y = x'Ax$ آنگاه

$$\frac{\partial y}{\partial x} = 2Ax$$

در معادله ۵-۱۹ مجهول ما بردار $\hat{\beta}$ است. اگر $X'y$ را a و $X'X$ را برابر A بگیریم، خواهیم داشت:

$$\frac{\partial (e'e)}{\partial \hat{\beta}} = -2X'y + 2X'X\hat{\beta}$$

با توجه به اینکه k پارامتر داریم؛ k معادله از این رابطه به دست می‌آید. مقادیری از این k پارامتر که معادله فوق را صفر می‌کند، تخمینهای ما خواهد بود. بدین ترتیب

$$-2X'y + 2X'X\hat{\beta} = 0$$

$$(X'X)\hat{\beta} = X'y \quad (5-20)$$

معادله‌های فوق در واقع سیستم معادله‌های نرمال برای یک مدل رگرسیون با k متغیر توضیحی است. دو طرف معادله ۵-۲۰ را در $(X'X)^{-1}$ ضرب می‌کنیم، خواهیم داشت:

$$\hat{\beta} = (X'X)^{-1}X'y \quad (5-21)$$

بدیهی است برای اینکه معادله ۵-۲۱ برقرار باشد، باید $(X'X)$ معکوس داشته باشد. دومین شرط مربوط به متغیرهای توضیحی در مدل‌های رگرسیون چندمتغیره، یعنی معادله ۵-۱۲ دلالت بر استقلال کامل ستونهای ماتریس X از یکدیگر دارد؛ بنابراین با توجه به معادله ۵-۱۲ داریم:

$$\rho(X'X) = \rho(X) = k$$

و همین امر تضمینی است که درمینال $(X'X)$ صفر نباشد و معکوس داشته باشد.

اکنون که فرمول $\hat{\beta}$ را در مدل‌های رگرسیون چندمتغیره به دست آوردیم، بهتر است به ساختار ماتریسهای موجود در آن توجه بیشتری داشته باشیم تا چگونگی استفاده از این فرمول در حل مسائل روشن شود. ابتدا به مطالعه ماتریس $X'X$ می‌پردازیم. ماتریس X را در معادله ۵-۳ می‌نویسیم و آن را از سمت راست در ترانهاد خود ضرب می‌کنیم،

$$X'X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{11} & X_{12} & & X_{1n} \\ X_{21} & X_{22} & & X_{2n} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ X_{k1} & X_{k2} & & X_{kn} \end{bmatrix} \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1n} \\ 1 & X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{nn} \end{bmatrix},$$

و در نتیجه بعد از ضرب کردن این دو ماتریس داریم

$$X'X = \begin{bmatrix} n & \sum X_{1i} & \sum X_{2i} & \dots & \sum X_{ki} \\ \sum X_{1i} & \sum X_{1i}^2 & \sum X_{1i} X_{2i} & \dots & \sum X_{1i} X_{ki} \\ \sum X_{2i} & \sum X_{2i} X_{1i} & \sum X_{2i}^2 & \dots & \sum X_{2i} X_{ki} \\ \vdots & \vdots & \vdots & & \vdots \\ \sum X_{ki} & \sum X_{ki} X_{1i} & \sum X_{ki} X_{2i} & \dots & \sum X_{ki}^2 \end{bmatrix}.$$

(۵-۲۲)

ملاحظه می‌شود که با داشتن مقادیر X_i به راحتی می‌توان کمیت‌های لازم برای عناصر مختلف ماتریس $X'X$ را محاسبه و سپس این ماتریس را معکوس کرد. یادآوری می‌کنیم که $(X'X)^{-1}$ یک ماتریس $(k \times k)$ است.

به بررسی بردار $X'y$ می‌پردازیم، ابتدا می‌گوییم که

$$X'y \rightarrow (k \times n)(n \times 1) \rightarrow (k \times 1).$$

برای مطالعه عناصر $X'y$ ، باید ماتریس X' را در بردار y ضرب کنیم،

$$X'y = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{T1} & X_{T2} & & X_{Tn} \\ X_{r1} & X_{r2} & & X_{rn} \\ \vdots & \vdots & & \vdots \\ X_{k1} & X_{k2} & & X_{kn} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix}$$

بعد از ضرب کردن داریم،

$$X'y = \begin{bmatrix} \sum Y_t \\ \sum X_{Tt} Y_t \\ \sum X_{rt} Y_t \\ \vdots \\ \sum X_{kt} Y_t \end{bmatrix} \quad (5.23)$$

با داشتن مقادیر X_t, Y_t ، محاسبه عناصر بردار $X'y$ بسیار ساده است. با استفاده از ماتریس $(X'X)^{-1}$ و بردار $X'y$ ، محاسبه بردار $\hat{\beta}$ ، به راحتی با استفاده از فرمول 5.21 انجام خواهد شد.

دقیقاً مانند آنچه در رگرسیون ساده و در معادله‌های 1.23 و 1.24 دیدیم، می‌توان نشان داد که اولاً میانگین جمله‌های پسماند صفر است، و ثانیاً جمله‌های پسماند از متغیرهای توضیحی مستقل است. برای این منظور، ابتدا بردار e ، از معادله 5.17 را به صورت زیر می‌نویسیم:

$$y = e + X\hat{\beta}$$

و سپس آن را در معادله 5.20 قرار می‌دهیم

$$(X'X)\hat{\beta} = X'e + (X'X)\hat{\beta}$$

در نتیجه داریم:

$$X'e = 0, \quad (5.24)$$

که در آن $n \times 1 \rightarrow 0$ و شامل عناصر صفر است. معادله ۵-۲۴ را به صورت زیر بسط می‌دهیم:

$$X'e = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{r1} & X_{r2} & & X_{rn} \\ X_{k1} & X_{k2} & & X_{kn} \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} \sum e_i \\ \sum X_{ri} e_i \\ \sum X_{ki} e_i \\ \vdots \\ \sum X_{ki} e_i \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

بدین ترتیب ملاحظه می‌شود که $\sum e_i = 0$ یا $\bar{e} = 0$ و نیز $\sum X_{ri} e_i = 0$ ، $i = 1, 2, \dots, k$ است. یادآوری می‌کنیم که نتیجه تخمینهای حداقل مربعات معمولی همواره پسماندهایی است که میانگین آنها صفر است به شرط اینکه مدل رگرسیون دارای جمله ثابت باشد. قبل از پایان این مبحث دوباره به ساختار معادله‌های نرمال در رگرسیون چندمتغیره توجه می‌کنیم. معادله ۵-۲۰ را می‌نویسیم:

$$(X'X)\hat{\beta} = X'y.$$

مقادیر $(X'X)$ و $X'y$ را به ترتیب از معادله‌های ۵-۲۲ و ۵-۲۳ در معادله فوق قرار می‌دهیم. بعد از ضرب کردن داریم:

$$n\hat{\beta}_1 + \hat{\beta}_2 \sum X_{r1} + \hat{\beta}_3 \sum X_{r2} + \dots + \hat{\beta}_k \sum X_{rk} = \sum Y_i,$$

$$\hat{\beta}_1 \sum X_{r1} + \hat{\beta}_2 \sum X_{r1}^2 + \hat{\beta}_3 \sum X_{r1} X_{r2} + \dots + \hat{\beta}_k \sum X_{r1} X_{rk} = \sum X_{r1} Y_i,$$

$$\hat{\beta}_1 \sum X_{t1} + \hat{\beta}_2 \sum X_{t1} X_{t1} + \hat{\beta}_3 \sum X_{t1}^2 + \dots + \hat{\beta}_k \sum X_{t1} X_{t1} = \sum X_{t1} Y_t,$$

⋮

$$\hat{\beta}_1 \sum X_{tk} + \hat{\beta}_2 \sum X_{tk} X_{t1} + \hat{\beta}_3 \sum X_{tk} X_{t1} + \dots + \hat{\beta}_k \sum X_{tk}^2 = \sum X_{tk} Y_t.$$

بدیهی است از حل دستگاه k معادله k مجهول فوق مقادیر مجهولات $(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)$ به دست می‌آید. اگر معادله اول نرمال را در نظر گرفته و دو طرف آن را بر n تقسیم کنیم، خواهیم داشت:

$$\bar{Y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{X}_1 + \hat{\beta}_3 \bar{X}_1 + \dots + \hat{\beta}_k \bar{X}_k, \quad (5.25)$$

که دقیقاً تعمیم نتیجه‌ای است که در معادله ۱-۲۶ برای رگرسیون ساده به دست آوردیم. مثال ۵.۱ مدل رگرسیون زیر مفروض است:

$$Y_t = \beta_1 + \beta_2 X_{t1} + \beta_3 X_{t1} + U_t.$$

برای تخمین پارامترهای β_1, β_2 و β_3 ، مشاهدات Y_t, X_{t1} و X_{t1} را در جدول ۵.۱ نوشته‌ایم.

جدول ۵.۱

Y_t	۳	۱	۸	۳	۵
X_{t1}	۳	۱	۵	۲	۴
X_{t1}	۵	۴	۶	۴	۶

به کمک روش حداقل مربعات معمولی مقادیر $\hat{\beta}_1, \hat{\beta}_2$ و $\hat{\beta}_3$ را به دست آورید. با استفاده از فرمول ۵.۲۱ می‌دانیم:

$$\begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix} = \begin{bmatrix} n & \sum X_{t1} & \sum X_{t1} \\ \sum X_{t1} & \sum X_{t1}^2 & \sum X_{t1} X_{t1} \\ \sum X_{t1} & \sum X_{t1} X_{t1} & \sum X_{t1}^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum Y_t \\ \sum X_{t1} Y_t \\ \sum X_{t1} Y_t \end{bmatrix},$$

$$\hat{\beta} = (X'X)^{-1} X'y$$

چندان دشوار نیست و محاسبه آن در واقع یادآوری جبر ماتریسی نیز هست. در مورد این گونه ماتریسهای (3×3) اگر فرمول معکوس کردن آنها نیز فراموش شود، همواره می توان از راه معادله های نرمال استفاده کرد و با روش حذفی به راحتی یک دستگاه سه معادله سه مجهولی را حل کرد. این روشی است که در این قسمت از آن استفاده خواهیم کرد. اگر دو طرف رابطه فوق را در مقدار عددی $(X'X)$ ضرب کنیم خواهیم داشت

$$\begin{bmatrix} 0 & 10 & 20 \\ 10 & 50 & 81 \\ 20 & 81 & 129 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix} = \begin{bmatrix} 20 \\ 76 \\ 109 \end{bmatrix} ,$$

که در واقع همان معادله های نرمال ۲۰-۵ برای این مثال است.

کافی است سه برابر سطر اول ماتریس $(X'X)$ را از سطر دوم و پنج برابر سطر اول را از سطر سوم کم کنیم. سیستم جدید معادله های نرمال به صورت زیر خواهد بود،

$$\begin{bmatrix} 0 & 10 & 20 \\ 0 & 10 & 6 \\ 0 & 6 & 4 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix} = \begin{bmatrix} 20 \\ 16 \\ 9 \end{bmatrix} .$$

اکنون باید $0/6$ سطر دوم را از سطر سوم کم کنیم. در نتیجه

$$\begin{bmatrix} 0 & 10 & 20 \\ 0 & 10 & 6 \\ 0 & 0 & 0/4 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix} = \begin{bmatrix} 20 \\ 16 \\ -0/6 \end{bmatrix} .$$

معادله سوم سیستم معادله های نرمال فوق عبارت است از

$$0/4 \hat{\beta}_3 = -0/6 ,$$

$$\hat{\beta}_3 = -1/5 .$$

معادله دوم در سیستم معادله‌های نرمال فوق را می‌نویسیم،

$$۱۰\hat{\beta}_2 + \hat{\beta}_3 = ۱۶.$$

در نتیجه:

$$۱۰\hat{\beta}_2 + ۶(-۱/۵) = ۱۶,$$

$$\hat{\beta}_2 = ۲/۵.$$

از معادله اول نرمال هم داریم:

$$۵\hat{\beta}_1 + ۱۵\hat{\beta}_2 + ۲۵\hat{\beta}_3 = ۲۰,$$

$$۵\hat{\beta}_1 + ۱۵(۲/۵) + ۲۵(-۱/۵) = ۲۰,$$

$$\hat{\beta}_1 = ۴.$$

بدین ترتیب تخمین مدل رگرسیون مفروض عبارت است از

$$\hat{Y}_i = ۴ + ۲/۵X_{2i} - X_{3i}.$$

در پایان این مثال، به یک قرارداد در نمایش محاسبات مدل‌های رگرسیون چندمتغیره اشاره می‌کنیم. در مثال فوق دیدیم که برای محاسبه $\hat{\beta}_1$ ، $\hat{\beta}_2$ و $\hat{\beta}_3$ کمیت‌هایی که باید محاسبه شود عبارت است از:

الف) محاسباتی که باید برای $(X'X)$ انجام شود

$$\sum X_{2i} = ۱۵, \quad \sum X_{3i} = ۲۵, \quad \sum X_{2i}X_{3i} = ۸۱,$$

$$\sum X_{2i}^2 = ۵۵, \quad \sum X_{3i}^2 = ۱۲۹, \quad \sum X_{2i}X_{3i} = ۸۱.$$

یادآوری می‌کنیم که ماتریس $(X'X)$ یک ماتریس قرینه است و مقدار $\sum X_{2i}X_{3i}$ برابر با $\sum X_{3i}X_{2i}$ است.

ب) محاسباتی که برای $(X'Y)$ لازم است عبارتند از

$$\sum Y_i = ۲۰, \quad \sum X_{2i}Y_i = ۷۶, \quad \sum X_{3i}Y_i = ۱۰۹.$$

معمولاً در مسائل آموزشی، به جای مشاهدات خام متغیرهای Y_i ، X_{1i} و X_{2i} ، مقادیر محاسبه شده فوق را می آورند تا به طور مستقیم ماتریسهای $(X'X)$ و $X'y$ محاسبه شود. روش زیر برای نمایش مقادیر مجموع مربعات $\sum X_{1i}^2$ و مجموع حاصلضربهای $\sum X_{1i}X_{2i}$ ، به صورت قراردادی پذیرفته شده است.

	Y_i	X_{1i}	X_{2i}
Y_i	□	۷۶	۱۰۹
Y_{1i}	۷۶	۵۵	۸۱
Y_{2i}	۱۰۹	۸۱	۱۲۹

یادآوری می کنیم که در محاسبات این مثال، مقدار $\sum Y_i^2$ محاسبه نشده و مقدار آن در جدول فوق خالی مانده و با علامت □ مشخص شده است. قرینگی ماتریس $X'X$ نیز از جدول فوق به سهولت ملاحظه می شود. همچنین یادآوری می شود که جدول فوق فقط برای نشان دادن مقادیر مجموع مربعات و مجموع حاصلضربهاست؛ بنابراین $\sum X_{1i}$ یا $\sum Y_i$ در آن منعکس نیست. روش فوق برای حالتی که k متغیر توضیحی داریم نیز به همین صورت به کار می رود.

تخمین مدل رگرسیون خطی چندمتغیره برحسب انحراف از میانگین

در بسیاری موارد، مشاهدات متغیرهای برونزا و درونزا را برحسب انحراف از میانگین می نویسیم تا به این ترتیب اولاً محاسبات $(X'X)^{-1}$ و $X'y$ ساده تر انجام شود و ثانیاً همان گونه که بزودی خواهیم دید از ابعاد ماتریس $(X'X)$ یکی کم شده، معکوس کردن آن بسیار ساده تر خواهد شد؛ بنابراین ضروری است که تخمین مدل های رگرسیون چندمتغیره را در این حالت نیز به طور خلاصه مطالعه کنیم.

مدل ۵-۱ را یک بار دیگر می نویسیم؛

$$Y_i = \beta_1 + \beta_2 X_{1i} + \dots + \beta_k X_{ki} + U_i.$$

دو طرف را برای تمام مقادیر $i = 1, 2, \dots, n$ جمع کرده و بر n تقسیم می کنیم.

جمله‌های اختلال، یعنی Σ_{it} که در معادله‌های ۵-۹، ۵-۱۰ و ۵-۱۱ تعریف شده است - به همان ترتیب برای این حالت نیز تعریف خواهد شد. $e'e$ نیز به همان ترتیب معادله ۵-۱۶ است؛ با وجود این، ماتریس $(X'X)$ و $X'y$ تغییراتی خواهد کرد که بررسی خواهد شد. با استفاده از معادله ۵-۲۸، ماتریس $(X'X)$ را می‌نویسیم،

$$X'X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{k1} & x_{k2} & \dots & x_{kn} \end{bmatrix} \begin{bmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & & \vdots \\ x_{1n} & x_{2n} & & x_{nn} \end{bmatrix},$$

X' X

و بعد از ضرب کردن این دو ماتریس، داریم:

$$X'X = \begin{bmatrix} \sum x_{1t}^2 & \sum x_{1t} x_{2t} & \dots & \sum x_{1t} x_{kt} \\ \sum x_{2t} x_{1t} & \sum x_{2t}^2 & \dots & \sum x_{2t} x_{kt} \\ \vdots & \vdots & & \vdots \\ \sum x_{kt} x_{1t} & \sum x_{kt} x_{2t} & \dots & \sum x_{kt}^2 \end{bmatrix}, \quad (5.30)$$

بنابراین ماتریس $X'X$ ، در حالتی که مشاهدات برحسب انحراف از میانگین باشد، بسیار ساده است. عناصر قطری آن مجموع مربعات x_{it} و عناصر غیرقطری آن مجموع حاصلضربها، یعنی $\sum x_{it} x_{jt}$ ، است. ملاحظه می‌شود که $X'X$ یک ماتریس مربع و $(k-1) \times (k-1)$ است. در حالی که در حالت قبلی، یعنی وقتی مشاهدات برحسب انحراف از میانگین نبود، این ماتریس $(k \times k)$ است. نتیجه می‌گیریم که اگر مشاهدات را برحسب انحراف از میانگین بنویسیم از ابعاد ماتریس $X'X$ یکی کم می‌شود و معکوس کردن آن ساده‌تر خواهد بود. البته در عمل، از این امتیاز فقط در حالتی می‌توان استفاده کرد که بخواهیم سه پارامتر یا چهار پارامتر تخمین بزنیم؛ زیرا با این تبدیل، ابعاد ماتریس $X'X$ به جای اینکه (3×3) یا (4×4) بشود، (2×2) یا (3×3) خواهد بود که

معکوس کردن آن بسیار ساده تر است. اما برای مدل‌های رگرسیون که بیش از چهار پارامتر دارد، باید عملاً از کامپیوتر استفاده شود؛ بنابراین فرق چندانی نمی‌کند که مشاهدات را بر حسب انحراف از میانگین بنویسیم.

برای بررسی ماتریس $X'y$ به ترتیب زیر عمل می‌کنیم،

$$X'y = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{k1} & x_{k2} & & x_{kn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \sum x_{1i} y_i \\ \sum x_{2i} y_i \\ \vdots \\ \sum x_{ki} y_i \end{bmatrix} \quad (5.31)$$

$X' \qquad y$

ماتریس $X'y$ در این حالت، $1 \times (k-1)$ و عناصر آن $\sum x_{ii} y_i$ است. بردار جمله‌های پسماند نیز همانند معادله ۵-۱۶ تعریف می‌شود

$$e = y - \hat{y},$$

یا

$$e = y - X\hat{\beta}, \quad (5.32)$$

که در آن، $n \times 1 \rightarrow e$ و $1 \times (k-1) \rightarrow X$ و $1 \times (k-1) \rightarrow \hat{\beta}$. مجموع مربعات پسماند نیز دقیقاً همان رابطه ۵-۱۹ است، با این تفاوت که از تعداد ستونها و ردیفهای ماتریس X و بردار $\hat{\beta}$ به ترتیب یکی کم شده است. با روشی کاملاً مشابه با حالت قبلی می‌توان از $e'e$ نسبت به $\hat{\beta}$ مشتق گرفت و به معادله‌های نرمال به شرح زیر رسید،

$$(X'X)\hat{\beta} = X'y, \quad (5.33)$$

که در آن $(k-1) \times (k-1) \rightarrow X'X$ و $1 \times (k-1) \rightarrow X'y$ و $1 \times (k-1) \rightarrow \hat{\beta}$. معادله ۵-۷۸ بیان دقیقتری از این معادله است. فرمول $\hat{\beta}$ نیز عبارت خواهد بود از

$$\hat{\beta} = (X'X)^{-1} X'y, \quad (5.34)$$

که در آن $1 \times (k-1) \rightarrow \hat{\beta}$ ، $(k-1) \times (k-1) \rightarrow (X'X)^{-1}$ و $1 \times (k-1) \rightarrow X'y$. معادله ۵-۷۹ صورت دقیقتری از این فرمول است. بدین ترتیب تخمین مدل رگرسیون ۵-۲۹ عبارت است از

$$\hat{y} = X\hat{\beta}, \quad (5.35)$$

مثال ۵-۲ مدل رگرسیون زیر مفروض است،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

با استفاده از مشاهدات مندرج در جدول ۵-۲، پارامترهای α ، β و γ را با حداقل مربعات معمولی و با روش انحراف از میانگین تخمین بزنید.

جدول ۵-۲

Y_i	۳	۱	۸	۳	۵
X_i	۳	۱	۵	۲	۴
Z_i	۵	۴	۶	۴	۶

این مثال دقیقاً همان مثال ۵-۱ است، با این تفاوت که اولاً حروف گذاری متغیرهای توضیحی و نیز پارامترها را تغییر داده ایم و ثانیاً خواسته شده است که ابتدا مشاهدات را برحسب انحراف از میانگین بنویسیم تا کاربرد فرمول ۵-۳۴ روشن شود. معمولاً راه حل ساده این گونه مسائل این است که حروف متغیرها و پارامترها را چنان تغییر دهیم که با فرمولهای متعارف منطبق شود. اما در این مثال و به عنوان یک تمرین، با همین حروف جدید کاری کنیم تا مفهوم ماتریس $X'X$ و بردارهای $\hat{\beta}$ و $X'y$ موجود در فرمول ۵-۳۴ بیشتر روشن شود.

ابتدا مدل را برحسب انحراف از میانگین می نویسیم:

$$y_i = \beta x_i + \gamma z_i + u_i.$$

با استفاده از معادله ۵-۳۰ ماتریس $X'X$ را برای این مدل می نویسیم،

$$X'X = \begin{bmatrix} \sum x_i^2 & \sum x_i z_i \\ \sum z_i x_i & \sum z_i^2 \end{bmatrix}.$$

به کمک معادله ۵-۳۱ می توان بردار $X'y$ را برای این مدل به صورت زیر نوشت:

$$X'y = \begin{bmatrix} \sum x_i y_i \\ \sum z_i y_i \end{bmatrix}$$

و سرانجام بردار $\hat{\beta}$ برای این مثال، به این ترتیب است:

$$\hat{\beta} = \begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix}$$

برای محاسبه کمیت‌های لازم برای تشکیل ماتریس $X'X$ و بردار $X'y$ ، مشاهدات را برحسب انحراف از میانگین می نویسیم. کمیت‌های لازم در جدول ۵-۳ محاسبه شده است.

جدول ۵-۳

y_i	x_i	z_i	y_i	x_i	z_i	x_i^2	z_i^2	$x_i y_i$	$z_i y_i$	$x_i z_i$
۳	۲	۰	-۱	۰	۰	۰	۰	۰	۰	۰
۱	۱	۴	-۳	-۲	-۱	۴	۱	۶	۳	۲
۸	۵	۶	۴	۲	۱	۴	۱	۸	۴	۲
۲	۲	۴	-۱	-۱	-۱	۱	۱	۱	۱	۱
۵	۴	۶	۱	۱	۱	۱	۱	۱	۱	۱
Σ			۰	۰	۰	۱۰	۴	۱۶	۹	۶

$$\bar{y} = 4, \quad \bar{x} = 3, \quad \bar{z} = 0$$

ماتریس $X'X$ و $X'y$ را تشکیل می دهیم،

$$X'X = \begin{bmatrix} \sum x_i^2 & \sum x_i z_i \\ \sum z_i x_i & \sum z_i^2 \end{bmatrix} = \begin{bmatrix} 10 & 6 \\ 6 & 4 \end{bmatrix}$$

$$X'y = \begin{bmatrix} \sum x_i y_i \\ \sum z_i y_i \end{bmatrix} = \begin{bmatrix} 16 \\ 9 \end{bmatrix}$$

فرمول ۵-۳۴ را برای این مثال می نویسیم،

$$\begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix} = \begin{bmatrix} \sum x_1^2 & \sum x_1 z_1 \\ \sum z_1 x_1 & \sum z_1^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum x_1 y_1 \\ \sum x_1 z_1 \end{bmatrix} \\ = \begin{bmatrix} 10 & 6 \\ 6 & 4 \end{bmatrix}^{-1} \begin{bmatrix} 16 \\ 9 \end{bmatrix} = \begin{bmatrix} 1 & -\frac{3}{4} \\ -\frac{3}{4} & \frac{5}{4} \end{bmatrix}$$

و بعد از ضرب کردن، خواهیم داشت:

$$\begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix} = \begin{bmatrix} 2/5 \\ -1/5 \end{bmatrix}$$

یعنی $\hat{\beta} = 2/5$ و $\hat{\gamma} = -1/5$. برای محاسبه $\hat{\alpha}$ باید از معادله ۵-۲۵ استفاده کنیم. ابتدا این معادله را برای حالت خاص این مثال می نویسیم،

$$\bar{Y} = \hat{\alpha} + \hat{\beta}\bar{X} + \hat{\gamma}\bar{Z}$$

و در نتیجه خواهیم داشت:

$$4 = \hat{\alpha} + 2/5(3) - 1/5(5)$$

$$\hat{\alpha} = 4$$

بدین ترتیب، تخمین مدل مفروض به صورت زیر است:

$$\hat{Y}_i = 4 + 2/5 X_i - 1/5 Z_i$$

در پایان این قسمت به یک نکته اشاره می کنیم. فرمولهای ۵-۲۱ و ۵-۳۴ را یک بار دیگر ملاحظه کنید. هر دو فرمول به صورت زیر است،

$$\hat{\beta} = (X'X)^{-1} X'y$$

در حالی که فرمول ۵-۳۴ برای حالتی است که مشاهدات بر حسب انحراف از میانگین

است. در متون اقتصادسنجی کوشش بر این است که تنوع استفاده از حروف یا علامتگذاری به حداقل برسد. با اینکه به راحتی می‌توان از حروف دیگر برای فرمول ۵.۳۴ استفاده کرد، اما معمولاً آن را به همین صورت می‌نویسند؛ با وجود این، به این نکته توجه می‌کنند که ابعاد ماتریس $(X'X)$ و بردارهای $X'y$ و $\hat{\beta}$ با فرمول ۵.۲۱ متفاوت است و عناصر ماتریس و بردارهای فوق برحسب انحراف از میانگین محاسبه شده است.

۵.۴ ضریب تعیین

قبلاً از معادله ۱.۲۹ دیدیم که ضریب تعیین در رگرسیون ساده بنا به تعریف برابر است با نسبت تغییرات توضیح داده شده به کل تغییرات Y_t ،

$$r^2 = \frac{\sum \hat{y}_t^2}{\sum y_t^2} \quad (\text{ضریب تعیین})$$

در مبحث رگرسیون چندمتغیره و در معادله ۴-۱۶ نیز دیدیم که همین تعریف صادق است. با وجود این، آن را با R^2 نشان می‌دهیم تا از ضریب تعیین در رگرسیون ساده، یعنی r^2 ، قابل تمیز باشد. در این قسمت ابتدا به بررسی ضریب تعیین به زبان ماتریسی می‌پردازیم و سپس مفهوم «ضریب تعیین تعدیل شده»^۱ را بررسی می‌کنیم.

بیان ماتریسی R^2 مستلزم تعریف مجموع تغییرات توضیح داده شده، یعنی $\sum \hat{y}_t^2$ و مجموع کل تغییرات Y_t ، یعنی $\sum y_t^2$ به زبان ماتریسی است. می‌دانیم

$$\sum \hat{y}_t^2 = \sum (\hat{Y}_t - \bar{\hat{Y}})^2$$

که با توجه به معادله ۱-۳۰، یعنی $\bar{\hat{Y}} = \bar{Y}$ ، خواهیم داشت

$$\sum \hat{y}_t^2 = \sum (\hat{Y}_t - \bar{Y})^2.$$

اگر $n, \dots, 2, 1, t$ ، می‌توان نشان داد که

$$\sum \hat{y}_t^2 = \hat{y}' \hat{y}, \quad (5.36)$$

که در آن $n \times 1 \rightarrow \hat{y}$. برای اثبات، دو بردار \hat{y} و y را در هم ضرب کرده و نشان می‌دهیم که برابر $\sum \hat{y}_i^2$ است. ملاحظه می‌شود که

$$\hat{y}' \hat{y} = [\hat{y}_1 \hat{y}_2 \dots \hat{y}_n] \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \hat{y}_1^2 + \hat{y}_2^2 + \dots + \hat{y}_n^2 = \sum_1^n \hat{y}_i^2.$$

به همین ترتیب نشان می‌دهیم که مجموع کل تغییرات Y_1 برابر است با

$$\sum y_i^2 = y' y. \quad (5.37)$$

که در آن، $n \times 1 \rightarrow y$. برای اثبات کافی است که در بردار y و y را در هم ضرب کنیم،

$$y' y = [y_1 y_2 \dots y_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = y_1^2 + y_2^2 + \dots + y_n^2 = \sum_1^n y_i^2.$$

بدین ترتیب طبق معادله ۱-۲۹ می‌توان ضریب تعیین برای رگرسیون چندمتغیره را به صورت زیر تعریف کرد،

$$R^2 = \frac{\hat{y}' \hat{y}}{y' y}. \quad (5.38)$$

معادله ۵-۳۸ از نظر محاسباتی تا حدی مشکل است زیرا مستلزم محاسبه مقادیر \hat{y} برای تمام مقادیر t است. برای به دست آوردن یک فرمول ساده‌تر برای R^2 ، ابتدا معادله ۵-۳۵، یعنی $\hat{y} = X\hat{\beta}$ را می‌نویسیم و سپس آن را در ترانهاد خود از سمت راست ضرب می‌کنیم، خواهیم داشت:

$$\hat{y}' \hat{y} = (\hat{\beta}' X') (X\hat{\beta}).$$

مقدار $\hat{\beta}$ از معادله ۵-۳۴ را در معادله فوق قرار می‌دهیم،

$$\hat{y} = \hat{\beta}' X' X [(X' X)^{-1} X' y].$$

مقدار $(X' X)^{-1}$ داخل کروشه با مقدار $(X' X)$ خارج از کروشه حذف شده است و در نتیجه

$$\hat{y} = \hat{\beta}' X' y. \quad (5.39)$$

معادله ۵-۳۹ را در معادله ۵-۳۸ جایگزین می‌کنیم. داریم

$$R^1 = \frac{\hat{\beta}' X' y}{y' y} \quad (5.40)$$

معادله فوق، مشابه معادله ۱-۳۵ برای رگرسیون ساده و بیان ماتریسی معادله ۴-۱۹ است. از این معادله می‌توان در حل مسائل استفاده فراوانی کرد؛ زیرا وقتی یک مدل رگرسیون را تخمین زده‌ایم، مقادیر $\hat{\beta}$ و $X' y$ را داریم، بنابراین $y' y$ را که همان $\sum y_i^2$ است محاسبه می‌کنیم تا R^1 به سهولت و از طریق معادله ۵-۴۰ به دست آید.

محاسبه R^1 با فرمول فوق مستلزم داشتن مقادیر تخمین پارامترها است. می‌توان معادله ۵-۴۰ را تغییر داد به نحوی که R^1 فقط برحسب مشاهدات متغیرهای درون‌زا و برون‌زا نوشته شود. برای این منظور، باید مقدار $\hat{\beta}$ را در معادله ۵-۴۰ قرار دهیم. ابتدا معادله ۵-۳۴ را یک بار دیگر می‌نویسیم:

$$\hat{\beta} = (X' X)^{-1} X' y.$$

معادله فوق را ترانهاد می‌کنیم،

$$\hat{\beta}' = y' X (X' X)^{-1}. \quad (5.41)$$

یادآوری می‌شود که چون $(X' X)$ ، متقارن است، با ترانهاد خود برابر است؛ یعنی

$$(X' X)' = (X' X).$$

معادله ۵-۴۱ را در ۵-۴۰ قرار می‌دهیم،

$$R^1 = \frac{y' X (X' X)^{-1} X' y}{y' y}. \quad (5.42)$$

معادله ۵-۴۲ فرمول مستقیم محاسبه R^1 است.

در مدل رگرسیون ساده، قضیه مهمی را در قالب معادله ۱-۴۰ ثابت کردیم. اثبات این قضیه به زبان ماتریسی به قرار زیر است. می‌خواهیم ثابت کنیم که در یک مدل رگرسیون چندمتغیره، کل تغییرات متغیر درون‌زا برابر است با مجموع تغییرات توضیح داده شده، بعلاوه تغییرات توضیح داده نشده؛ به عبارت دیگر، می‌دانیم در رگرسیون ساده داریم

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2.$$

باید ابتدا مشابه معادله فوق را برای مدل‌های رگرسیون چندمتغیره بنویسیم. با استفاده از معادله‌های ۵-۳۶ و ۵-۳۷ می‌دانیم که

$$\sum y_i^2 = \mathbf{y}'\mathbf{y} \quad \text{و} \quad \sum \hat{y}_i^2 = \hat{\mathbf{y}}'\hat{\mathbf{y}}$$

همچنین می‌توان نشان داد که در یک مدل رگرسیون که متغیرها برحسب انحراف از میانگین اندازه‌گیری شده است، یعنی مدل ۵-۲۶، تغییرات توضیح داده نشده، یعنی مجموع مربعات پسماند، برابر است با

$$\sum e_i^2 = \mathbf{e}'\mathbf{e}. \quad (۵-۴۳)$$

برای اثبات کافی است به این نکته توجه کنیم که $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$ ؛ پس بودن جمله ثابت، یعنی β_1 ، در مدل رگرسیون ۵-۲۶ در ابعاد بردار \mathbf{y} و $\hat{\mathbf{y}}$ تأثیری ندارد، بلکه فقط مقدار $\hat{\mathbf{y}}$ را متأثر می‌کند؛ زیرا $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$. بنابراین:

$$\mathbf{e}'\mathbf{e} = [e_1 \ e_2 \ \dots \ e_n] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = e_1^2 + e_2^2 + \dots + e_n^2 = \sum e_i^2. \quad (۵-۴۴)$$

در نتیجه بیان ماتریسی قضیه‌ای که می‌خواهیم ثابت کنیم به شرح زیر خواهد بود،

$$\mathbf{y}'\mathbf{y} = \hat{\mathbf{y}}'\hat{\mathbf{y}} + \mathbf{e}'\mathbf{e} \quad (۵-۴۵)$$

یادآوری می‌شود که هر یک از جمله‌های معادله ۵-۴۵ یک اسکالر است.

برای اثبات، ابتدا تغییرات توضیح داده نشده، یعنی مجموع مربعات پسماند را می‌نویسیم،

$$e'e = (y - \hat{y})' (y - \hat{y}).$$

با جایگزینی معادله ۵-۳۵، یعنی $\hat{y} = X\hat{\beta}$ ، در معادله فوق داریم

$$e'e = (y - X\hat{\beta})' (y - X\hat{\beta}),$$

$$= (y' - \hat{\beta}' X') (y - X\hat{\beta}).$$

مقدار $\hat{\beta}$ از معادله ۵-۳۴ را در این معادله قرار می‌دهیم:

$$e'e = [y' - y'X(X'X)^{-1}X'] [y - X(X'X)^{-1}X'y].$$

و بعد از ضرب کردن داریم:

$$e'e = y'y - y'X(X'X)^{-1}X'y - y'X(X'X)^{-1}X'y + y'X(X'X)^{-1}X'y.$$

معادله فوق را ساده می‌کنیم،

$$e'e = y'y - y'X(X'X)^{-1}X'y. \quad (5-46)$$

حال نشان می‌دهیم که جمله آخر با $\hat{y}'\hat{y}$ برابر است. برای این منظور با استفاده از معادله ۵-۳۵، یعنی $\hat{y} = X\hat{\beta}$ ، چنین می‌نویسیم،

$$\hat{y}'\hat{y} = \hat{\beta}'X'X\hat{\beta} = [y'X(X'X)^{-1}] (X'X) [(X'X)^{-1}X'y],$$

در نتیجه

$$\hat{y}'\hat{y} = y'X(X'X)^{-1}X'y. \quad (5-47)$$

البته قبلاً نیز به نتیجه فوق رسیده بودیم؛ زیرا معادله فوق در واقع صورت کسر ۵-۴۲ است. به هر حال با جایگزینی معادله ۵-۴۷ در ۵-۴۶ خواهیم داشت

$$y'y = \hat{y}'\hat{y} + e'e,$$

که اثبات معادله ۵-۴۵ است.

با استفاده از معادله فوق می‌توان سه نکته دیگر را نیز بررسی کرد. اولاً، قلمرو تغییرات R^1 ، ثانیاً، فرمول دیگری برای R^1 ، و ثالثاً به دست آوردن مجموع مربعات پسماند $(e'e)$. برای به دست آوردن قلمرو تغییرات R^1 باید دو طرف معادله ۵-۴۵ را بر کل تغییرات متغیر درون‌زا، یعنی $y'y$ تقسیم کنیم:

$$\frac{\hat{y}'\hat{y}}{y'y} = \frac{y'y}{y'y} - \frac{e'e}{y'y}$$

با توجه به تعریف R^2 در رابطه ۵-۳۸ معادله فوق را به صورت زیر می‌نویسیم؛

$$R^2 = 1 - \frac{e'e}{y'y} \quad (5.48)$$

این معادله، نه تنها فرمول دیگری برای R^2 محسوب می‌شود، بلکه می‌تواند قلمرو تغییرات R^2 را نیز تعیین کند. بهترین حالت موقعی است که پسماندی نداشته باشیم؛ یعنی تخمین ما توانسته باشد تمام تغییرات y_t در مدل ۵-۲۶ را توضیح دهد. در چنین حالتی مجموع مربعات پسماند برابر صفر می‌شود و $e'e = 0$ ، در نتیجه $R^2 = 1$. بدترین حالت موقعی است که تخمین ما هیچ قدرت توضیحی نداشته باشد، یعنی $\hat{y}'\hat{y} = 0$ و در نتیجه $e'e = y'y$ باشد، در این صورت $R^2 = 0$ است؛ بنابراین

$$0 \leq R^2 \leq 1. \quad (5.49)$$

استفاده دیگری که معمولاً از معادله ۵-۴۸ می‌شود، این است که با داشتن R^2 ، مثلاً از معادله ۵-۴۲ و نیز با محاسبه $y'y$ ، می‌توان $e'e$ را محاسبه کرد:

$$e'e = (1 - R^2)y'y. \quad (5.50)$$

همان‌گونه که خواهیم دید، محاسبه مجموع مربعات پسماند، یعنی $e'e$ ، شرط لازم برای تخمین واریانس جمله اختلال است.

قبلاً در فصل اول و در مسئله ۱-۱۹ دیدیم که می‌توان r^2 را از فرمول زیر به دست آورد،

$$r^2 = \frac{\sum (y_i \hat{y}_i)^2}{\sum y_i^2 \sum \hat{y}_i^2}$$

یعنی r^2 با مجذور ضریب همبستگی بین Y_i و \hat{Y}_i برابر است. در اینجا می‌خواهیم ثابت کنیم که برای یک رگرسیون چندمتغیره نیز این خصوصیت صادق است؛ یعنی R^2 برابر است با مجذور ضریب همبستگی بین y و \hat{y} .

برای اثبات، از تعریف ضریب همبستگی بین y و \hat{y} شروع کرده، نشان خواهیم داد که این ضریب برابر R^2 است. می‌دانیم

$$\text{ضریب همبستگی بین } y \text{ و } \hat{y} = \frac{y' \hat{y}}{\sqrt{y' y} \sqrt{\hat{y}' \hat{y}}}$$

دو طرف را مجذور می‌کنیم،

$$\text{مجدور ضریب همبستگی بین } y \text{ و } \hat{y} = \frac{(y' \hat{y})^2}{(y' y) (\hat{y}' \hat{y})} \quad (0.01)$$

با توجه به $y = \hat{y} + e$ ، می‌توان صورت کسر را به صورت زیر نوشت

$$y' \hat{y} = (\hat{y}' + e') \hat{y} = \hat{y}' \hat{y} + e' \hat{y}.$$

با توجه به $\hat{y} = X\hat{\beta}$ داریم

$$y' \hat{y} = \hat{y}' \hat{y} + e' X \hat{\beta}.$$

با توجه به معادله ۰.۲۴ می‌دانیم $e' X = 0$ ، بنابراین با جایگزینی در معادله فوق خواهیم داشت

$$y' \hat{y} = \hat{y}' \hat{y}. \quad (0.02)$$

معادله ۰.۰۲ را در معادله ۰.۰۱ قرار می‌دهیم،

$$\text{مجدور ضریب همبستگی بین } y \text{ و } \hat{y} = \frac{(\hat{y}' \hat{y})^2}{(y' y) (\hat{y}' \hat{y})} = \frac{\hat{y}' \hat{y}}{y' y}$$

اما با توجه به معادله ۵-۳۸ می‌دانیم که

$$\frac{\hat{y}}{\hat{y}} = R^2,$$

بنابراین خواهیم داشت

$$R^2 = \text{مجدور ضریب همبستگی بین } \hat{y} \text{ و } y. \quad (5.53)$$

R^2 در مدل‌های رگرسیون بدون جمله ثابت

در فصل اول و مسأله ۱-۲۱ دیدیم که اگر بخواهیم یک رگرسیون ساده و بدون جمله ثابت، یعنی

$$Y_i = \beta X_i + U_i$$

را با فرمول معمولی ضریب تعیین، یعنی

$$r^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2},$$

تخمین بزنیم، ممکن است r^2 منفی شود. در این قسمت می‌خواهیم این نکته را برای حالت عمومی مدل‌های رگرسیون که فاقد جمله ثابت هستند ثابت کنیم. برای این منظور کافی است نشان دهیم که قضیه ۵-۴۵ برای مدل‌های فاقد جمله ثابت نیز صادق است. مدل رگرسیون ۵-۱ را بدون جمله ثابت β_1 می‌نویسیم،

$$Y_i = \beta_2 X_{i2} + \beta_3 X_{i3} + \dots + \beta_k X_{ik} + U_i. \quad (5.54)$$

اگر معادله ۵-۵۴ را دقیقاً مانند سیستم معادله‌های ۵-۳ نوشته و سپس آن را به صورت ماتریسی نمایش دهیم، خواهیم داشت:

$$y = X\beta + u, \quad (5.55)$$

که در آن $n \times 1 \rightarrow y$ و $n \times (k-1) \rightarrow X$ و $(k-1) \times 1 \rightarrow \beta$ و $n \times 1 \rightarrow u$. توجه داریم که

تفاوت معادله ۵-۵۵ با معادله ۵-۴ این است که در معادله ۵-۵۵ اولین ستون ماتریس X که در معادله ۵-۴ شامل عناصر یک است - حذف شده است و بردار β نیز فاقد اولین عنصر، یعنی β_1 ، است. دقیقاً مشابه معادله ۵-۱۹ مقدار $e'e$ را نوشته و از آن نسبت به $\hat{\beta}$ مشتق می‌گیریم. در نتیجه

$$\frac{\partial e'e}{\partial \hat{\beta}} = -2X'y + 2X'X\hat{\beta}.$$

مشتق را مساوی صفر قرار داده و برای $\hat{\beta}$ حل می‌کنیم،

$$\hat{\beta} = (X'X)^{-1} X'y. \quad (5-56)$$

فرمول ۵-۵۶ که برای مدل‌های رگرسیون فاقد جمله ثابت است دقیقاً مشابه فرمول ۵-۲۱ است، با این تفاوت که $(X'X)$ یک ماتریس $(k-1) \times (k-1)$ و $X'y$ نیز یک بردار $(k-1) \times 1$ است. عناصر این ماتریس و بردار برحسب مقادیر اصلی است نه انحراف از میانگین. می‌دانیم $y = \hat{y} + e$ و با مراجعه به $\hat{y} = X\hat{\beta}$ داریم

$$y = X\hat{\beta} + e.$$

از ضرب معادله فوق در ترانهاد خود، خواهیم داشت

$$\begin{aligned} y'y &= (\hat{\beta}'X' + e')(X\hat{\beta} + e), \\ &= \hat{\beta}'X'X\hat{\beta} + \hat{\beta}'X'e + e'X\hat{\beta} + e'e. \end{aligned}$$

دقیقاً مشابه معادله ۵-۲۴، می‌توان نشان داد که از معادله‌های نرمال برای مدل ۵-۵۴ نتایج زیر به دست می‌آید،

$$X'e = 0 \quad \text{یا} \quad e'X = 0,$$

در نتیجه با توجه به $\hat{y} = X\hat{\beta}$ داریم

$$y'y = \hat{y}'\hat{y} + e'e.$$

ملاحظه می‌شود قضیه ۵-۴۵ در مورد مدل‌های رگرسیون فاقد جمله ثابت نیز صدق

می‌کند؛ بنابراین معادله فوق دقیقاً مشابه این قضیه برای حالت عمومی است. با این تفاوت که اولاً، با یک روش تقریباً متفاوتی به دست آمده، ثانیاً تمام عناصر بردارهای آن برحسب مقادیر اصلی است نه انحراف از میانگین. با توجه به اینکه در ادامه بحث، باید مرتباً به عناصر ماتریسها توجه کنیم که آیا برحسب انحراف از میانگین است یا مشاهدات اصلی؛ بنابراین تا جایی که امکان دارد، عبارتها را به زبان غیرماتریسی می‌نویسیم. ابتدا معادله فوق را به صورت غیرماتریسی می‌نویسیم،

$$\sum Y_t^* = \sum \hat{Y}_t^* + \sum e_t^* \quad (0.07)$$

اگر R^* را به صورت زیر تعریف کنیم،

$$R^* = 1 - \frac{\sum e_t^*}{\sum Y_t^*} = \frac{\sum \hat{Y}_t^*}{\sum Y_t^*},$$

آنگاه با توجه به معادله 0.07 خواهیم داشت $0 < R^* < 1$. اما اگر بخواهیم R^* را به روال معمول، یعنی مطابق معادله 0.48 تعریف کنیم:

$$R^* = 1 - \frac{\sum e_t^*}{\sum y_t^*},$$

آنگاه مقدار R^* می‌تواند منفی شود. با توجه به اینکه مقدار e_t^* برای مشاهدات اصلی و انحراف از میانگین یکی است، ابتدا R^* را به صورت زیر می‌نویسیم،

$$R^* = 1 - \frac{\sum Y_t^* - \sum \hat{Y}_t^*}{\sum y_t^*}.$$

می‌دانیم

$$\sum y_t^* = \sum (Y_t - \bar{Y})^* = \sum Y_t^* - n\bar{Y}^*$$

بنابراین

$$R^* = \frac{\sum \hat{Y}_t^* - n\bar{Y}^*}{\sum Y_t^* - n\bar{Y}^*} \quad (0.08)$$

مخرج کسر همواره مثبت است، زیرا برابر با $\sum y_i^2$ است. اما در مواردی که $\sum Y_i^2$ از $n\bar{Y}^2$ کمتر بشود، مقدار R^2 منفی خواهد بود. یادآوری می‌کنیم که:

$$\sum \hat{Y}_i^2 - n\bar{Y}^2 \neq [\sum \hat{y}_i^2 = \sum (\hat{Y}_i - \bar{Y})^2],$$

زیرا در مدل‌هایی که فاقد جمله ثابت هستند اولاً، $\bar{Y} \neq \hat{Y}$ و ثانیاً، $\bar{e} \neq 0$ است. در واقع این دو شرط، همان‌گونه که در معادله‌های ۱-۳۰ و ۱-۳۸ دیدیم، حاصل معادله اول نرمال است که خود از مشتق‌گیری e_i^2 نسبت به جمله ثابت به دست آمده است. بدیهی است در مدل‌های فاقد جمله ثابت، این شرایط نیز محقق نخواهد بود. در پایان، به این نکته نیز توجه داریم که بحث از R^2 در مدل‌هایی است که جمله ثابت ندارد. نکته‌های مربوط به بیان ماتریسی R^2 ، برحسب مشاهدات اصلی در قسمت ۵-۷ مطرح خواهد شد.

۵-۵ ضریب تعیین تعدیل شده یا \bar{R}^2

می‌دانیم با افزایش متغیرهای توضیحی در یک مدل رگرسیون، مقدار R^2 معمولاً زیاد شده، یا حداقل کاهش نمی‌یابد؛ برای مثال، یک مدل رگرسیون شامل یک متغیر توضیحی X_{1t} را در نظر گرفته و مقادیر e_i^2 و $\sum y_i^2$ آن را حساب می‌کنیم. به این مدل یک متغیر توضیحی دیگر، مانند X_{2t} اضافه می‌کنیم. اقتضای روش حداقل مربعات معمولی این است که مجموع مربعات پسماند را حداقل کند. اگر مجموع مربعات پسماند را برای مدل جدید $\sum e_i^{*2}$ بنامیم، می‌توان چنین نوشت

$$RSS^* = \sum e_i^{*2} \leq \sum e_i^2 \quad (5.51)$$

زیرا مکانیسم حداقل‌سازی مجموع مربعات پسماند در مدل دوم به گونه‌ای است که اگر حذف متغیر جدید، یعنی X_{2t} در به حداقل رساندن $\sum e_i^{*2}$ کمکی باشد، قطعاً در مشتق‌گیری، ضریب متغیر X_{2t} ، صفر خواهد شد. دلیل این امر شمولیت مدل دوم نسبت به مدل اول است؛ یعنی مدل دوم شامل مدل اول نیز بوده، بنابراین

مکانیسم حداقل سازی $\sum e_i^2$ در مدل اول را نیز به طور ضمنی دربردارد؛ بنابراین نتیجه می‌گیریم که به موازات افزایش متغیرهای توضیحی، مقدار $\sum e_i^2$ هیچگاه افزایش نخواهد یافت.

یا توجه به معادله ۴-۱۶، فرمول ضریب تعیین برای مدل جدید را به صورت زیر

می‌نویسیم،

$$R^1 = 1 - \frac{\sum e_i^2}{\sum y_i^2}$$

ضریب تعیین برای مدل رگرسیون اول که فقط شامل یک متغیر توضیحی X_{11} است از معادله ۴-۴۳ به دست می‌آید،

$$r^1 = 1 - \frac{\sum e_i^2}{\sum y_i^2}$$

از طرف دیگر می‌دانیم که مقدار کل تغییرات $(\sum y_i^2)$ برای هر دو مدل یکسان است؛ بنابراین با توجه به رابطه ۵-۵۹ خواهیم داشت:

$$R^1 \geq r^1$$

به همین ترتیب می‌توان نشان داد که اگر به مدل دوم، متغیرهای توضیحی جدید اضافه کنیم، مقادیر ضریب تعیین آنها مرتباً زیاد شده، یا حداقل کاهش نمی‌یابد.

از نظر محاسباتی، می‌توان مسأله را از زاویه دیگری نیز حداقل ملاحظه کرد. یک

متغیر توضیحی جدید، فقط صورت کسر معادله ۴-۱۶ را تغییر می‌دهد و مخرج آن را مطلقاً متأثر نمی‌کند. علت این امر، تعریف خاصی است که از ضریب تعیین شده است. در واقع این ضریب را به صورت نسبت دو مجموعه از تغییرات معرفی کرده‌ایم، به گونه‌ای که تغییرات منعکس شده در مخرج، تابعی از تعداد متغیرهای توضیحی نیست. حال اگر به جای تغییر، از مفهوم واریانس استفاده کنیم، دیگر مشکل فوق را نداریم و یک متغیر توضیحی جدید، می‌تواند صورت و مخرج کسر R^1 را همزمان تغییر دهد.

از تعریف ضریب تعیین ۴-۱۶ شروع می‌کنیم،

$$R^2 = \frac{ESS}{TSS} = \frac{\sum \hat{y}_i^2}{\sum y_i^2}$$

صورت و مخرج آن را بر n تقسیم کرده، خواهیم داشت:

$$R^2 = \frac{\sum \hat{y}_i^2 / n}{\sum y_i^2 / n} = 1 - \frac{\sum e_i^2 / n}{\sum y_i^2 / n} \quad (0.60)$$

در معادله ۲-۲۸ ثابت کردیم که اگر بخواهیم واریانس پسماندها یک تخمین نارایب از جمله اختلال باشد، باید $\sum e_i^2$ را بر درجات آزادی آن تقسیم کنیم. مدل جدید را که شامل X_{21} است می‌نویسیم،

$$Y_i = \beta_1 + \beta_2 X_{21} + \beta_3 X_{21} + U_i$$

می‌دانیم برای محاسبه $\sum e_i^2$ در این مدل، سه درجه آزادی از دست می‌دهیم؛ زیرا محاسبه $\sum e_i^2$ ، مستلزم تخمینهای β_1 ، β_2 و β_3 است؛ بنابراین درجات آزادی $\sum e_i^2$ در معادله ۰-۶۰ و برای این مدل برابر $(n-3)$ خواهد بود.

البته بیان دقیقتر مفهوم درجات آزادی برای $\sum e_i^2$ در مدل فوق‌الذکر به این صورت است که اگر بخواهیم سه پارامتر β_1 ، β_2 و β_3 را تخمین بزنیم، آنگاه با توجه به معادله‌های نرمال، یعنی معادله‌های ۴-۱۲، ۴-۱۳ و ۴-۱۴، خواهیم داشت

$$\sum e_i = 0, \sum e_i X_{21} = 0, \sum e_i X_{21} = 0$$

بنابراین مقادیر پسماندها به گونه‌ای تعیین می‌شود که سه شرط فوق در آنها صدق کند. در واقع فقط $(n-3)$ مقدار e_i می‌تواند به «آزادی» مقدار بگیرد، بنابراین می‌گوییم درجات آزادی $\sum e_i^2$ برابر $(n-3)$ است؛ به عبارت دیگر، به ازای هر مجموعه از مقادیری که $(n-3)$ پسماند اختیار می‌کند، مقادیر سه پسماند باقیمانده، فقط با حل یک دستگاه معادله، شامل سه شرط فوق امکانپذیر است؛ یعنی در واقع برای محاسبه $\sum e_i^2$ ، سه درجه آزادی از دست داده‌ایم. نتیجه می‌گیریم که باید $\sum e_i^2$ در معادله ۰-۶۰ را به جای n بر $(n-3)$ تقسیم کنیم. همچنین برای اینکه بتوان به تخمین نارایبی از واریانس

متغیر درون‌زا رسید، باید $\sum y_t^2$ در معادله ۵-۶۰ را بر $(n-1)$ تقسیم کنیم؛ زیرا محاسبه $\sum y_t^2$ مستلزم محاسبه \bar{Y} است؛ بنابراین یک درجه آزادی از دست خواهیم داد. اصطلاحاً می‌گوییم R^2 را نسبت به درجات آزادی تعدیل کرده‌ایم. به ضریب تعیینی که بدین ترتیب به دست می‌آید، «ضریب تعیین تعدیل شده» گفته و آن را با \bar{R}^2 نشان می‌دهند.

$$\bar{R}^2 = 1 - \frac{\sum e_t^2 / (n-3)}{\sum y_t^2 / (n-1)},$$

یا

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-3}\right) \frac{\sum e_t^2}{\sum y_t^2}. \quad (5-61)$$

با استفاده از تعریف ضریب تعیین، یعنی معادله ۴-۱۶، می‌دانیم که:

$$\frac{\sum e_t^2}{\sum y_t^2} = 1 - R^2.$$

با جایگزینی معادله فوق در معادله ۵-۶۱، خواهیم داشت

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-3}\right) (1 - R^2), \quad (5-62)$$

که در آن R^2 ، ضریب تعیین مدلی است که دارای دو متغیر توضیحی است. معادله ۵-۶۲ رابطه بین ضریب تعیین معمولی و ضریب تعیین تعدیل شده را نشان می‌دهد. بحث فوق را برای یک مدل رگرسیون با k متغیر توضیحی تعمیم می‌دهیم. مدل ۵-۱ را یک بار دیگر می‌نویسیم؛

$$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \dots + \beta_k X_{kt} + U_t.$$

با استفاده از تعریف ضریب تعیین در رابطه ۴-۱۶ می‌توان نوشت،

$$R^2 = \frac{ESS}{TSS} = \frac{\sum \hat{y}_t^2}{\sum y_t^2} = 1 - \frac{\sum e_t^2}{\sum y_t^2}.$$

می‌دانیم برای محاسبه مجموع مربعات پسماند، یعنی $\sum e_i^2$ ، در مدل فوق باید $\hat{\beta}_1$ ، $\hat{\beta}_2$ تا $\hat{\beta}_k$ را داشته باشیم؛ بنابراین k درجه آزادی از دست داده، در نتیجه درجات آزادی برای $\sum e_i^2$ برابر $(n - k)$ است؛ با وجود این، درجات آزادی برای $\sum y_i^2$ همان $(n - 1)$ خواهد بود. با تقسیم $\sum e_i^2$ و $\sum y_i^2$ بر درجات آزادی آنها به \bar{R}^2 می‌رسیم.

$$\bar{R}^2 = 1 - \frac{\sum e_i^2 / (n - k)}{\sum y_i^2 / (n - 1)},$$

یا

$$\begin{aligned} \bar{R}^2 &= 1 - \left(\frac{n-1}{n-k}\right) \frac{\sum e_i^2}{\sum y_i^2}, \\ &= 1 - \left(\frac{n-1}{n-k}\right) \frac{e'e}{y'y}. \end{aligned} \quad (5-63)$$

با توجه به معادله ۴-۱۶ داریم:

$$\frac{\sum e_i^2}{\sum y_i^2} = \frac{e'e}{y'y} = 1 - R^2,$$

با جایگزینی در معادله ۵-۶۳، خواهیم داشت:

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k}\right) (1 - R^2), \quad (5-64)$$

که در آن n ، تعداد مشاهدات و k ، تعداد پارامترهایی است که در مدل رگرسیون باید تخمین زد. یادآوری می‌شود که در این مدل $(k - 1)$ متغیر توضیحی داریم. با ملاحظه معادله ۵-۶۴ به سهولت می‌توان گفت که

$$\bar{R}^2 \leq R^2. \quad (5-65)$$

\bar{R}^2 ، در واقع آماره‌ای است که خودمان ساخته‌ایم و همبستگی هیچیک از دو متغیر را نشان نمی‌دهد. به همین دلیل، برای مقادیر کوچک R^2 ، می‌توان نشان داد که \bar{R}^2 ممکن است منفی شود. برای این منظور، ابتدا معادله ۵-۶۴ را به صورت زیر می‌نویسیم،

$$1 - \bar{R}^2 = \left(\frac{n-1}{n-k}\right) (1 - R^2). \quad (5-66)$$

به ازای $R^2 < \frac{k-1}{n-1}$ ، خواهیم داشت:

$$(1-R^2) > 1 - \frac{k-1}{n-1},$$

یا

$$(1-R^2) > \frac{n-k}{n-1}.$$

با جایگزینی رابطه فوق در معادله ۵-۶۶، داریم:

$$1 - \bar{R}^2 > 1,$$

یا

$$\bar{R}^2 < 0.$$

(۵-۶۷)

برای مثال، فرض کنید، می خواهیم مدلی را تخمین بزنیم که ۴ پارامتر دارد. با ۲۱ مشاهده روی متغیرها داریم: $k=4$ و $n=21$. اگر R^2 برابر ۰/۱ باشد، چون

$$R^2 < \frac{k-1}{n-1},$$

$$0/1 < \frac{4-1}{21-1},$$

بنابراین \bar{R}^2 منفی خواهد بود. برای ارزیابی صحت این نتیجه می توان مستقیماً از معادله ۵-۶۴ نیز استفاده کرد

$$\bar{R}^2 = 1 - \left(\frac{21-1}{21-4} \right) (1-0/1),$$

$$= 1 - \frac{20}{17} (0/1) = 1 - \frac{18}{17} = -0/058.$$

تأثیر متغیرهای توضیحی در \bar{R}^2 و $\hat{\sigma}_u^2$

می توان ثابت کرد که آن دسته از متغیرهای توضیحی که واریانس جمله اختلال، یعنی $\hat{\sigma}_u^2$ را حداقل کنند، ضریب تعیین تعدیل شده، یعنی \bar{R}^2 ، را حداکثر خواهد کرد. برای اثبات کافی است دو طرف معادله ۵-۶۶ را بر $(n-1)$ تقسیم کنیم،

$$\frac{(1-\bar{R}^2)}{(n-1)} = \frac{(1-R^2)}{(n-k)}.$$

دو طرف معادله فوق را در $\sum y_i^2$ ضرب می‌کنیم،

$$\frac{(1 - \bar{R}^2) \sum y_i^2}{(n-1)} = \frac{(1 - R^2) \sum y_i^2}{(n-k)} \quad (5-68)$$

با توجه به معادله ۵-۱۶، می‌دانیم $\sum y_i^2 = \sum e_i^2 (1 - R^2)$ است. همچنین در معادله‌های ۵-۲۸ و ۵-۶۱ دیدیم که اگر مجموع مربعات پسماند، یعنی $\sum e_i^2$ را بر درجات آزادی آن، یعنی $(n-k)$ تقسیم کنیم، یک تخمین نااریب از واریانس جمله اختلال خواهیم داشت؛ بنابراین

$$\frac{(1 - R^2) \sum y_i^2}{(n-k)} = \frac{\sum e_i^2}{(n-k)} = \hat{\sigma}_u^2 \quad (5-69)$$

معادله ۵-۶۹ را در معادله ۵-۶۸ جایگزین می‌کنیم،

$$\frac{(1 - \bar{R}^2) \sum y_i^2}{(n-1)} = \frac{\sum e_i^2}{(n-k)} = \hat{\sigma}_u^2 \quad (5-70)$$

مشاهده می‌شود که به موازات افزایش تعداد متغیرهای توضیحی، یعنی با افزایش k ، مقدار $(n-k)$ کمتر شده؛ بنابراین مخرج کسر $\frac{\sum e_i^2}{n-k}$ کاهش می‌یابد. می‌دانیم با افزایش متغیرهای توضیحی، مقدار $\sum e_i^2$ کاهش یافته یا حداقل ثابت می‌ماند؛ بنابراین در بهترین حالت، یک متغیر توضیحی جدید می‌تواند، صورت و مخرج کسر $\hat{\sigma}_u^2$ را همزمان کاهش دهد. اما اینکه آیا $\hat{\sigma}_u^2$ کم شده یا افزایش می‌یابد، در فصل هفتم و در قسمت ۳-۷ بررسی خواهد شد. در این قسمت فقط می‌خواهیم رابطه بین $\hat{\sigma}_u^2$ و \bar{R}^2 را بررسی کنیم. با فرض ثبات n و $\sum y_i^2$ و با توجه به معادله ۵-۷۰، می‌توان نتیجه گرفت که جهت تغییرات $\hat{\sigma}_u^2$ و $(1 - \bar{R}^2)$ به ازای افزایش یا کاهش متغیرهای توضیحی یکسان است. این نتیجه را می‌توان به این صورت نیز بیان کرد که به ازای افزایش متغیرهای توضیحی، تغییرات $\hat{\sigma}_u^2$ در جهت مخالف تغییرات \bar{R}^2 خواهد بود؛ یعنی آن دسته از متغیرهای توضیحی که $\hat{\sigma}_u^2$ را کاهش می‌دهند \bar{R}^2 را افزایش خواهند داد.

۵-۶ نکته‌هایی دربارهٔ R^2 و \bar{R}^2

به نظر می‌رسد قبل از پایان دادن به مباحث ضریب تعیین لازم باشد بعضی مسائل را فهرست وار مطرح کنیم.

۱. با توجه به معادلهٔ ۴-۱۶ داریم:

$$R^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2}$$

می‌دانیم تخمین زننده حداقل مربعات معمولی، مجموع مربعات پسماند، یعنی $\sum e_i^2$ ، را حداقل می‌کند. با توجه به فرمول فوق، ملاحظه می‌شود که وقتی $\sum e_i^2$ حداقل باشد، R^2 حداکثر خواهد بود. بدین ترتیب می‌توان گفت که معیار روش حداقل مربعات معمولی در تخمین پارامترها حداکثر نمودن R^2 است. نتیجه می‌گیریم که معیارهای حداقل کردن مجموع مربعات پسماند و حداکثر نمودن R^2 در واقع یکی هستند.

مسئلهٔ حداکثرسازی R^2 یک مبحث بسیار مهم در اقتصادسنجی است. وقتی در تخمین یک مدل رگرسیون به مقداری از ضریب تعیین می‌رسیم که نزدیک به یک است، معمولاً این امر بدین صورت تفسیر می‌شود که پارامترهای مدل بسیار خوب تخمین زده شده است. در بین نکاتی که در این قسمت مطرح خواهیم کرد، سعی می‌کنیم به این سؤال پاسخ دهیم که آیا واقعاً چنین استنتاجی صحیح است؟ ابتدا باید دو حالت زیر را از یکدیگر تفکیک کرد.

در حالت اول، فرض بر این است که مدل رگرسیون از هر نظر دقیقاً تعریف شده است به گونه‌ای که ما حق هیچگونه تصرفی را در آن نداریم؛ به عبارت دقیقتر، اولاً شکل ریاضی تابع رگرسیون کاملاً مشخص است که مثلاً خطی است یا غیرخطی، ثانیاً تعداد و نوع متغیرهای توضیحی کاملاً شناخته شده است به نحوی که نه می‌توان متغیری به مدل اضافه کرد و نه متغیر یا متغیرهایی از آن را حذف کرد، و ثالثاً فرضهای جملهٔ اختلال را نیز به طور دقیق و غیرقابل تغییری قبول کرده‌ایم. در چنین شرایطی از ما خواسته شده است که پارامترهای این مدل را تخمین بزنیم. نکتهٔ مهم این است که فقط در این حالت

می توان گفت که معیار روش حداقل مربعات معمولی همان حداکثرسازی R^2 است؛ یعنی باید پارامترهای مدل را چنان تخمین زد که \hat{Y}_i بتواند در بالاترین سطح ممکن، تغییرات Y_i را توضیح دهد. بدیهی است این معیار چیزی نیست جز حداقل سازی مجموع تغییرات توضیح داده نشده یا پسماندها، که همان معیار متعارف حداقل مربعات معمولی است.

در حالت دوم، فرض می کنیم که می خواهیم پارامترهای یک مدل رگرسیون را چنان تخمین بزنیم که R^2 حداکثر شود. اما برای رسیدن به این هدف نه تنها اجازه داریم شکل ریاضی تابع رگرسیون را به هر نحوی که مقتضی باشد تغییر دهیم، بلکه می توانیم هر متغیر توضیحی جدیدی را که بخواهیم به مدل اضافه کنیم یا هر متغیر دیگری را نیز هنگامی که لازم است، حذف کنیم. همچنین می توانیم هر فرض جدیدی را نیز که موجب افزایش R^2 شود مطرح کنیم، یا در فرضهایی که قبلاً ارائه شده است به سهولت تصرف کنیم. به طور خلاصه، هر اقدامی که بتواند به افزایش R^2 منجر شود، مجاز خواهد بود. همه بحث ما این است که چنین تفسیری از حداکثر نمودن R^2 کاملاً غلط است و هیچگونه ارزشی از نظر کاربردی ندارد؛ زیرا ماهیت این نحوه برخورد با R^2 این است که سعی کنیم مدلی بسازیم که بتواند از مجموعه مشاهدات موجود در نمونه مفروض یک R^2 بسیار بالا نتیجه دهد^۱. در صورت موفقیت در ساختن چنین مدل مصنوعی، بدیهی است اگر نمونه را تغییر داده و با مجموعه دیگری از مشاهدات کار کنیم، چه بسا مدلی که با اینهمه زحمت ساخته ایم، R^2 بسیار پایینی را نشان دهد.

نتیجه ای که می گیریم، این نیست که اولاً نباید به R^2 توجه چندانی کرد و ثانیاً نباید در طراحی مدل هیچگونه تغییری داد. بلکه برعکس، باید با توجه به مقادیر جملات پسماند، یعنی تغییرات توضیح داده نشده، سعی نمود مدل را به نحوی تغییر داد که قدرت توضیحی بیشتری داشته باشد. که در واقع دقیقاً به معنای افزایش R^2 است. اما نکته مهم این است که این تغییراتی که در مدل می دهیم، باید مجوز نظری داشته باشد. نارسایی عملکرد تخمین ما در تفسیر تغییرات مشاهده شده در متغیر درون‌زا، به طور

۱. در اقتصادسنجی اصطلاحاً به این روش، Data Mining می گویند.

کلی روشن کننده ضعف نظری مدل رگرسیون است. همه کوشش باید این باشد که بتوان در ساختار ریاضی مدل و فرضهای آن، یا فرضهای جمله اختلال چنان تصرف کرد که این ضعف نظری کاهش یابد. بدیهی است هر اقدام مؤثری در این زمینه ممکن است به افزایش R^2 منجر شود، اما چه بسا ضرورتاً R^2 بسیار بالایی نتیجه ندهد. بنابراین، اگر دو مدل در تفسیر تغییرات یک متغیر درون‌زا داشته باشیم که اولی R^2 بسیار بالا مثلاً $0/99$ و دومی R^2 مثلاً $0/52$ داشته باشد، اما مدل اول چنان «ساخته» شده باشد که R^2 بالا نتیجه دهد، در حالی که مدل دوم از مبنای نظری قوی تری برخوردار باشد، آنگاه قطعاً مدل دوم مرجح است. در بین نکاتی که مطرح شد دوباره به این بحث اشاره خواهد شد. اما به این نکته اشاره می‌کنیم که یکی از راههای تشخیص اینکه آیا یک مدل، تنها برای رسیدن به R^2 بالا طراحی شده یا خیر، این است که ابتدا به موازین نظری آن توجه شود و سپس از آزمونهای مختلف آماری برای معنی دار بودن پارامترهای آن استفاده کرده و سرانجام با ایجاد تغییر مختصر در نمونه، تغییرات R^2 را ارزیابی کنیم.

۲. معمولاً این سؤال مطرح می‌شود که « R^2 بالا» دقیقاً به چه معنی است؟ بدیهی است نمی‌توان به این سؤال به طور دقیق پاسخ داد، اما توجه به نکات کلی زیر مفید است. در مواردی که مشاهدات ما در مورد متغیرهای درون‌زا و برون‌زای یک مدل به صورت سریهای زمانی است، R^2 مقدار بسیاری نشان می‌دهد؛ زیرا تغییر متغیرها در زمان، روند تقریباً مشابهی دارند. می‌توان گفت^۱ که اگر یک متغیر اقتصادی را به طور تصادفی انتخاب کنیم، مثلاً Y_t و مقدار این متغیر را در دوره قبل در نظر بگیریم، Y_{t-1} ، و سپس مدل رگرسیون Y_t بر Y_{t-1} را بسازیم، به طور متوسط R^2 حدود $0/7$ خواهد بود. همچنین اگر یک متغیر اقتصادی را به طور تصادفی انتخاب کرده و رگرسیون آن را با دو تا شش متغیر اقتصادی دیگر که آنها نیز به طور تصادفی انتخاب می‌شود، در نظر بگیریم، مقدار R^2 به طور متوسط از $0/5$ بیشتر خواهد شد به شرط اینکه مشاهدات مربوط به متغیرها برحسب سری زمانی باشد. با وجود این، باید توجه کرد که وقتی در یک مدل

۱. به مقاله (۱۹۶۱) E. Ames and S. Reiter مراجعه شود.

رگرسیون، مشاهدات به صورت مقطعی باشد، مقادیر R^2 در مقایسه با حالت سری زمانی بسیار کمتر خواهد بود.

برای تبیین این مسأله مثالی می‌زنیم. فرض کنید مقادیر یک متغیر اقتصادی به طور مداوم در زمان، یک سیر نزولی را نشان می‌دهد، به گونه‌ای که می‌توان مسیر تغییرات زمانی آن را با معادله $Y_t = 50 - 0.2t$ نشان داد که در آن t زمان است. همچنین فرض می‌کنیم یک متغیر دیگری نیز، مانند X_t ، داریم که اساساً هیچگونه ارتباطی با Y_t ندارد، اما ملاحظه می‌کنیم که مقادیر آن در زمان - به طور مرتب زیاد می‌شود. اگر مسیر تغییرات زمانی X_t را از نظر ریاضی با معادله $X_t = 5 + 0.4t$ ، تقریب کنیم، به راحتی می‌توان از این معادله، t را برحسب X_t به دست آورده و در معادله Y_t قرار داد تا بدین ترتیب معادله Y_t برحسب X_t به دست آید،

$$t = 2/5 X_t - 12/5,$$

$$Y_t = 50 - 0.2(2/5 X_t - 12/5),$$

یا:

$$Y_t = 52/5 - 0/5 X_t.$$

بدیهی است با چنین ساختار ریاضی بین متغیرها، اگر با استفاده از سری زمانی تغییرات Y_t و X_t ، مدل رگرسیون

$$Y_t = \alpha + \beta X_t + U_t,$$

را تخمین بزنیم، قطعاً R^2 بسیار بالایی را به دست خواهیم آورد، و حال آنکه می‌دانیم با توجه به موازین نظری، اساساً هیچگونه رابطه‌ای بین X_t و Y_t موجود نیست.

نتیجه دیگری که از مثال فوق می‌توان گرفت این است که در یک مدل رگرسیون Y_t روی X_t ، اگر مقدار عددی R^2 حتی نزدیک به یک باشد، مطلقاً دلالت نمی‌کند که ضرورتاً یک رابطه علت و معلولی قوی بین X_t و Y_t وجود دارد. اثبات وجود رابطه علت و معلولی، در حوزه تحلیل‌های نظری صورت می‌پذیرد. در مثال قبل دیدیم که اساساً هیچگونه رابطه‌ای بین X_t و Y_t وجود نداشت؛ در چنین مواردی، مقدار عددی R^2 ، تنها

منعکس کننده روند موجود در متغیرهای X_t و Y_t است. راه حلی که معمولاً برای خنثی کردن تأثیر روند بر مقدار R^2 می توان پیشنهاد کرد، این است که به جای استفاده از مقادیر مطلق متغیرهای درونزا و برونزا از نسبت یا نرخ تغییرات آنها استفاده شود، زیرا وقتی با نرخ تغییرات یک متغیر کار می کنیم، تأثیر روند را تا حد بسیاری از بین خواهیم برد. برای مثال، در اکثر کشورهای صنعتی ملاحظه می شود که بعد از جنگ جهانی دوم، تغییرات تولید و تورم، یک روند صعودی را نشان داده است، در حالی که این روند در نرخ رشد تولید یا نرخ رشد تورم ملاحظه نمی شود. با وجود این، کار کردن با نرخ تغییر متغیرها نیز مسائل حاصل خود را دارد که بحث بیشتر درباره آنها از موضوع این قسمت خارج است و باید در اقتصادسنجی کاربردی بررسی شود.

۳. می دانیم که یکی از اهداف مهم روشهای تخمین در اقتصادسنجی این است که بتوان به تخمینهای «خوبی» از پارامترهای مدل رگرسیون رسید. با وجود این، نباید «خوبی» تخمین پارامترها را با توجه به معیار «مقدار R^2 » تفسیر کرد؛ بدین معنی که اگر R^2 زیاد باشد نتیجه بگیریم که تخمین پارامترها خوب است و برعکس البته تخمینهای خوب از پارامترهای یک مدل ممکن است موجب افزایش مقدار R^2 شود، اما عکس آن صحیح نیست. بنابراین، می توان نتیجه گرفت که معیار R^2 نقش چندانی در استنتاج نهایی درباره خوبی تخمین مدل رگرسیون ندارد.

با اینکه در اقتصادسنجی نسبت به عدم اهمیت زیاد مقدار R^2 ، اتفاق نظر وجود دارد، اما در تمام محاسبات کامپیوتری در تخمین پارامترهای یک مدل رگرسیون، حتماً مقدار R^2 منعکس می شود و همواره مقادیر بالای R^2 به عنوان معیاری در موفقیت تخمین یک مدل معرفی می شود.^۱

۴. با توجه به معادله ۵-۴۸ می دانیم R^2 تابعی از $\sum e_t^2$ و $\sum y_t^2$ و در نتیجه تابعی از e_t و y_t است؛ بنابراین در نمونه گیریهای تکراری، R^2 از یک نمونه به نمونه دیگر تغییر می کند. در واقع، چون e_t یک متغیر تصادفی است، R^2 نیز یک متغیر تصادفی خواهد بود

۱. برای توضیحات بیشتر به مقاله (۱۹۸۷) J. S. Cramer مراجعه شود.

که دارای توزیع احتمال است و می‌توان فرضیه‌های مختلف در مورد آن را آزمود. در معادله ۲-۵۴ نیز دیدیم که می‌توان آماره F را برحسب r^2 نوشت. در فصل آینده و در قسمت ۶-۴ نیز دوباره به این مسأله مراجعه نموده و آن را برای حالت رگرسیون چندمتغیره مطرح خواهیم کرد. در این قسمت به مفهوم آزمون فرضیه $H_0: R^2 = 0$ می‌پردازیم:

در مدل رگرسیون

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i,$$

فرض بر این است که تمام پارامترهای β غیر از β_1 صفر است؛ یعنی متغیر درون‌زای Y_i مطلقاً تحت تأثیر متغیرهای توضیحی X_2, X_3, \dots, X_k نیست. چنین فرضی دلالت بر این می‌کند که در جامعه آماری، مقدار مجموع تغییرات توضیح داده شده، یعنی $\sum \hat{y}_i^2$ ، برابر صفر می‌شود؛ بنابراین مقدار R^2 جامعه نیز صفر خواهد شد. اما در عمل ما فقط یک نمونه داریم و با اینکه در واقع، تمام پارامترهای مدل غیر از β_1 صفر است در تخمین حاصل از نمونه مفروض، برای اکثر پارامترها، تخمین‌هایی غیر صفر خواهیم داشت؛ بنابراین مقداری نیز برای R^2 به دست خواهد آمد که ضرورتاً غیر صفر است. سؤال این است که چگونه به کمک R^2 به دست آمده از نمونه می‌توان درباره این فرضیه - که R^2 جامعه صفر است - قضاوت آماری کرد؟ در فصل آینده و در معادله ۶-۴۱ خواهیم دید که کمیت

$$\left(\frac{n-k}{k-1} \right) \left(\frac{R^2}{1-R^2} \right),$$

دارای توزیع F با $(k-1, n-k)$ درجه آزادی است.

۵. در ارائه بحث فوق می‌توان به یک مورد دیگر نیز اشاره کرد. فرض کنید در یک مدل رگرسیون چند متغیره یک یا چند متغیر توضیحی دیگر وارد مدل کنیم. سؤال این است که با چه معیاری می‌توان ضرورت یا اهمیت متغیرهای جدید در مدل رگرسیون مفروض را ثابت کرد؟ با توجه به مباحثی که در ابتدای قسمت ۵-۵ و در تعریف \bar{R}^2 گفتیم، به نظر می‌رسد که چون ورود هر متغیر توضیحی جدید باعث می‌شود که R^2

افزایش یافته و یا حداقل کاهش نیابد بنابراین با معیار R^2 ، مدلی که شامل تعداد بیشتری از متغیرهای توضیحی شده است بر مدل قبلی مرجح خواهد بود. همچنین در مباحث قبلی دیدیم که ورود متغیر یا متغیرهای توضیحی، درجات آزادی $\sum e_i^2$ و $\sum y_i^2$ را تغییر می دهد؛ پس معیار \bar{R}^2 نسبت به R^2 از برتری نسبی برخوردار می شود. در واقع، می توان حالتی را در نظر گرفت که اضافه کردن یک متغیر توضیحی جدید، موجب افزایش R^2 و کاهش \bar{R}^2 می شود.

برای توضیح بیشتر این نکته، معادله ۵-۶۳ را یک بار دیگر می نویسیم،

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k} \right) \frac{\sum e_i^2}{\sum y_i^2}.$$

فرض کنید $n=21$ ، $k=6$ ، $\sum e_i^2=10$ و $\sum y_i^2=100$ است؛ بنابراین:

$$\bar{R}^2 = 1 - \left(\frac{21-1}{21-6} \right) \frac{10}{100} = 1 - 0/20 = 0/80.$$

برای همین مثال، مقدار R^2 را نیز محاسبه می کنیم:

$$R^2 = 1 - \frac{\sum e_i^2}{\sum y_i^2} = 1 - \frac{10}{100} = 0/80.$$

یک متغیر توضیحی دیگر وارد مدل می کنیم. فرض بر این است که این متغیر جدید، قدرت توضیحی بسیاری ندارد؛ با وجود این، توانسته است مجموع مربعات پسماند را تا حدی کاهش دهد، به گونه ای که داریم، $\sum e_i^2 = 14/5$. مقدار \bar{R}^2 در این حالت به قرار زیر خواهد بود،

$$\bar{R}^2 = 1 - \left(\frac{21-1}{21-7} \right) \frac{14/5}{100} = 1 - 0/207 = 0/793,$$

یعنی مقدار \bar{R}^2 کاهش یافته است. مفید است مقدار R^2 را نیز حساب کنیم،

$$R^2 = 1 - \frac{14/5}{100} = 0/850,$$

ملاحظه می شود که R^2 افزایش نشان می دهد.

نتیجه می‌گیریم که معیار \bar{R}^2 از R^2 دقیقتر است. با وجود این، برای گرفتن تصمیم نهایی در مورد حفظ یا حذف متغیرهای توضیحی جدید، باید علاوه بر دقت در موازن نظری، از آزمونهای آماری نیز استفاده کرد که موضوع بحث فصل آینده خواهد بود.

۶. در فصل نهم و در قسمت ۹-۳ به این نکته اشاره خواهیم کرد که اگر یک متغیر توضیحی جدید به مدل رگرسیون اضافه کنیم، موقعی \bar{R}^2 زیاد می‌شود که مقدار t مربوط به تخمین پارامتر این متغیر توضیحی جدید بیشتر از یک باشد. بنابراین، ملاحظه می‌شود که مسأله حداکثرسازی \bar{R}^2 کاملاً با قاعده حفظ متغیرهای توضیحی در یک مدل رگرسیون متفاوت است؛ زیرا می‌دانیم موقعی می‌توان یک متغیر توضیحی مانند X_{it} را مفید و ضروری تشخیص داد که فرضیه $H_0: \beta_i = 0$ در سطح معنی‌دار، مثلاً ۵ درصد رد شود. اگر این فرضیه با $t < 1$ رد شود، نتیجه می‌گیریم که متغیر توضیحی مورد نظر باید در مدل حفظ شود، در حالی که می‌دانیم ورود این متغیر، مقدار \bar{R}^2 را افزایش نداده است؛ زیرا مقدار t متعلق به تخمین پارامتر آن از یک کمتر است. نتیجه می‌گیریم که در مورد حفظ یا حذف یک متغیر توضیحی، دو معیار آزمونهای آماری و \bar{R}^2 ممکن است در مواردی به یک استنتاج منتهی نشود.

مثال ۵-۳ مدل رگرسیون را - که در مثال ۵-۲ مطرح کردیم - یک بار دیگر در نظر بگیرید:

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

با استفاده از جدول ۵-۱،

اولاً، R^2 را یک بار با فرمول ۵-۴۰ و بار دیگر با استفاده از فرمول ۵-۴۸ به دست آورید.

ثانیاً، مقدار \bar{R}^2 را محاسبه کنید.

۱. ابتدا فرمول ۵-۴۰ را می‌نویسیم،

$$R^2 = \frac{\hat{\beta}' X' y}{y' y}.$$

می دانیم عناصر بردارهای موجود در فرمول فوق باید بر حسب انحراف از میانگین باشد. بنابراین، ضرورتاً از محاسبات مربوط به مثال ۵-۲ استفاده می کنیم. با توجه به جدول ۵-۳ می دانیم

$$\sum x_i y_i = 16 \quad , \quad \sum z_i y_i = 9 .$$

همچنین در محاسبات مثال ۵-۲ دیدیم

$$\hat{\beta} = 2/5 \quad , \quad \hat{\gamma} = -1/5 .$$

برای محاسبه R^2 به $\sum y_i^2$ نیز نیاز داریم که باید محاسبه کنیم،

$$\sum y_i^2 = y'y = 1 + 9 + 16 + 1 + 1 = 28 .$$

صورت کسر R^2 ، یعنی $\hat{\beta}' X' y$ عبارت است از

$$\begin{aligned} \hat{\beta}' X' y &= [\hat{\beta} \quad \hat{\gamma}] \begin{bmatrix} \sum x_i y_i \\ \sum z_i y_i \end{bmatrix} = \hat{\beta} \sum x_i y_i + \hat{\gamma} \sum z_i y_i , \\ &= 2/5 (16) - 1/5 (9) = 40 - 13/5 = 26/5 . \end{aligned}$$

بنابراین

$$R^2 = \frac{26/5}{28} = 0.9286 ,$$

یعنی تخمین ما از مدل توانسته است حدود ۹۵ درصد تغییرات Y_i را توضیح دهد. اگر بخواهیم R^2 را با استفاده از فرمول ۵-۴۸ به دست آوریم، باید چنین نوشت،

$$R^2 = 1 - \frac{e'e}{y'y} .$$

برای محاسبه $e'e$ ابتدا معادله ۵-۴۵ را دوباره می نویسیم،

$$y'y = \hat{\gamma}' \hat{\gamma} + e'e ,$$

$$28 = 26/5 + e'e ,$$

در نتیجه

$$e'e = 1/0.$$

بنابراین مجموع مربعات پسماند یا تغییرات توضیح داده نشده برابر ۱/۵ است؛ در نتیجه

$$R^2 = 1 - \frac{1/0}{28} = \frac{28 - 1/0}{28} = 0/9464.$$

۲. طبق معادله ۵-۶۴ می دانیم که

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k} \right) (1 - R^2),$$

در نتیجه خواهیم داشت

$$\begin{aligned} \bar{R}^2 &= 1 - \left(\frac{5-1}{5-3} \right) (1 - 0/9464), \\ &= 1 - 2 (0/0536) = 0/8928. \end{aligned}$$

۵-۷ بیان ماتریسی R^2 برحسب مشاهدات اصلی*

تا کنون در تمام فرمولهایی که برای رگرسیون چند متغیره استخراج کرده ایم هیچگاه علامت خاصی را به کار نبرده ایم که بیان کننده این واقعیت باشد که محاسبات برحسب انحراف از میانگین انجام شده است. دلیل این امر، پرهیز از پیچیده تر کردن علائم در بیان فرمولهاست. تا زمانی که می توان از متن بحث به این نکته پی برد که آیا محاسبات با استفاده از مشاهدات اصلی است یا انحراف از میانگین، ضرورتی برای معرفی علائم جدید نیست. اما مطرح کردن این سؤال مفید است که آیا روشی وجود دارد که بتوان به آن وسیله ماتریسهایی را که برحسب مشاهدات اصلی محاسبه شده است، به مقادیر انحراف از میانگین تبدیل کرد؟ در صورتی که پاسخ مثبت باشد به راحتی می توان در مواردی که به مقادیر انحراف از میانگین نیاز است، از این تبدیل استفاده کرده و آن را مستقیماً روی مشاهدات اصلی اجرا کرد. این مسأله، بویژه در محاسبه R^2 بسیار مهم است؛ زیرا می دانیم در فرمول

$$R^2 = \frac{\hat{y}'\hat{y}}{y'y},$$

مقادیر صورت و مخرج باید برحسب انحراف از میانگین محاسبه شود. سؤال این است که چگونه می توان در محاسبه R^1 مستقیماً از ماتریسهای استفاده کرد که برحسب مشاهدات اصلی تنظیم شده است.

برای پاسخ به این سؤال، ابتدا ماتریس Λ را معرفی می کنیم و آن را ماتریس تبدیل می نامیم،

$$\Lambda = I - \frac{1}{n} \mathbf{1}\mathbf{1}' \quad (5.71)$$

که در آن $n \times n \rightarrow \Lambda$ یک بردار شامل n عدد یک و I یک ماتریس $n \times n$ و تمام عناصر قطری آن یک و عناصر غیرقطری آن صفر است (ماتریس یکه)،

$$\Lambda = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} - \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad (5.72)$$

نشان می دهیم که اگر ماتریس تبدیل (Λ) را در هر برداری از سمت چپ ضرب کنیم، بردار به دست آمده برحسب انحراف از میانگین خواهد بود. برای این منظور، بردار ستونی y ، شامل n مقدار را در نظر می گیریم. می دانیم y برحسب مشاهدات اصلی است، یعنی

$$y = [Y_1 \quad Y_2 \quad \dots \quad Y_n]$$

ماتریس تبدیل (Λ) را از سمت چپ در y ضرب می کنیم.

$$\Lambda y = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} - \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

$$= \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} - \begin{bmatrix} \bar{Y} \\ \bar{Y} \\ \vdots \\ \bar{Y} \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} .$$

ملاحظه می‌شود Λy ، بردار y را به برداری تبدیل کرد که عناصر آن انحراف از میانگین مقادیر Y است.

در اینجا اشاره به دو مورد خاص ضروری است. نکته اول این است که اگر Λ را در هر برداری ضرب کنیم که عناصر آن با یکدیگر مساوی است، یک بردار صفر خواهیم داشت؛ مثلاً

$$\Lambda \mathbf{1} = \mathbf{0} . \quad (0.73)$$

نکته دوم این است که چون میانگین مقادیر پسماند صفر است، بنابراین

$$\Lambda \mathbf{e} = \mathbf{e} , \quad (0.74)$$

یعنی بردار \mathbf{e} همواره برحسب انحراف از میانگین است. همچنین می‌توان نشان داد که Λ یک ماتریس متقارن و «خودتوان» است؛ یعنی $\Lambda \Lambda = \Lambda' = \Lambda$ و $\Lambda' = \Lambda$.

اکنون به بررسی مدل‌های رگرسیون می‌پردازیم. مدل رگرسیون ۰-۱ یا ۰-۴ مفروض است

$$y = X\beta + u ,$$

که در آن، $y \rightarrow n \times 1$ ، $X \rightarrow (n \times k)$ ، $\beta \rightarrow k \times 1$ و $u \rightarrow n \times 1$. فرض بر این است که مقادیر تمام متغیرهای این مدل برحسب مشاهدات اصلی بیان شده است. با توجه به معادله ۰-۱۷ می‌دانیم

$$y = X\hat{\beta} + e . \quad (0.75)$$

ماتریس X را به صورت زیر افراز می‌کنیم،

$$X = [X_1 \quad X_2],$$

که در آن، $X_1 = i$ ، یک بردار ستونی شامل اعداد یک است که منعکس کننده جمله ثابت در مدل رگرسیون ۵-۱ است. X_2 نیز یک ماتریس $n \times (k-1)$ از مشاهدات متغیرهای توضیحی X_2 ، X_3 تا X_k است. با جایگزینی معادله فوق در معادله ۵-۷۵ داریم

$$y = i\hat{\beta}_1 + X_2\hat{\beta}_2 + e, \quad (5.76)$$

که در آن $\hat{\beta}' = [\hat{\beta}_1 \quad \hat{\beta}_2]$ بدین معنی است که بردار $\hat{\beta}$ به دو قسمت افراز شده است که قسمت اول شامل یک ضریب $\hat{\beta}_1$ و قسمت دوم شامل بردار $\hat{\beta}_2$ با $(k-1)$ عضو است. دو طرف معادله ۵-۷۶ را از سمت چپ در Λ ضرب می‌کنیم،

$$\Lambda y = \Lambda i\hat{\beta}_1 + \Lambda X_2\hat{\beta}_2 + \Lambda e.$$

با توجه به معادله‌های ۵-۷۳ و ۵-۷۴، معادله فوق را می‌توان چنین نوشت،

$$\Lambda y = \Lambda X_2\hat{\beta}_2 + e. \quad (5.77)$$

دو طرف معادله فوق را در X_2' ضرب می‌کنیم،

$$X_2' \Lambda y = X_2' \Lambda X_2 \hat{\beta}_2 + X_2' e.$$

از معادله ۵-۷۴ می‌دانیم که جمله‌های پسماند از متغیرهای توضیحی مستقل هستند؛ یعنی $X_2' e = 0$. بنابراین:

$$X_2' \Lambda y = X_2' \Lambda X_2 \hat{\beta}_2. \quad (5.78)$$

با توجه به تقارن و خودتوانی ماتریس Λ ، معادله فوق را می‌توان چنین نوشت،

$$(\Lambda X_2)' (\Lambda y) = (\Lambda X_2)' (\Lambda X_2) \hat{\beta}_2. \quad (5.79)$$

معادله ۵-۷۹ در واقع همان معادله‌های نرمال برای مدل رگرسیون چند متغیره ۵-۴ است و

صورت دیگری از معادله ۵-۳۳ می‌باشد. برای توضیح بیشتر معادله ۵-۳۳ را دوباره می‌نویسیم،

$$(X'X)\hat{\beta} = X'y.$$

فرق بین معادله فوق با معادله ۵-۷۹ این است که در معادله ۵-۷۹ ماتریس Λ توانسته است X_1 و y را برحسب انحراف از میانگین بنویسد. بدیهی است معادله ۵-۷۹ را می‌توان برای $\hat{\beta}_1$ به این صورت حل کرد،

$$\hat{\beta}_1 = [(\Lambda X_1)' (\Lambda X_1)]^{-1} (\Lambda X_1)' (\Lambda y). \quad (5-80)$$

معادله فوق مقادیر $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ را در یک مدل رگرسیون چند متغیره محاسبه می‌کند. یادآوری می‌کنیم که در این معادله، $X_1 \rightarrow n \times (k-1)$ و $y \rightarrow n \times 1$ و $\Lambda \rightarrow (k-1) \times 1$. این معادله دقیقاً صورت دیگری از معادله ۵-۳۴ است.

حال به بررسی R^2 می‌پردازیم. دیدیم که Λy بردار مقادیر متغیر درون‌زا و ΛX_1 نیز ماتریس مقادیر $(k-1)$ متغیر برون‌زا برحسب انحراف از میانگین هستند. ابتدا مجموع تغییرات متغیر درون‌زا ($y' y$) را برحسب انحراف از میانگین به دست می‌آوریم. کافی است Λy را ترانهاد نموده و از سمت چپ در Λy ضرب کنیم، خواهیم داشت

$$TSS = y' \Lambda' \Lambda y,$$

که با توجه به تقارن و خودتوانی، ماتریس Λ عبارت خواهد بود از

$$TSS = y' \Lambda y$$

در معادله ۵-۴۵ دیدیم

$$y' y = \hat{y}' \hat{y} + e' e.$$

باید رابطه فوق را برای این حالت نیز ثابت کنیم. کافی است معادله ۵-۷۷ را ترانهاد نموده،

$$y' \Lambda' = \hat{\beta}'_1 X_1' \Lambda' + e',$$

و سپس در معادله ۵-۷۷ ضرب کنیم؛ در نتیجه

$$y' \Lambda' \Lambda y = \hat{\beta}'_1 X'_1 \Lambda' (\Lambda X_1 \hat{\beta}_1) + \hat{\beta}'_1 X'_1 \Lambda' e + e' \Lambda X_1 \hat{\beta}_1 + e' e.$$

با مراجعه به ۵-۷۴ و استفاده از تقارن Λ' ، داریم

$$\Lambda' e = e \quad , \quad e' \Lambda = e,$$

و همچنین با توجه به معادله ۵-۲۴ می دانیم که

$$X_1 e = 0 \quad , \quad e' X = 0.$$

بنابراین جمله‌های دوم و سوم از سمت راست معادله فوق حذف شده، خواهیم داشت

$$y' \Lambda \Lambda y = \hat{\beta}'_1 X'_1 \Lambda' (\Lambda X_1 \hat{\beta}_1) + e' e,$$

که با استفاده از تقارن و خودتوانی Λ به صورت زیر نوشته می شود،

$$y' \Lambda y = \hat{\beta}'_1 X'_1 \Lambda X_1 \hat{\beta}_1 + e' e. \quad (5.81)$$

مفهوم معادله فوق عبارت است از

$$TSS = ESS + RSS.$$

می دانیم R^2 بنا بر تعریف برابر است با $\frac{\text{تغییرات توضیح داده شده (ESS)}}{\text{کل تغییرات (TSS)}}$ و

$$R^2 = \frac{\hat{\beta}'_1 X'_1 \Lambda X_1 \hat{\beta}_1}{y' \Lambda y}, \quad (5.82)$$

معادله ۵-۷۸ را در معادله فوق جایگزین می کنیم،

$$R^2 = \frac{\hat{\beta}'_1 X'_1 \Lambda y}{y' \Lambda y}. \quad (5.83)$$

فرمول فوق، مقدار R^2 را برای حالتی محاسبه می کند که y و X_1 بر حسب مشاهدات اصلی - نه انحراف از میانگین - محاسبه شده است.

روش فوق در استخراج R^2 اهمیت نظری بسیاری دارد؛ زیرا یک روش عمومی در حل این گونه مسائل را نشان می‌دهد. یکی از موارد کاربرد این روش، در فصل ششم و در مسأله ۶-۱۰ مطرح شده است؛ با وجود این، برای حل مسائل R^2 - وقتی محاسبات برحسب مقادیر اصلی انجام شده است - می‌توان از روش بسیار ساده زیر استفاده کرد: فرض می‌کنیم محاسبات برحسب مقادیر اصلی است،

$$y' y = \sum_1^n Y_i^2.$$

می‌دانیم

$$\begin{aligned} y' \Lambda y &= \sum (Y_i - \bar{Y})^2, \\ &= \sum Y_i^2 - n \bar{Y}^2, \\ &= y' y - n \bar{Y}^2. \end{aligned} \quad (5.14)$$

کافی است از $y' y$ که برحسب مقادیر اصلی محاسبه شده است مقدار $n \bar{Y}^2$ را کم کنیم تا مجموع کل تغییرات (TSS) به دست آید. به همین ترتیب اگر \hat{y} را برحسب مشاهدات اصلی حساب کرده باشیم، می‌توان، $n \bar{Y}^2$ را از آن کم کرد تا تغییرات توضیح داده شده (ESS)، یعنی $\hat{\beta}'_1 X_1' \Lambda X_1 \hat{\beta}_1$ ، نتیجه شود؛ زیرا

$$\hat{\beta}'_1 X_1' \Lambda X_1 \hat{\beta}_1 = \sum (\hat{Y}_i - \bar{Y})^2.$$

از معادله ۱-۳۰، می‌دانیم که $\bar{\hat{Y}} = \bar{Y}$ ، بنابراین

$$\begin{aligned} \hat{\beta}'_1 X_1' \Lambda X_1 \hat{\beta}_1 &= \sum \hat{Y}_i^2 + n \bar{Y}^2 - 2 \bar{Y} \sum \hat{Y}_i, \\ &= \sum \hat{Y}_i^2 + n \bar{Y}^2, \\ &= \hat{y}' y - n \bar{Y}^2. \end{aligned} \quad (5.15)$$

با استفاده از معادله‌های ۵-۱۴ و ۵-۱۵ مقدار R^2 به سهولت محاسبه می‌شود:

$$R^2 = \frac{\hat{y}' \hat{y} - n \bar{Y}^2}{y' y - n \bar{Y}^2}. \quad (5.16)$$

فرمولهای دیگری نیز می توان برای R^2 به دست آورد. می دانیم $\hat{y} = X\hat{\beta}$ ، در نتیجه

$$R^2 = \frac{\hat{\beta}' X' X \hat{\beta} - n \bar{Y}^2}{y' y - n \bar{Y}^2} \quad (5.87)$$

با قراردادن $\hat{\beta} = (X' X)^{-1} X' y$ در فرمول 5.87 خواهیم داشت

$$R^2 = \frac{\hat{\beta}' X' y - n \bar{Y}^2}{y' y - n \bar{Y}^2} \quad (5.88)$$

یادآوری می شود که در فرمولهای 5.86، 5.87 و 5.88 مقادیر y و X برحسب مشاهدات اصلی است. همچنین باید توجه داشت که اگر بخواهیم R^2 را از فرمولی به دست آوریم که شامل $e'e$ است، دیگر ضرورتی ندارد که $n\bar{Y}^2$ را از $e'e$ کسر نماییم؛ زیرا $e'e$ همواره به صورت انحراف از میانگین است؛ چون میانگین مقادیر پسماندها صفر است.

مثال 5.4 مدلی را که در مثال 5.2 مطرح شد، یک بار دیگر ملاحظه می کنیم؛

$$Y_i = \alpha_i + \beta X_i + \gamma Z_i + U_i$$

می دانیم بردار y و ماتریس X برابر است با:

$$y = \begin{bmatrix} 3 \\ 1 \\ 8 \\ 3 \\ 5 \end{bmatrix}, \quad X = \begin{bmatrix} 1 & 3 & 5 \\ 1 & 1 & 4 \\ 1 & 5 & 6 \\ 1 & 2 & 4 \\ 1 & 4 & 6 \end{bmatrix}$$

اولاً، با استفاده از ماتریس تبدیل Λ ، مشاهدات را برحسب انحراف از میانگین

نوشته، سپس با به دست آوردن معادله های نرمال، پارامترهای β و γ را تخمین بزنید.

ثانیاً، R^2 را با استفاده از معادله 5.83 محاسبه کنید.

ثالثاً، R^2 را به روش ساده و با استفاده از مقادیر اصلی و به کمک یکی از

فرمولهای 5.86 یا 5.87 یا 5.88 به دست آورید.

۱. ماتریس Λ ، تعریف شده در معادله 5.72 را از سمت چپ در y و X ضرب

می‌کنیم تا مشاهدات متغیر درون‌زا و متغیرهای برون‌زا برحسب انحراف از میانگین، به دست آید،

$$\Delta y = \begin{bmatrix} -1 \\ -3 \\ 4 \\ -1 \\ 1 \end{bmatrix}, \quad \Delta x_T = \begin{bmatrix} 0 & 0 \\ -2 & -1 \\ 2 & 1 \\ -1 & -1 \\ 1 & 1 \end{bmatrix}.$$

با توجه به معادله ۵-۷۹، می‌دانیم معادله‌های نرمال عبارت است از

$$(\Delta x_T)' (\Delta y) = (\Delta x_T)' (\Delta x_T) \hat{\beta}_T,$$

در نتیجه داریم

$$\begin{bmatrix} 0 & -2 & 2 & -1 & 1 \\ 0 & -1 & 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ -3 \\ 4 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 & -2 & 2 & -1 & 1 \\ 0 & -1 & 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ -2 & -1 \\ 2 & 1 \\ -1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix},$$

$$\begin{bmatrix} 16 \\ 9 \end{bmatrix} = \begin{bmatrix} 10 & 6 \\ 6 & 4 \end{bmatrix} \begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix}.$$

از حل این دستگاه معادلات نرمال، مقادیر $\hat{\beta} = 2/5$ و $\hat{\gamma} = -1/5$ به دست می‌آید.
۲. برای محاسبه R^2 از فرمول ۵-۸۳، ابتدا فرمول مزبور را می‌نویسیم،

$$R^2 = \frac{\hat{\beta}'_T X'_T \Delta y}{y' \Delta y}.$$

بردار Δy را در ترانهاد خود ضرب کرده و با توجه به تقارن و خودتوانی Λ داریم

$$y' \Lambda' \Delta y = y' \Delta y = \text{TSS},$$

در نتیجه

$$TSS = \begin{bmatrix} -1 & -3 & 4 & -1 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ -3 \\ 4 \\ -1 \\ 1 \end{bmatrix} = 28.$$

برای به دست آوردن صورت کسر R^2 ، ابتدا $X' \Lambda y$ را محاسبه می‌کنیم،

$$X' \Lambda y = X' \Lambda' (\Lambda y),$$

$$= \begin{bmatrix} 0 & -2 & 2 & -1 & 1 \\ 0 & -1 & 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ -3 \\ 4 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 16 \\ 9 \end{bmatrix}.$$

صورت کسر R^2 ، یعنی تغییرات توضیح داده شده (ESS)، برابر است با

$$ESS = \hat{\beta}' (X' \Lambda y),$$

$$= \begin{bmatrix} 2/5 & -1/5 \end{bmatrix} \begin{bmatrix} 16 \\ 9 \end{bmatrix} = 26/5,$$

در نتیجه

$$RSS = TSS - ESS = 28 - 26/5 = 1/5.$$

بدین ترتیب:

$$R^2 = \frac{26/5}{2} = 0.9476.$$

۳. برای محاسبه R^2 ، به کمک مقادیر اصلی مشاهدات، از معادله ۵-۸۸

استفاده می‌کنیم،

$$R^2 = \frac{\hat{\beta}' X' y - n \bar{Y}^2}{y' y - n \bar{Y}^2}.$$

یادآوری می‌کنیم که در فرمول فوق y و X برحسب مشاهدات اصلی است. در مثال ۵-۱ دیدیم که

$$\bar{y} = 4 \quad \text{و} \quad \hat{\beta}' = [4 \quad 2/5 \quad -1/5] \quad \text{و} \quad X'y = \begin{bmatrix} 20 \\ 76 \\ 109 \end{bmatrix} .$$

همچنین در همین مثال می‌توان نشان داد که $y'y = 108$.

$$\hat{\beta}' X'y = [4 \quad 2/5 \quad -1/5] \begin{bmatrix} 20 \\ 76 \\ 109 \end{bmatrix} = 106/5 ,$$

در نتیجه خواهیم داشت

$$R^2 = \frac{106/5 - 5(4)^2}{108 - 5(4)^2} = \frac{26/5}{28} = 0.9664 ,$$

که دقیقاً برابر همان مقدار R^2 است که با استفاده از مقادیر انحراف از میانگین محاسبه شده است.

مسائل فصل پنجم

۵-۱ مسأله شماره ۴-۱ را یک بار دیگر مطرح می‌کنیم،

$$C_t = \hat{\alpha}_1 + 0.92 Y_t + e_{1t},$$

$$C_t = \hat{\alpha}_2 + 0.84 C_{t-1} + e_{2t},$$

$$C_{t-1} = \hat{\alpha}_3 + 0.78 Y_t + e_{3t},$$

$$C_t = \hat{\alpha}_4 + 0.55 C_{t-1} + e_{4t},$$

که در آن $\hat{\alpha}_i$ ، تخمین جمله ثابت و e_{it} مقادیر پسماند است. با استفاده از اطلاعات فوق، پارامترهای β_1 و β_2 را در مدل رگرسیون زیر فقط از راه ماتریسی تخمین بزنید،

$$C_t = \alpha + \beta_1 Y_t + \beta_2 C_{t-1} + U_t.$$

۵-۲ مدل رگرسیون

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i,$$

مفروض است. با استفاده از یک نمونه، شامل ۲۳ مشاهده، محاسبات زیر را برحسب انحراف از میانگین انجام داده‌ایم،

$$\sum x_i^2 = 12, \quad \sum x_i z_i = 8,$$

$$\sum z_i^2 = 12, \quad \sum y_i x_i = 10,$$

$$\sum y_i^2 = 10, \quad \sum y_i z_i = 8.$$

اولاً، پارامترهای β و γ را تخمین بزنید.
ثانیاً، R^2 و \bar{R}^2 را محاسبه کنید.

۵-۳ مدل زیر را ملاحظه کنید.

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i.$$

با استفاده از یک نمونه شامل ۱۰ مشاهده روی Y_i ، X_{1i} و X_{2i} ، محاسبات زیر را انجام داده ایم،

$$\begin{aligned} \sum Y_i &= 20, & \sum Y_i^2 &= 88/2, & \sum X_{1i} Y_i &= 59, \\ \sum X_{1i} &= 30, & \sum X_{1i}^2 &= 92, & \sum X_{2i} Y_i &= 88, \\ \sum X_{2i} &= 40, & \sum X_{2i}^2 &= 163, & \sum X_{1i} X_{2i} &= 119. \end{aligned}$$

اولاً، پارامترهای α ، β_1 و β_2 را تخمین بزنید.
ثانیاً، R^2 و \bar{R}^2 را به دست آورید.
ثالثاً، مجموع مربعات پسماند را حساب کنید.

۵-۴ مدل

$$\ln Y_i = \beta_1 + \beta_2 \ln X_{2i} + \beta_3 \ln X_{3i} + U_i,$$

مفروض است. نشان دهید که تخمین پارامترهای این مدل $(\hat{\beta}_2, \hat{\beta}_3)$ به ترتیب برابر با کشش Y_i نسبت به X_{2i} و X_{3i} است. آیا این کششها ثابت هستند؟
۵-۵ مدل رگرسیون زیر را مشاهده کنید،

$$Y_i = \beta_1 + \beta_2 X_{1i} + \beta_3 X_{2i} + \beta_4 (X_{2i} - X_{3i}) + \beta_5 X_{0i} + U_i.$$

به نظر شما کدامیک از پارامترهای مدل فوق را می توان تخمین زد؟ توضیح دهید.

۵-۶ فرض کنید تقاضا برای کالای a ، تابعی از قیمت آن کالا (P_{ai}) ، قیمت k کالاهای دیگر $(P_{1i}, P_{2i}, \dots, P_{ki})$ ، سطح عمومی قیمتها (\bar{P}_i) و درآمد (Y_i) است؛ بنابراین

$$D_{ai} = \alpha_i + \beta_1 P_{ai} + \beta_2 P_{1i} + \beta_3 P_{2i} + \dots + \beta_k P_{ki} + \gamma \bar{P}_i + \lambda Y_i + U_i,$$

که در آن D_{ai} ، مقدار تقاضا شده برای کالای a در دوره t است. فرض کنید سطح

عمومی قیمت‌ها به صورت بسیار ابتدایی چنین تعریف شده است،

$$\bar{P}_i = \frac{\sum_{k=1}^m P_{ik}}{M}$$

که در آن $i = 1, 2, \dots, k, k+1, k+2, \dots, m$. به نظر شما چه مشکل اساسی در تخمین این مدل وجود دارد؟ چرا نمی‌توان تمام پارامترهای آن را تخمین زد؟
۵-۷ به نظر شما در مدل رگرسیون

$$Y_i = \alpha_0 + \alpha_1 X_{1i} + \alpha_2 X_{2i} + \alpha_3 (X_{3i} - X_{2i}) + \alpha_4 X_{1i} X_{2i} + U_i$$

کدامیک از پارامترهای مدل را نمی‌توان تخمین زد؟
۵-۸ مدل رگرسیون زیر مفروض است:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + U_i$$

اولاً، آیا تمام پارامترهای این مدل را می‌توان با روش حداقل مربعات معمولی تخمین زد؟
ثانیاً، با استفاده از مشاهدات زیر

Y_i	-۱	-۱	۲	۴	۵
X_i	۰	۱	۲	۵	۶

معادله‌های نرمال را بنویسید.

۵-۹ مدل رگرسیون

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + U_i$$

مفروض است. این مدل را تخمین زده‌ایم و \bar{R}^2 برابر ۰/۸۶ به دست آمده است. متغیر توضیحی X_{2i} را به مدل اضافه کرده و دوباره آن را تخمین می‌زنیم. این باره \bar{R}^2 برابر ۰/۸۲ محاسبه شده است. نشان دهید که حتماً اشتباهی در محاسبات صورت گرفته است، مگر اینکه حجم نمونه ۸ یا کمتر از آن باشد.

۵-۱۰ می‌خواهیم مدل رگرسیون زیر را تخمین بزنیم،

$$Y_t = \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + U_t, \quad (1)$$

که در آن جمله ثابت وجود ندارد. فرض کنید کامپیوتری که در اختیار داریم، چنان برنامه‌ریزی شده است که هر مجموعه مشاهداتی که به آن بدهیم، یک جمله ثابت را نیز در نظر می‌گیرد؛ یعنی کامپیوتر همواره مدل زیر را تخمین می‌زند،

$$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + \dots + \beta_k X_{kt} + U_t. \quad (2)$$

نشان دهید که اگر به ازای هر زوج مشاهده که به کامپیوتر می‌دهیم، یعنی (Y_t, X_{it}) ، $i = 2, 3, \dots, k$ و $t = 1, 2, \dots, n$ یک زوج مشاهده دیگر با علامت مخالف، یعنی $(-Y_t, -X_{it})$ ، نیز وارد کامپیوتر کنیم، آنگاه

اولاً، تخمینهای به دست آمده از معادله (۲) دقیقاً برابر تخمینهای مطلوب ما برای معادله (۱) خواهد بود.

ثانیاً، مجموع مربعات پسماند برای مدل (۲) دقیقاً دو برابر مجموع مربعات پسماند برای مدل (۱) است.

ثالثاً، نسبت انحراف معیار تخمین (SEE) در مدل (۲) نسبت به مدل (۱) برابر است با:

$$\frac{SEE(2)}{SEE(1)} = \frac{\sqrt{2(n-k+1)}}{\sqrt{(2n-k)}}$$

حل مسائل فصل پنجم

۵-۱ معادله ۵-۲۰ یعنی $(X'X)\beta = X'Y$ را - که همان سیستم معادله‌های نرمال است - برای این مسأله می‌نویسیم،

$$\begin{bmatrix} \sum y_t^2 & \sum y_t c_{t-1} \\ \sum c_{t-1} y_t & \sum c_{t-1}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} \sum c_t y_t \\ \sum c_{t-1} c_t \end{bmatrix}$$

با استفاده از اطلاعات موجود می‌دانیم

$$\frac{\sum c_t y_t}{\sum y_t^2} = 0/92 \rightarrow \sum c_t y_t = 0/92 \sum y_t^2, \quad (1)$$

$$\frac{\sum c_t c_{t-1}}{\sum c_{t-1}^2} = 0/84 \rightarrow \sum c_t c_{t-1} = 0/84 \sum c_{t-1}^2, \quad (2)$$

$$\frac{\sum c_{t-1} y_t}{\sum y_t^2} = 0/78 \rightarrow \sum c_{t-1} y_t = 0/78 \sum y_t^2, \quad (3)$$

$$\frac{\sum y_t c_{t-1}}{\sum c_{t-1}^2} = 0/00 \rightarrow \sum y_t c_{t-1} = 0/00 \sum c_{t-1}^2. \quad (4)$$

از معادله (۴) داریم

$$\sum y_t c_{t-1} = \frac{\sum y_t c_{t-1}}{0/00} \cdot 0/00$$

معادله (۳) را در معادله فوق قرار می‌دهیم،

$$\sum c_{t-1}^2 = \frac{0/78 \sum y_t^2}{0/00}. \quad (5)$$

بدین ترتیب، تمام عناصر ماتریس $X'X$ را برحسب $\sum y_t^2$ به دست آورده‌ایم. با معادله‌های

(۳) و (۵) این ماتریس را می‌سازیم،

$$\begin{bmatrix} \sum y_i^2 & \sum y_i c_{i-1} \\ \sum c_{i-1} y_i & \sum c_{i-1}^2 \end{bmatrix} = \begin{bmatrix} \sum y_i^2 & \cdot/\sqrt{8} \sum y_i^2 \\ \cdot/\sqrt{8} \sum y_i^2 & \frac{\cdot/\sqrt{8} \sum y_i^2}{\cdot/50} \end{bmatrix}$$

به ترتیبی مشابه، بردار $X'Y$ را برحسب $\sum y_i^2$ می‌نویسیم. در معادله (۱) عنصر اول این بردار را برحسب $\sum y_i^2$ به دست آورده‌ایم. باید عنصر دوم آن، یعنی $\sum c_{i-1} y_i$ را برحسب $\sum y_i^2$ بنویسیم. کافی است معادله (۵) را در (۲) قرار دهیم،

$$\sum c_{i-1} y_i = \cdot/84 \left(\frac{\cdot/\sqrt{8} \sum y_i^2}{\cdot/50} \right) = 1/191 \sum y_i^2$$

بنابراین ماتریس $X'Y$ برحسب $\sum y_i^2$ به ترتیب زیر خواهد بود،

$$\begin{bmatrix} \sum c_{i-1} y_i \\ \sum c_{i-1}^2 \end{bmatrix} = \begin{bmatrix} \cdot/92 \sum y_i^2 \\ 1/191 \sum y_i^2 \end{bmatrix}$$

معادله‌های نرمال ۵-۲۰ را نوشته، بعد از حذف $\sum y_i^2$ خواهیم داشت

$$\begin{bmatrix} 1 & \cdot/\sqrt{8} \\ \cdot/\sqrt{8} & \frac{\cdot/\sqrt{8}}{\cdot/50} \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} \cdot/92 \\ 1/191 \end{bmatrix}$$

که با حل این دستگاه در معادله دو مجهولی داریم $\hat{\beta}_1 = \cdot/585$ و $\hat{\beta}_2 = \cdot/464$

۵-۲ الف) ابتدا ماتریس $(X'X)$ را تشکیل می‌دهیم،

$$(X'X) = \begin{bmatrix} \sum x_i^2 & \sum x_i z_i \\ \sum z_i x_i & \sum z_i^2 \end{bmatrix} = \begin{bmatrix} 12 & 8 \\ 8 & 12 \end{bmatrix}$$

$(X'X)$ را معکوس می‌کنیم. داریم

$$(X'X)^{-1} = \begin{bmatrix} 0/10 & -0/1 \\ -0/1 & 0/10 \end{bmatrix}$$

بردار $X'y$ را برای این مسأله محاسبه می‌کنیم،

$$X'y = \begin{bmatrix} \sum x_i y_i \\ \sum z_i y_i \end{bmatrix} = \begin{bmatrix} 10 \\ 8 \end{bmatrix}$$

بردار $\hat{\beta}$ به راحتی قابل محاسبه است،

$$\begin{aligned} \hat{\beta} = \begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \end{bmatrix} &= \begin{bmatrix} \sum x_i^2 & \sum x_i z_i \\ \sum z_i x_i & \sum z_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum x_i y_i \\ \sum z_i y_i \end{bmatrix} \\ &= \begin{bmatrix} 0/10 & -0/1 \\ -0/1 & 0/10 \end{bmatrix} \begin{bmatrix} 10 \\ 8 \end{bmatrix} = \begin{bmatrix} 0/7 \\ 0/2 \end{bmatrix} \end{aligned}$$

بنابراین $\hat{\beta} = 0/7$ و $\hat{\gamma} = 0/2$

(ب) می‌دانیم

$$R^2 = \frac{\hat{\beta}' X' y}{y' y}$$

صورت کسر R^2 برای این مسأله عبارت است از

$$\begin{aligned} \hat{\beta}' X' y &= [\hat{\beta} \quad \hat{\gamma}] \begin{bmatrix} \sum x_i y_i \\ \sum z_i y_i \end{bmatrix} = \hat{\beta} \sum x_i y_i + \hat{\gamma} \sum z_i y_i \\ &= 0/7 (10) + 0/2 (8) = 8/6 \end{aligned}$$

مخرج کسر R^2 برای این مسأله برابر است با

$$y'y = \sum y_i^2 = 10$$

در نتیجه R^2 به راحتی محاسبه می‌شود،

$$R^2 = \frac{8/7}{10} = 0/86$$

برای محاسبه \bar{R}^2 از فرمول زیر استفاده می‌کنیم،

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k}\right) (1-R^2),$$

$$= 1 - \left(\frac{23-1}{23-3}\right) (1-0/86) = 0/846.$$

۵-۳ الف) با توجه به اینکه کمیت‌های محاسبه شده برحسب مشاهدات اصلی است، باید آنها را به انحراف از میانگین تبدیل نماییم. کمیت‌های لازم را به شرح زیر محاسبه می‌کنیم.

$$\bar{y} = 2 \quad , \quad \bar{x}_1 = 3 \quad , \quad \bar{x}_2 = 4.$$

$$\sum x_{1i}^2 = \sum X_{1i}^2 - n\bar{x}_1^2 = 92 - 10(3)^2 = 2,$$

$$\sum x_{2i}^2 = \sum X_{2i}^2 - n\bar{x}_2^2 = 163 - 10(4)^2 = 3,$$

$$\sum x_{1i} x_{2i} = \sum X_{1i} X_{2i} - n\bar{x}_1 \bar{x}_2 = 119 - 10(3)(4) = -1,$$

$$\sum x_{1i} y_i = \sum X_{1i} Y_i - n\bar{x}_1 \bar{y} = 59 - 10(3)(2) = -1,$$

$$\sum x_{2i} y_i = \sum X_{2i} Y_i - n\bar{x}_2 \bar{y} = 88 - 10(4)(2) = 8,$$

$$\sum y_i^2 = \sum Y_i^2 - n\bar{y}^2 = 88/2 - 10(2)^2 = 48/2.$$

ماتریس $(X'X)$ و بردار $X'y$ را محاسبه می‌کنیم. خواهیم داشت

$$(X'X) = \begin{bmatrix} \sum x_{1i}^2 & \sum x_{1i} x_{2i} \\ \sum x_{2i} x_{1i} & \sum x_{2i}^2 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix},$$

$$(\mathbf{X}'\mathbf{y}) = \begin{bmatrix} \sum x_{1t} y_t \\ \sum x_{2t} y_t \end{bmatrix} = \begin{bmatrix} -1 \\ 8 \end{bmatrix} .$$

$\hat{\beta}$ به صورت زیر محاسبه می شود،

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \begin{bmatrix} 2 & -1 \\ -1 & 3 \end{bmatrix}^{-1} \begin{bmatrix} -1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} ,$$

یعنی $\hat{\beta}_1 = 1$ و $\hat{\beta}_2 = 3$. برای تخمین α به صورت زیر عمل می کنیم،

$$\bar{y} = \hat{\alpha} + \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2 ,$$

$$2 = \hat{\alpha} + (1)(3) + (3)(4) ,$$

$$\hat{\alpha} = -13$$

بدین ترتیب تخمین مدل عبارت است از

$$\hat{Y}_t = -13 + X_{1t} + 3X_{2t} .$$

(ب) برای محاسبه R^2 ، ابتدا فرمول ۵-۴۰ را می نویسیم،

$$R^2 = \frac{\hat{\mathbf{y}}'\hat{\mathbf{y}}}{\mathbf{y}'\mathbf{y}} = \frac{\hat{\beta}'\mathbf{X}'\mathbf{y}}{\mathbf{y}'\mathbf{y}} .$$

صورت کسر R^2 برای این مسأله عبارت است از

$$\hat{\beta}'\mathbf{X}'\mathbf{y} = [\hat{\beta}_1 \quad \hat{\beta}_2] \begin{bmatrix} \sum x_{1t} y_t \\ \sum x_{2t} y_t \end{bmatrix} ,$$

$$= [1 \quad 3] \begin{bmatrix} -1 \\ 8 \end{bmatrix} = 23 .$$

مخرج کسر R^2 برابر است با

$$\mathbf{y}'\mathbf{y} = \sum y_i^2 = ۴۸/۲,$$

بدین ترتیب داریم

$$R^2 = \frac{۲۳}{۴۸/۲} = ۰/۴۷۷۱.$$

می‌دانیم

$$\begin{aligned} \bar{R}^2 &= ۱ - \left(\frac{n-1}{n-k}\right) (1-R^2), \\ &= ۱ - \left(\frac{۱۰-1}{۱۰-۳}\right) (1-۰/۴۷۷۱), \\ &= ۱ - ۰/۶۷۲۳ = ۰/۳۲۷۷۷. \end{aligned}$$

ج) برای به دست آوردن مجموع مربعات پسماند $(\sum e_i^2)$ از فرمول ۵-۴۸ استفاده می‌کنیم.

$$R^2 = ۱ - \frac{\mathbf{e}'\mathbf{e}}{\mathbf{y}'\mathbf{y}},$$

یا

$$\begin{aligned} \mathbf{e}'\mathbf{e} &= (1-R^2)\mathbf{y}'\mathbf{y}, \\ &= (1-۰/۴۷۷۱)(۴۸/۲) = ۲۵/۲۰. \end{aligned}$$

۵-۴ می‌دانیم:

$$\beta_{\gamma} = \frac{\partial \ln Y_i}{\partial \ln X_{\gamma i}} = \frac{\left(\frac{1}{Y_i}\right) d Y_i}{\left(\frac{1}{X_{\gamma i}}\right) d X_{\gamma i}} = \frac{d Y_i}{d X_{\gamma i}} \cdot \frac{X_{\gamma i}}{Y_i}.$$

$$\beta_{\gamma} = \frac{\partial \ln Y_i}{\partial \ln X_{\gamma i}} = \frac{\left(\frac{1}{Y_i}\right) d Y_i}{\left(\frac{1}{X_{\gamma i}}\right) d X_{\gamma i}} = \frac{d Y_i}{d X_{\gamma i}} \cdot \frac{X_{\gamma i}}{Y_i}.$$

برای اثبات رابطه‌های فوق به توضیحات موجود در پاراگرافی مبحث مدل‌های

خطی - لگاریتمی، مندرج در فصل سوم قسمت ۳-۶ مراجعه کنید. بنابراین، $\hat{\beta}_3$ و $\hat{\beta}_4$ به ترتیب تخمین کشش Y_t نسبت به X_{t1} و X_{t2} است. همچنین این کششها ثابت نیز هستند زیرا می‌دانیم تخمین پارامترهای یک مدل رگرسیون با روش حداقل مربعات معمولی مقادیر ثابتی را نتیجه می‌دهد.

۵-۵ مدل را به صورت زیر می‌نویسیم،

$$Y_t = \beta_1 + (\beta_2 + \beta_3) X_{t2} + (\beta_2 - \beta_3) X_{t1} + \beta_0 X_{0t} + U_t.$$

با داشتن مشاهدات X_{0t} ، X_{t1} ، X_{t2} ، Y_t می‌توان مدل را به صورت زیر تخمین زد،

$$\hat{Y}_t = \hat{\beta}_1 + \hat{\gamma} X_{t2} + \hat{\lambda} X_{t1} + \hat{\beta}_0 X_{0t}.$$

ملاحظه می‌شود برای پارامترهای β_0 و β_1 می‌توان به تخمینهای مستقلی رسید، اما نمی‌توان $\hat{\beta}_2$ ، $\hat{\beta}_3$ و $\hat{\beta}_4$ را به دست آورد. آنچه می‌توان تخمین زد، عبارت است از

$$\hat{\gamma} = \hat{\beta}_2 + \hat{\beta}_3, \quad \hat{\lambda} = \hat{\beta}_2 - \hat{\beta}_3.$$

بدیهی است که دو معادله سه مجهولی داریم. اگر بتوان یک رابطه دیگر بین $\hat{\beta}_2$ و $\hat{\beta}_3$ به دست آورد، تخمینهای مستقل از این پارامترها ممکن خواهد بود. در این گونه موارد معمولاً باید از نظریه‌های اقتصادی کمک گرفت تا بتوان به چنین رابطه‌هایی رسید.

۵-۶ با توجه به صورت مسأله می‌دانیم

$$\bar{P}_t = \frac{1}{k} [P_{at} + P_{1t} + P_{2t} + \dots + P_{kt} + P_{(k-1)t} + \dots + P_{mt}].$$

بنابراین \bar{P}_t با متغیرهای P_{at} ، P_{1t} ، \dots ، P_{kt} همخطی کامل خواهد داشت و به همین دلیل نمی‌توان پارامترهای متغیرهایی را که همخطی دارند تخمین زد. برای تبیین بیشتر این نکته می‌توانیم \bar{P}_t را در مدل مفروض جایگزین کنیم. بعد از ساده کردن خواهیم داشت

$$D_{at} = \alpha_0 + \left(\beta_0 + \frac{\gamma}{M}\right) P_{at} + \left(\beta_1 + \frac{\gamma}{M}\right) P_{1t} + \dots + \left(\beta_k + \frac{\gamma}{M} P_{kt}\right) \\ + \gamma^* \bar{P}_t + \lambda Y_t + U_t,$$

که در آن

$$\bar{P}_i^* = \frac{1}{k} [P_{(k+1)t} + P_{(k+2)t} + \dots + P_{mt}] .$$

ملاحظه می شود که فقط می توان α_1 ، λ و γ را تخمین زد. البته γ مطلوب مانیست؛ بنابراین از $k+1$ پارامتری که در مدل اولیه وجود دارد، فقط میتوان به تخمین دو پارامتر رسید.

۵.۷ پارامترهایی را که نمی توان با روش حداقل مربعات معمولی تخمین زد، عبارت است از α_1 ، α_2 و α_3 ، زیرا سومین متغیر توضیحی مدل مفروض، یعنی $(X_{1t} - X_{2t})$ یک ترکیب خطی از دو متغیر اول و دوم، یعنی X_{1t} و X_{2t} است؛ به عبارت دیگر، بین این سه متغیر توضیحی همخطی کامل برقرار است. مدل مفروض را می توان به صورت زیر نوشت

$$Y_t = \alpha_0 + (\alpha_1 + \alpha_2) X_{1t} + (\alpha_2 - \alpha_3) X_{2t} + \alpha_4 X_{1t} X_{2t} + U_t .$$

بنابراین فقط می توان به تخمینهای مستقلی از α_1 و α_2 رسید. توجه داریم که چهارمین متغیر توضیحی مدل مفروض، یعنی $(X_{1t} X_{2t})$ ، یک ترکیب «غیرخطی» از متغیرهای اول و دوم است و به همین دلیل مشکلی را در امر تخمین ایجاد نکرده است و α_1 را می توان به راحتی تخمین زد.

۵.۸ الف) هیچ مشکل نظری در تخمین پارامترهای این مدل وجود ندارد؛ زیرا X_1 و X_1^2 یک رابطه غیرخطی با یکدیگر دارند.

ب) می دانیم برای نوشتن معادله های نرمال، باید از $\sum e_i^2$ نسبت به β_0 ، β_1 و β_2 مشتق گرفته و آنها را مساوی صفر قرار دهیم. اما قبلاً در فصل چهارم و در ادامه بحث معادله های نرمال ۵-۴ دیدیم که یک روش ساده برای نوشتن معادله های نرمال وجود دارد. با استفاده از این روش داریم

$$\sum Y_i = n \hat{\beta}_0 + \hat{\beta}_1 \sum X_i + \hat{\beta}_2 \sum X_i^2 .$$

$$\sum X_i Y_i = \hat{\beta}_0 \sum X_i + \hat{\beta}_1 \sum X_i^2 + \hat{\beta}_2 \sum X_i^3 ,$$

$$\sum X_i^2 Y_i = \hat{\beta}_0 \sum X_i^2 + \hat{\beta}_1 \sum X_i^3 + \hat{\beta}_2 \sum X_i^4 ,$$

با استفاده از مشاهدات داده شده، می توان معادله های نرمال فوق را به صورت زیر نوشت،

$$۹ = ۵\hat{\beta}_0 + ۱۴\hat{\beta}_1 + ۶۶\hat{\beta}_2,$$

$$۵۳ = ۱۴\hat{\beta}_0 + ۶۶\hat{\beta}_1 + ۳۵۰\hat{\beta}_2,$$

$$۲۸۷ = ۶۶\hat{\beta}_0 + ۳۵۰\hat{\beta}_1 + ۱۹۲۸\hat{\beta}_2.$$

۵.۹ می دانیم اگر آماره t ، متعلق به تخمین یک پارامتر از یک کمتر باشد، ورود آن متغیر توضیحی در مدل می تواند \bar{R}^2 را کاهش دهد. در این مسأله ملاحظه می شود که بعد از اضافه کردن X_{4t} مقدار \bar{R}^2 کاهش یافته است؛ بنابراین اگر t مربوط به X_{4t} کمتر از یک باشد، این نتیجه قابل قبول است. اما با اطلاعات بسیار محدودی که در مسأله داده شده است، نمی توان این نکته را ارزیابی کرد که واقعاً t از یک کمتر است. باید به دنبال معیار دیگری برویم. مسلماً \bar{R}^2 در مدل جدید کاهش یافته، بنابراین سؤال این است که چه شرطی لازمه تحقق چنین نتیجه ای است؟ با استفاده از معادله ۵-۶۳ یعنی

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k}\right) (1 - R^2),$$

یا مستقیماً از معادله ۵-۶۶ می دانیم

$$\frac{1 - \bar{R}^2}{n-1} = \frac{1 - R^2}{n-k},$$

یا

$$(1 - R^2) = \frac{(n-k)}{(n-1)} (1 - \bar{R}^2). \quad (۱)$$

واضح است که با ورود یک متغیر توضیحی جدید به مدل رگرسیون، مقدار R^2 اضافه شده یا حداقل ثابت می ماند؛ بنابراین باید انتظار داشته باشیم که در مدل رگرسیون جدید، یعنی بعد از اضافه کردن X_{4t} ، مقدار $(1 - R^2)$ از حالت قبل کمتر بشود. در نتیجه، همواره باید رابطه زیر برقرار باشد،

$$(1 - R^2) \leq (1 - \bar{R}^2) \text{ در مدل قدیم}. \quad (۲)$$

رابطه (۱) را در (۲) قرار می‌دهیم،

$$\frac{(n-k)}{(n-1)} (1-\bar{R}^1) \leq \frac{(n-k)}{(n-1)} (1-\bar{R}^2),$$

و یا حذف (n-1) خواهیم داشت

$$(n-k) (1-\bar{R}^1) \leq (n-k) (1-\bar{R}^2). \quad (3)$$

می‌دانیم تعداد پارامترهایی که تخمین زده می‌شود، یعنی k، در مدل‌های جدید و قدیم به ترتیب برابر با ۵ و ۴ است. برای اینکه محاسبات ارائه شده در مسأله صحیح باشد، باید در رابطه (۳) صدق کند؛ بنابراین باید داشته باشیم

$$(n-5) (1-0/82) \leq (n-4) (1-0/86),$$

$$0/18(n-5) \leq 0/14(n-4),$$

$$n \leq 1/0.$$

یعنی شرط صحت محاسبات ارائه شده این است که حجم نمونه ۸ یا کمتر از آن باشد.
 ۵-۱۰ الف) وقتی به ازای هر زوج مشاهده (Y_i, X_{ij}) یک زوج مشاهده با علامت مخالف، یعنی $(-Y_i, -X_{ij})$ وارد کامپیوتر کنیم، ماتریس X و بردار Y به صورت زیر خواهد بود،

$$y^* = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \\ -Y_1 \\ -Y_2 \\ \vdots \\ -Y_n \end{bmatrix}, \quad X^* = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \\ 1 & -X_{11} & -X_{21} & \dots & -X_{k1} \\ 1 & -X_{12} & -X_{22} & \dots & -X_{k2} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & -X_{1n} & -X_{2n} & \dots & -X_{kn} \end{bmatrix}$$

که $\mathbf{y}^* \rightarrow 2n \times 1$ و $\mathbf{X}^* \rightarrow 2n \times k$. در واقع مدل (۲) با \mathbf{y}^* و \mathbf{X}^* تخمین زده می شود. باید نشان دهیم مقادیر $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ که با استفاده از \mathbf{y}^* و \mathbf{X}^* به دست می آید دقیقاً برابر تخمینهایی است که از مدل (۱) و به کمک \mathbf{y} و \mathbf{X} حاصل می شود. برای این منظور $\mathbf{X}^* \mathbf{X}^*$ و $\mathbf{y}^* \mathbf{y}^*$ را تشکیل می دهیم. ابتدا \mathbf{y}^* و \mathbf{X}^* را به صورت زیر تعریف می کنیم،

$$\mathbf{y}^* = \begin{bmatrix} \mathbf{y} \\ -\mathbf{y} \end{bmatrix}, \quad \mathbf{X}^* = \begin{bmatrix} \mathbf{I} & \mathbf{X} \\ \mathbf{I} & -\mathbf{X} \end{bmatrix},$$

که در آن \mathbf{I} یک بردار ستونی $(n \times 1)$ و شامل عدد یک است و $\mathbf{X} \rightarrow n \times (k-1)$ ؛ بنابراین

$$\mathbf{X}^* \mathbf{X}^* = \begin{bmatrix} \mathbf{I}' & \mathbf{I} \\ \mathbf{X}' & -\mathbf{X}' \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{X} \\ \mathbf{I} & -\mathbf{X} \end{bmatrix} = \begin{bmatrix} 2n & \mathbf{0} \\ \mathbf{0} & 2\mathbf{X}'\mathbf{X} \end{bmatrix},$$

که در آن n ، تعداد مشاهدات اولیه در مدل (۱) است. به همین ترتیب می توان نوشت

$$\mathbf{X}^* \mathbf{y}^* = \begin{bmatrix} \mathbf{I}' & \mathbf{I} \\ \mathbf{X}' & -\mathbf{X}' \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ -\mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 2\mathbf{X}'\mathbf{y} \end{bmatrix}.$$

اگر تخمین $\hat{\beta}$ را در مدل (۲)، $\hat{\beta}^*$ بنامیم، خواهیم داشت

$$\begin{aligned} \hat{\beta}^* &= (\mathbf{X}^* \mathbf{X}^*)^{-1} (\mathbf{X}^* \mathbf{y}^*), \\ &= \begin{bmatrix} 2n & \mathbf{0} \\ \mathbf{0} & 2\mathbf{X}'\mathbf{X} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{0} \\ 2\mathbf{X}'\mathbf{y} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \frac{1}{n} & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}'\mathbf{X})^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ 2\mathbf{X}'\mathbf{y} \end{bmatrix}. \end{aligned}$$

بنابراین داریم:

$$\hat{\beta}^* = \begin{bmatrix} \mathbf{0} \\ (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \end{bmatrix} \quad (3)$$

اما اگر خواهیم مدل (۱) را بدون جمله ثابت تخمین بزنیم، باید از فرمول زیر استفاده کنیم،

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \quad (4)$$

ملاحظه می‌شود که $\hat{\beta}^*$ در معادله (۳) برای جمله ثابت مدل (۲) تخمینی برابر صفر و برای پارامترهای $\beta_1, \beta_2, \beta_3$ تا β_k تخمینی برابر $X'X^{-1}X'y$ نتیجه می‌دهد که دقیقاً برابر $\hat{\beta}$ در معادله (۴) است. بدین ترتیب توانستیم از مدلی که جمله ثابت دارد، تخمین پارامترهای مدل فاقد جمله ثابت را، با استفاده از مشاهدات اصلی، به دست آوریم.

$$\hat{\beta}^* = \begin{bmatrix} 0 \\ \hat{\beta} \end{bmatrix} \quad (5)$$

ب) بردار جمله‌های پسماند برای مدل (۲) عبارت است از

$$e^* = y^* - \hat{y}^* = y^* - X^* \hat{\beta}^* .$$

با استفاده از معادله (۵)، داریم:

$$\begin{aligned} e^* &= \begin{bmatrix} y \\ -y \end{bmatrix} - \begin{bmatrix} 1 & X \\ 1 & -X \end{bmatrix} \begin{bmatrix} 0 \\ \hat{\beta} \end{bmatrix} \\ &= \begin{bmatrix} (y - X\hat{\beta}) \\ -(y - X\hat{\beta}) \end{bmatrix} = \begin{bmatrix} e \\ -e \end{bmatrix} . \end{aligned}$$

با داشتن e^* می‌توان مجموع مربعات پسماند برای مدل (۲) را حساب کرد:

$$e^{*'} e^* = [e' \ -e'] \begin{bmatrix} e \\ -e \end{bmatrix} = e'e + e'e = 2e'e, \quad (6)$$

که در آن $e'e$ ، مجموع مربعات پسماند برای مدل (۱) است.

ج) با توجه به معادله ۱-۴۵ می‌دانیم که انحراف معیار تخمین در یک مدل رگرسیون ساده شامل دو پارامتر، عبارت است از

$$SEE = \sqrt{\frac{\sum e_i^2}{n-2}} .$$

درجات آزادی برای $\sum e_i^2$ در حالتی که k پارامتر داریم، برابر است با $(n-k)$ ، بنابراین خطای معیار تخمین برای یک رگرسیون چند متغیره با k پارامتر عبارت است از

$$SEE = \sqrt{\frac{e'e}{n-k}}$$

در مدل (۱) مقدار پارامترها $(k-1)$ است؛ بنابراین

$$SEE (1) = \sqrt{\frac{e'e}{n-(k-1)}}, \quad (7)$$

و برای مدل (۲) داریم

$$SEE (2) = \sqrt{\frac{e^*e^*}{\gamma n-k}}$$

با استفاده از معادله (۶)، می توان خطای معیار تخمین (۲) (SEE) را به صورت زیر نوشت،

$$SEE (2) = \sqrt{\frac{\gamma e'e}{\gamma n-k}}, \quad (8)$$

در نتیجه، خواهیم داشت

$$\frac{SEE (2)}{SEE (1)} = \sqrt{\frac{\gamma e'e (n-k+1)}{(\gamma n-k) e'e}} = \frac{\sqrt{\gamma (n-k+1)}}{\sqrt{\gamma n-k}}$$

فصل ششم

آزمون پارامترهای مدل رگرسیون خطی چندمتغیره

۶-۱ مقدمه

بعد از تخمین پارامترهای یک مدل رگرسیون چندمتغیره در فصل پنجم، در این فصل به آزمون این پارامترها می پردازیم. بردار پارامترهای β شامل k پارامتر را در نظر می گیریم. β را با روش حداقل مربعات معمولی و به عنوان تخمین بردار β به دست آورده ایم. سؤال این است که چگونه به کمک $\hat{\beta}$ می توان درباره β استنباط آماری داشت؟

با توجه به اینکه β یک بردار تصادفی است باید ابتدا تابع توزیع احتمال و میانگین و واریانس آن را مطالعه کنیم. این نکات در قسمت ۶-۲ با نام خصوصیات آماری β بررسی شده است. در این قسمت قضیه گاس - مارکف به طور خلاصه و نیز خصوصیت «بهترین تخمین زننده نااریب خطی» برای $\hat{\beta}$ مطرح می شود. همچنین نشان داده ایم که اگر مجموع مربعات پسماند را بر درجات آزادی آن تقسیم کنیم، یک تخمین نااریب از واریانس جمله اختلال خواهیم داشت. آزمون واریانس جمله اختلال نیز در همین قسمت مطرح شده است.

با شناخت خصوصیات آماری $\hat{\beta}$ به بررسی آزمون β می پردازیم. قسمت ۶-۳ آزمون هر یک از پارامترهاست. آزمونهای $\beta_1 = a$ که a عدد ثابت و معلومی است و نیز آزمون $\beta_1 = 0$ در این قسمت ارزیابی خواهد شد. فرضیه صفر بودن همزمان زیرمجموعه ای از پارامترهای یک مدل رگرسیون چندمتغیره اهمیت فراوان دارد و کاربردهای متعددی در اقتصادسنجی کاربردی دارد. نکات آزمون این فرضیه که $\beta_1 = \beta_{1+1} = \dots = \beta_j = 0$ در همین قسمت مطرح می شود. آزمون معنی دار بودن کل مدل

رگرسیون که در واقع آزمون اعتبار مدل است در قسمت ۶-۴ بررسی می‌شود. در این قسمت فرضیه $\beta_1 = \beta_2 = \dots = \beta_k = 0$ را به کمک آنالیز واریانس، آزمون خواهیم کرد. ملاحظه می‌شود که در صورت صحت فرضیه H_0 ، مدل رگرسیون چندمتغیره مفروض، به مدل $Y_i = \beta + U_i$ تبدیل خواهد شد و بدین معنی است که تغییرات Y_i کاملاً تصادفی بوده و تنها تابعی از U_i است. این سؤال مهم معمولاً در بحث آزمونها مطرح می‌شود: آیا بین آزمونهای معنی‌دار بودن هریک از پارامترها و آزمون همزمان آنها تفاوت وجود دارد؟ به عبارت دیگر، اگر برای هریک از مقادیر $k, \dots, 3, 2 =$ فرضیه‌های $\beta_k = 0$ را به طور جداگانه انجام دهیم و به این نتیجه برسیم که تمام آنها قبول می‌شود، آیا ضرورتاً فرضیه $\beta_1 = \beta_2 = \dots = \beta_k = 0$ نیز قبول خواهد شد. در همین قسمت نشان خواهیم داد که ضرورتاً این گونه نیست.

آزمون یک ترکیب خطی از پارامترها موضوع قسمت ۶-۵ است. در این قسمت خواهیم دید که یک ترکیب خطی را می‌توان با آماره t آزمون کرد، ولی برای بیش از یک ترکیب خطی باید ضرورتاً از آزمون F استفاده کنیم که در قسمت ۶-۸ مطرح شده است. سؤال بسیار مهم دیگری مطرح می‌شود که اگر برای دو یا چند پارامتر به طور جداگانه فواصل اطمینان بسازیم، آیا دقیقاً به همان جوابی می‌رسیم که از ساختن همزمان ناحیه اطمینان برای همان پارامترها به دست می‌آید؟ در قسمت ۶-۶ خواهیم دید که پاسخ به این سؤال همواره مثبت نیست و به موازات افزایش همبستگی بین متغیرهای توضیحی، این تفاوت به مراتب بیشتر می‌شود.

در بسیاری از تحلیلهای اقتصادسنجی، موضوع مورد علاقه پژوهشگر می‌تواند این باشد که آیا پارامترهای چندمتغیر توضیحی در یک مدل رگرسیون چندمتغیره با یکدیگر برابر است؟ این سؤال با نام آزمون فرضیه $\beta_1 = \beta_{1+1} = \dots = \beta_j$ در قسمت ۶-۷ مطرح شده است. سرانجام در قسمت ۶-۸ آزمون همزمان چند ترکیب خطی از پارامترها را به زبان ماتریسی بررسی کرده‌ایم. در این قسمت، عمومی‌ترین آماره آزمون ارائه شده است و نشان داده‌ایم که بسیاری از آزمونهایی که در قسمتهای قبل مطرح گردیده است، می‌تواند حالت‌های خاصی از این آماره باشد. با اینکه مباحث این فصل مدلهای رگرسیون

با چند متغیر توضیحی است، اما به علت مشکلات محاسباتی، سعی شده است در مثالهای متعددی که آورده‌ایم، از مدل‌های بسیار کوچک استفاده شود.

۶-۲ خصوصیات آماری $\hat{\beta}$

مدل رگرسیون چندمتغیره ۵-۴ را یک بار دیگر می‌نویسیم:

$$y = X\beta + u,$$

که در آن، $y \rightarrow n \times 1$ ، $X \rightarrow n \times k$ ، $\beta \rightarrow k \times 1$ و $u \rightarrow n \times 1$ است. می‌دانیم تخمین پارامترهای این مدل، بر طبق معادله ۵-۲۱، عبارت است از

$$\hat{\beta} = (X'X)^{-1} X'y.$$

با جایگزینی معادله ۵-۴ در ۵-۲۱ داریم

$$\hat{\beta} = (X'X)^{-1} X'(X\beta + u),$$

یا

$$\hat{\beta} = \beta + (X'X)^{-1} X'u. \quad (6.1)$$

ملاحظه می‌شود $\hat{\beta}$ در معادله ۶-۱ تابعی از جمله اختلال (u) است؛ بنابراین یک بردار تصادفی است. سؤال این است که تابع توزیع، میانگین، واریانس و کوواریانس $\hat{\beta}$ چیست؟ با توجه به اینکه $\hat{\beta}$ یک تابع خطی از u است، اگر u توزیع نرمال داشته باشد، $\hat{\beta}$ نیز توزیع نرمال خواهد داشت. در رابطه ۵-۱۱ دیدیم که فرض چهارم از فرضهای کلاسیک جمله اختلال، در واقع نرمال بودن توزیع آن با میانگین 0 و واریانس $\sigma^2 I$ است؛ بنابراین نتیجه می‌گیریم که $\hat{\beta}$ نیز دارای توزیع نرمال است. حال به بررسی میانگین و واریانس آن می‌پردازیم.

۱. میانگین $\hat{\beta}$

معادله ۶-۱ را ملاحظه کنید. می‌دانیم که اولین فرض از فرضهای کلاسیک متغیرهای

توضیحی این است که هر یک از متغیرهای توضیحی غیر تصادفی بوده و در آزمایشهای فرضی تکراری ثابت حفظ شده‌اند. با توجه به اینکه، از دو طرف معادله ۶-۱ امید ریاضی می‌گیریم. خواهیم داشت

$$E(\hat{\beta}) = E(\beta) + E[(X'X)^{-1} X'u],$$

$$= \beta + (X'X)^{-1} X'E(u).$$

بنابر فرض اول از فرضهای کلاسیک، یعنی معادله ۵-۶، می‌دانیم $E(u) = 0$ ؛ بنابراین

$$E(\hat{\beta}) = \beta. \quad (6.2)$$

معادله ۶-۲ دلالت بر این می‌کند که $\hat{\beta}$ یک تخمین نااریب از β است.

۲. واریانس و کوواریانس $\hat{\beta}$

چون $\hat{\beta}$ یک بردار $(k \times 1)$ است، باید یک ماتریس $(k \times k)$ تعریف کنیم که عناصر قطری آن واریانس مقادیر $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ و عناصر غیرقطری آن کوواریانس β_1 و β_2 به ازای تمام مقادیر z و $z \neq 1$ باشد. برای این منظور می‌گوییم ماتریس $E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)']$ مطلوب است، زیرا $(\hat{\beta} - \beta)$ یک بردار $(k \times 1)$ بوده. و بنابراین ترانهاد آن، یعنی $(\hat{\beta} - \beta)'$ یک بردار $(1 \times k)$ می‌شود. بدین ترتیب حاصلضرب $(\hat{\beta} - \beta)(\hat{\beta} - \beta)'$ یک ماتریس $(k \times k)$ خواهد بود. اگر از این ماتریس امید ریاضی بگیریم، ماتریس واریانس - کوواریانس $\hat{\beta}$ به دست می‌آید،

$$E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] = E \begin{bmatrix} (\hat{\beta}_1 - \beta_1) \\ (\hat{\beta}_2 - \beta_2) \\ \vdots \\ (\hat{\beta}_k - \beta_k) \end{bmatrix} [(\hat{\beta}_1 - \beta_1) (\hat{\beta}_2 - \beta_2) \dots (\hat{\beta}_k - \beta_k)].$$

این دو بردار را در هم ضرب کرده، یک ماتریس $(k \times k)$ به دست می‌آوریم. اگر عملگر E را دخالت دهیم، خواهیم داشت

$$E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' = \begin{bmatrix} E(\hat{\beta}_1 - \beta_1)' E(\hat{\beta}_1 - \beta_1)(\hat{\beta}_2 - \beta_2) \dots E(\hat{\beta}_1 - \beta_1)(\hat{\beta}_k - \beta_k) \\ E(\hat{\beta}_2 - \beta_2)(\hat{\beta}_1 - \beta_1) E(\hat{\beta}_2 - \beta_2)' \dots E(\hat{\beta}_2 - \beta_2)(\hat{\beta}_k - \beta_k) \\ \vdots \\ E(\hat{\beta}_k - \beta_k)(\hat{\beta}_1 - \beta_1) E(\hat{\beta}_k - \beta_k)(\hat{\beta}_2 - \beta_2) \dots E(\hat{\beta}_k - \beta_k)' \end{bmatrix} \quad (6.3)$$

بنابر تعریف می‌دانیم

$$\text{Var}(\hat{\beta}_i) = E[\hat{\beta}_i - E(\hat{\beta}_i)] = E(\hat{\beta}_i - \beta_i)'$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = E[\hat{\beta}_i - E(\hat{\beta}_i)][\hat{\beta}_j - E(\hat{\beta}_j)] = E(\hat{\beta}_i - \beta_i)(\hat{\beta}_j - \beta_j)$$

با جایگزینی تعاریف فوق در معادله ۶.۳ داریم

$$E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' = \begin{bmatrix} \text{Var}(\hat{\beta}_1) & \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) & \dots & \text{Cov}(\hat{\beta}_1, \hat{\beta}_k) \\ \text{Cov}(\hat{\beta}_2, \hat{\beta}_1) & \text{Var}(\hat{\beta}_2) & \dots & \text{Cov}(\hat{\beta}_2, \hat{\beta}_k) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\hat{\beta}_k, \hat{\beta}_1) & \text{Cov}(\hat{\beta}_k, \hat{\beta}_2) & \dots & \text{Var}(\hat{\beta}_k) \end{bmatrix} \quad (6.4)$$

بدین ترتیب نشان دادیم که $E(\hat{\beta} - \beta)(\hat{\beta} - \beta)'$ یک ماتریس $(k \times k)$ و قرینه است که عناصر قطری آن واریانسها و عناصر غیرقطری آن کوواریانسهاست. با این حال هنوز نشان نداده‌ایم که فرمول به دست آوردن مقادیر عددی واریانسها یا کوواریانسها چیست؟ برای به دست آوردن این فرمول، به ترتیب زیر عمل می‌کنیم.

معادله ۶.۱ را به صورت زیر می‌نویسیم،

$$(\hat{\beta} - \beta) = (X'X)^{-1} X'u \quad (6.5)$$

دو طرف معادله فوق را ترانهاد می‌کنیم. می‌دانیم ترانهاد معکوس برابر معکوس ترانهاد

است. بدین ترتیب، ترانهاد $(X'X)^{-1}$ برابر معکوس ترانهاد $(X'X)$ است. چون $(X'X)$ یک ماتریس متقارن است؛ بنابراین با ترانهاد خود برابر خواهد بود. نتیجه می‌گیریم که ترانهاد $(X'X)^{-1}$ برابر $(X'X)^{-1}$ است؛ بنابراین

$$(\hat{\beta} - \beta)' = u' X (X'X)^{-1}.$$

دو طرف معادله فوق را در معادله ۶-۵ ضرب می‌کنیم،

$$(\hat{\beta} - \beta) (\hat{\beta} - \beta)' = (X'X)^{-1} X' u [u' X (X'X)^{-1}].$$

از دو طرف معادله فوق امیدریاضی می‌گیریم. می‌دانیم X ثابت فرض شده است؛ بنابراین

$$E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' = (X'X)^{-1} X' E(uu') X (X'X)^{-1}. \quad (6.6)$$

سمت راست معادله ۶-۶ شامل ۵ ماتریس به شرح زیر است،

$$(X'X)^{-1} \rightarrow (k \times k), \quad E(uu') \rightarrow (n \times n),$$

$$X' \rightarrow (k \times n), \quad X \rightarrow (n \times k).$$

حاصلضرب این ۵ ماتریس، یک ماتریس مربع و متقارن $(k \times k)$ خواهد بود. این ماتریس در واقع همان ماتریس واریانس-کوواریانس تخمینهای $\hat{\beta}$ است که در معادله ۶-۴ تعریف شده است. اگر سمت راست معادله ۶-۴ را که یک ماتریس $(k \times k)$ است، مساوی سمت راست معادله ۶-۶ - که آن هم یک ماتریس $(k \times k)$ است - بگیریم؛ ضرورتاً عناصر متناظر این دو ماتریس یک به یک با هم برابرند. می‌دانیم عناصر قطری ماتریس سمت راست معادله (۶-۴) واریانس تخمینهای $\hat{\beta}$ است؛ بنابراین عناصر قطری ماتریس سمت راست معادله ۶-۶ نیز دقیقاً مقادیر همین واریانسها خواهد بود. به همین ترتیب می‌توان برای کوواریانس استدلال کرد. پس

$$\text{Var}(\hat{\beta}_i) = \{ [(X'X)^{-1} X' E(uu') X (X'X)^{-1}] \text{ماتریس} \}$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = \text{Cov}(\hat{\beta}_j, \hat{\beta}_i)$$

$$= \{ \text{عناصر } a_{ij} \text{ یا } a_{ji} \text{ از ماتریس } [(X'X)^{-1} X' E(uu') X (X'X)^{-1}] \}.$$

برای مثال، در مدل رگرسیون زیر،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k x_{ki} + U_i,$$

اگر بخواهیم واریانس β_2 را به دست آوریم، کافی است که سومین عنصر قطری ماتریس

$$(X'X)^{-1} X' E(uu') X (X'X)^{-1}$$

را به دست آوریم. به همین ترتیب کوواریانس $\hat{\beta}_2$ و $\hat{\beta}_6$ برابر با عنصر ردیف ۲ و ستون ۶ یا عنصر ردیف ۶ و ستون ۲ ماتریس فوق است.

محاسبه واریانس و کوواریانس تخمینهای $\hat{\beta}$ ، برای حالتی که جمله اختلال دارای دو خصوصیت واریانس همسانی و عدم خودهمبستگی باشد، بسیار ساده است. در واقع، در چنین حالتی $E(uu')$ برابر $\sigma^2 I$ می شود و این نکته را در معادله ۵-۱۱ نشان داده ایم. معادله

$$E(uu') = \sigma^2 I,$$

را که در آن σ^2 واریانس جمله اختلال و I ماتریس $(n \times n)$ شامل عناصر قطری یک و غیرقطری صفر است، در معادله ۶-۶ قرار می دهیم،

$$\begin{aligned} E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' &= (X'X)^{-1} X' (\sigma^2 I) X (X'X)^{-1}, \\ &= \sigma^2 (X'X)^{-1} X' X (X'X)^{-1}, \end{aligned}$$

در نتیجه خواهیم داشت

$$E(\hat{\beta} - \beta)(\hat{\beta} - \beta)' = \sigma^2 (X'X)^{-1}. \quad (6.7)$$

ملاحظه می شود در این حالت، محاسبه واریانسها و کوواریانسها بسیار ساده است، مثلاً برای به دست آوردن $\text{Var}(\hat{\beta}_2)$ ، سومین عنصر قطری $(X'X)^{-1}$ را گرفته، در σ^2 ضرب می کنیم. به همین ترتیب $\text{Cov}(\hat{\beta}_2, \hat{\beta}_3)$ برابر با حاصلضرب σ^2 در عنصر a_{23} یا a_{32} از $(X'X)^{-1}$ خواهد بود. در نتیجه می توان چنین نوشت،

$$\text{Var}(\hat{\beta}_i) = \sigma^2 [i \text{ امین عنصر قطری ماتریس } (X'X)^{-1}], \quad (6.8)$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = \sigma^2 [\text{عنصر } a_{ij} \text{ از ماتریس } (\mathbf{X}'\mathbf{X})^{-1}] \quad (6.9)$$

در محاسبه واریانس و کوواریانس تخمینهای $\hat{\beta}$ باید به این نکته توجه کنیم که اگر ماتریس $\mathbf{X}'\mathbf{X}$ بر حسب مشاهدات اصلی و نه انحراف از میانگین محاسبه شده و مدل رگرسیون مفروض نیز دارای جمله ثابت باشد، می توان از معادله های 6.8 و 6.9 استفاده کرد. اما معمولاً ماتریس $(\mathbf{X}'\mathbf{X})$ را به کمک مقادیر انحراف از میانگین متغیرهای توضیحی حساب می کنند. در این حالت اگر مدل رگرسیون مفروض جمله ثابت داشته باشد، دیگر نمی توان از معادله های فوق استفاده کرد؛ زیرا در مدل k پارامتر وجود دارد، در حالی که ماتریس $(\mathbf{X}'\mathbf{X})^{-1}$ ، یک ماتریس مربع $(k-1) \times (k-1)$ است. بدیهی است باید از فرمولهای زیر استفاده شود،

$$\text{Var}(\hat{\beta}_i) = \sigma^2 [(i-1)\text{امین عنصر قطری ماتریس } (\mathbf{X}'\mathbf{X})^{-1}] \quad (6.10)$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = \sigma^2 [\text{عنصر } a_{(i-1)(j-1)} \text{ از ماتریس } (\mathbf{X}'\mathbf{X})^{-1}] \quad (6.11)$$

۳. قضیه گاس - مارکف *

می دانیم تخمین $\hat{\beta}$ در مدل $y = \mathbf{X}\beta + u$ ، طبق معادله های 5.21 و 6.1 برابر است با

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'y,$$

$$\hat{\beta} = \beta + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'u.$$

اگر $\mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$ باشد، معادله های 5.21 و 6.1 را می توان به ترتیب به صورت زیر نوشت،

$$\hat{\beta} = \mathbf{A}y, \quad \hat{\beta} = \mathbf{A}u,$$

با توجه به ثابت بودن \mathbf{X} در آزمایشهای تکراری، نتیجه می گیریم که \mathbf{A} ثابت است. بنابراین می توان گفت که $\hat{\beta}$ یک ترکیب خطی از u یا y است. نااریب بودن $\hat{\beta}$ ، یعنی $E(\hat{\beta}) = \beta$ را نیز در معادله 6.2 ملاحظه کردیم. بنابراین، تا این قسمت از بحث به این

نتیجه رسیدیم که $\hat{\beta}_{OLS}$ یک تخمین زنده ناریب و خطی از β است. برای اینکه $\hat{\beta}_{OLS}$ بتواند بهترین تخمین زنده ناریب خطی باشد، باید ثابت کنیم که کمترین واریانس را نیز دارد؛ یعنی هیچ تخمین زنده دیگری را نمی توان یافت که در قلمرو تخمین زنده های خطی ناریب، واریانس کمتری از واریانس $\hat{\beta}_{OLS}$ داشته باشد.

برای اثبات، فرض می کنیم که یک تخمین زنده خطی ناریب غیر OLS مانند $\hat{\beta}_*$ وجود دارد. با توجه به خطی بودن این تخمین زنده، خواهیم داشت

$$\hat{\beta}_* = Dy.$$

اگر D برابر A باشد، آنگاه $\hat{\beta}_* = \hat{\beta}_{OLS}$. فرض می کنیم که $D = A + C$ ؛ بنابراین

$$\hat{\beta}_* = (A + C)y.$$

با جایگزینی معادله ۵-۴ در معادله فوق داریم:

$$\begin{aligned}\hat{\beta}_* &= (A + C)(X\beta + u), \\ &= (A + C)X\beta + (A + C)u.\end{aligned}$$

با توجه به مقدار ماتریس A ، ملاحظه می شود که

$$AX = (X'X)^{-1} X'X = I,$$

در نتیجه

$$\begin{aligned}\hat{\beta}_* &= \beta + CX\beta + (A + C)u, \\ &= (I + CX) + (A + C)u.\end{aligned}$$

از دو طرف معادله فوق امید ریاضی می گیریم،

$$E(\hat{\beta}_*) = (I + CX)\beta.$$

اگر قرار است $\hat{\beta}_*$ ناریب باشد، باید ضرورتاً

$$CX = 0$$

بنابراین بر فرض نااریب بودن $\hat{\beta}_*$ خواهیم داشت

$$\hat{\beta}_* = \beta + (A + C) u.$$

با توجه به تعریف واریانس، می توان چنین نوشت

$$\begin{aligned} \text{Var}(\hat{\beta}_*) &= E[(\hat{\beta}_* - \beta)(\hat{\beta}_* - \beta)'], \\ &= E\{[(A + C)u][(A + C)u]'\}, \\ &= E[(A + C)uu'(A + C)'], \\ &= (A + C)E(uu')(A + C)'. \end{aligned}$$

بر اساس فرضهای واریانس همسانی و عدم خودهمبستگی، داریم $E(uu') = \sigma^2 I$ ؛ بنابراین

$$\text{Var}(\hat{\beta}_*) = \sigma^2 (A + C)(A + C)'. \quad \text{اما می دانیم که}$$

$$(A + C)(A + C)' = A'A + CA' + AC' + CC'.$$

مقدار A را در رابطه فوق قرار می دهیم، خواهیم داشت

$$\begin{aligned} (A + C)(A + C)' &= (X'X)^{-1}X'X(X'X)^{-1} + CX(X'X)^{-1} \\ &+ (X'X)^{-1}X'C' + CC'. \end{aligned}$$

با توجه به اینکه $CX = 0$ بنابراین $X'C = 0$ ، در نتیجه

$$(A + C)(A + C)' = (X'X)^{-1} + CC'.$$

با جایگزینی معادله فوق در واریانس $\hat{\beta}_*$ ، داریم

$$\text{Var}(\hat{\beta}_*) = \sigma^2 [(X'X)^{-1} + CC'] = \sigma^2 (X'X)^{-1} + \sigma^2 CC'.$$

اما می‌دانیم $(X'X)^{-1} \sigma^2$ ، در واقع همان واریانس $\hat{\beta}_{OLS}$ است؛ بنابراین

$$\text{Var}(\hat{\beta}_*) = \text{Var}(\hat{\beta}_{OLS}) + \sigma^2 CC'$$

همچنین می‌دانیم CC' یک «ماتریس نیمه معین مثبت»^۱ است. برای اینکه یک فرم درجه دوم این ماتریس صفر باشد، باید $C=0$ ، یعنی تمام عناصر C باید صفر شود. حال می‌گوییم $\text{Var}(\hat{\beta}_*)$ برابر مجموع دو مقدار غیرمنفی شده است. واریانس $\hat{\beta}_{OLS}$ نمی‌تواند صفر باشد؛ بنابراین برای رسیدن به مقدار حداقل برای واریانس $\hat{\beta}_*$ ، ضرورتاً جمله دوم، یعنی $\sigma^2 CC'$ ، باید مقداری برابر صفر اختیار کند. می‌دانیم لازمه این کار صفر شدن عناصر C است. اما اگر $C=0$ ، آنگاه

$$D = A + C = A,$$

یعنی $\hat{\beta}_* = Dy$ یا $\hat{\beta}_{OLS} = Ay$ دقیقاً متساوی خواهد بود. به عبارت دیگر، در قلمرو تخمین زنده‌های خطی ناریب، تخمین زنده‌ای که حداقل واریانس را دارد، ضرورتاً حداقل مربعات معمولی خواهد بود. بدین ترتیب ثابت شد که $\hat{\beta}_{OLS}$ بهترین تخمین زنده ناریب خطی است، که در واقع اثبات قضیه گاس - مارکف است.

دیدیم که $\hat{\beta}$ توزیع نرمال دارد و میانگین و ماتریس واریانس - کوواریانس آن را نیز به دست آوردیم. این اطلاعات را با توجه به فرضهای کلاسیک می‌توان به صورت زیر نوشت،

$$\hat{\beta} \sim N[\beta, \sigma^2 (X'X)^{-1}]. \quad (7.12)$$

علاوه بر اینکه $\hat{\beta}$ بهترین تخمین زنده ناریب خطی است، می‌توان نشان داد که سازگار نیز هست. می‌دانیم عناصر ماتریس $(X'X)$ مقادیر ثابتی است؛ بنابراین $(X'X)^{-1} \sigma^2$ را می‌توان به صورت $\frac{\sigma^2}{n} (\frac{X'X}{n})^{-1}$ نوشت. برای بررسی خصوصیت سازگاری، باید حد واریانس $\hat{\beta}$ را به دست آوریم، وقتی n به سمت بی‌نهایت میل

1. Positive Semidefinite Matrix

به پیوست «۵ - الف» مراجعه شود.

می‌کند. خواهیم داشت

$$\lim_{n \rightarrow \infty} \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} = 0,$$

زیرا اولاً مقدار $\frac{\sigma^2}{n}$ در حد وقتی n به سمت بی‌نهایت میل کند برابر صفر می‌شود؛ یعنی

$$\lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0,$$

و ثانیاً با توجه به فرضهای کلاسیک در مورد متغیرهای توضیحی، می‌دانیم که ماتریس \mathbf{X} تصادفی نیست و عناصر آن در آزمایشهای تکراری ثابت فرض شده است. همچنین به علت عدم خودهمبستگی کامل، ماتریس $\frac{1}{n}(\mathbf{X}'\mathbf{X})$ در میانه غیر صفر دارد و وارون پذیر است و عناصر آن به ازای $n \rightarrow \infty$ ، مقادیر نامعین ندارد؛ یعنی

$$\lim_{n \rightarrow \infty} \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} = \text{مقدار معین}. \quad (6.13)$$

می‌توان نشان داد که اگر فرض نرمال بودن تابع توزیع احتمال جمله اختلال را کنار بگذاریم، تخمین زنده $\hat{\beta}_{OLS}$ هنوز خصوصیت بهترین تخمین زنده ناریب خطی را حفظ می‌کند و مقدار واریانس آن نیز تغییری نخواهد کرد، اما دیگر یک تخمین زنده «کارآ» نخواهد بود.^۱

۴. تخمین و آزمون واریانس جمله اختلال*

در فرمولهای ۶۸ و ۶۹ و همچنین ۶۱۰ و ۶۱۱ دیدیم که محاسبه واریانس و کوواریانس $\hat{\beta}$ مستلزم داشتن مقدار واریانس جمله اختلال، یعنی σ^2 است. می‌دانیم واریانس جمله اختلال به علت تصادفی بودن ϵ ، باید تخمین زده شود. قبلاً در فصل دوم و معادله ۲-۲۸ فرمول تخمین جمله اختلال را برای یک مدل رگرسیون ساده با دو

۱. به کتاب اقتصادسنجی تألیف (۱۹۷۶) P. Schmidt مراجعه شود.

پارامتر ملاحظه کردیم،

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-k}.$$

در مواردی نیز این نتیجه را به حالت کلی با k پارامتر تعمیم داده‌ایم و از فرمول زیر استفاده کرده‌ایم،

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{n-k} = \frac{\mathbf{e}'\mathbf{e}}{n-k}.$$

در این قسمت می‌خواهیم نه تنها صحت فرمول فوق را ثابت کنیم، بلکه نشان دهیم که اگر از این فرمول استفاده شود یک تخمین نااریب از σ^2 خواهیم داشت. باید از مجموع مربعات پسماند، یعنی $\mathbf{e}'\mathbf{e}$ امید ریاضی بگیریم. برای این منظور ابتدا از معادله ۵-۱۷ شروع می‌کنیم،

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}.$$

مقدار $\hat{\boldsymbol{\beta}}$ از معادله ۵-۲۱ را در معادله فوق قرار می‌دهیم،

$$\mathbf{e} = \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y},$$

یا

$$\mathbf{e} = [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}. \quad (6-14)$$

ماتریس $(n \times n)$ زیر را تعریف می‌کنیم،

$$\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'. \quad (6-15)$$

ماتریس \mathbf{M} بسیار مهم است؛ زیرا، اولاً یک ماتریس متقارن $(n \times n)$ است، ثانیاً $\mathbf{M} \cdot \mathbf{M} = \mathbf{M}' = \mathbf{M}$ ، یعنی یک ماتریس «خودتوان» است، یعنی وقتی \mathbf{M} را در خودش ضرب کنیم، حاصل با خودش برابر است و ثالثاً، وقتی \mathbf{M} را از سمت چپ در \mathbf{X} ضرب کنیم یک ماتریس صفر به دست می‌آوریم، یعنی

$$\mathbf{MX} = \mathbf{0} \quad (6-16)$$

$$MX = [I - X(X'X)^{-1}X']X, \quad \text{زیرا}$$

$$= X - X(X'X)^{-1}(X'X) = 0.$$

با این مقدمات، معادله ۶-۱۵ را در ۶-۱۴ قرار می‌دهیم، خواهیم داشت

$$e = My,$$

$$= M(X\beta + u),$$

$$= MX\beta + Mu,$$

با توجه به معادله ۶-۱۶ داریم

$$e = Mu. \quad (6.17)$$

مجموع مربعات پسماند، یعنی $\sum e_i^2$ ، را می‌نویسیم و معادله فوق را در آن جایگزین می‌کنیم،

$$\sum e_i^2 = e'e,$$

$$= u'M'Mu.$$

می‌دانیم چون M متقارن است؛ پس $M'M = M$ بنابراین

$$M'M = MM = M.$$

در نتیجه

$$e'e = u'Mu.$$

همچنین ملاحظه می‌شود که $u'Mu$ یک شکل درجه دوم است که مقدار آن برابر است با

$$e'e = [U_1 \ U_2 \ \dots \ U_n] \begin{bmatrix} M_{11} & M_{12} & \dots & M_{1n} \\ M_{21} & M_{22} & \dots & M_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ M_{n1} & M_{n2} & \dots & M_{nn} \end{bmatrix} \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_n \end{bmatrix},$$

$$= \sum_{i=1}^n \sum_{j=1}^n U_i U_j M_{ij}.$$

از رابطه فوق امید ریاضی می‌گیریم،

$$E(\mathbf{e}'\mathbf{e}) = E\left[\sum_{i=1}^n \sum_{j=1}^n U_i U_j M_{ij}\right].$$

با توجه به فرضهای واریانس همسانی و عدم خودهمبستگی، می‌دانیم،

$$\text{Var}(U_i) = E(U_i)^2 = \sigma^2,$$

$$\text{Cov}(U_i, U_j) = E(U_i U_j) = 0,$$

در نتیجه خواهیم داشت

$$E(\mathbf{e}'\mathbf{e}) = \sigma^2 \sum_i M_{ii}. \quad (6.18)$$

اما می‌دانیم که $\sum_i M_{ii}$ با مجموع عناصر قطری ماتریس \mathbf{M} که «اثر ماتریس» نام دارد و با $\text{tr}(\mathbf{M})$ نشان داده می‌شود برابر است. در پیوست، ماتریسها (پیوست ۵-الف) دو خاصیت مهم tr عبارت است از

$$\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B}),$$

که \mathbf{A} و \mathbf{B} ماتریسهای مربع بوده و هم‌رتبه هستند و همچنین

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}),$$

که \mathbf{A} و \mathbf{B} به ترتیب $(m \times n)$ و $(n \times m)$ هستند. با استفاده از این خصوصیات، معادله ۶-۱۸ را چنین می‌نویسیم،

$$\begin{aligned} E(\mathbf{e}'\mathbf{e}) &= \sigma^2 \text{tr}(\mathbf{M}), \\ &= \sigma^2 \text{tr}[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'], \\ &= \sigma^2 \{\text{tr}(\mathbf{I}) - \text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\}. \quad (6.19) \end{aligned}$$

اما می دانیم $\text{tr}(\mathbf{I}) = n$ ، زیرا \mathbf{I} یک ماتریس $(n \times n)$ است که عناصر قطری آن عدد یک است. همچنین می دانیم

$$\begin{aligned} \text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] &= \text{tr}\{[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{X}'\}, \\ &= \text{tr}\mathbf{X}'[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}], \\ &= \text{tr}(\mathbf{I}) \\ &= k, \end{aligned}$$

زیرا $(\mathbf{X}'\mathbf{X})^{-1}$ یک ماتریس $(k \times k)$ است. بنابراین معادله ۶-۱۹ را می توان چنین نوشت،

$$E(\mathbf{e}'\mathbf{e}) = \sigma^2(n-k),$$

با تقسیم دو طرف معادله فوق بر $(n-k)$ و توجه به این نکته که $(n-k)$ عدد ثابتی است، داریم

$$E\left(\frac{\mathbf{e}'\mathbf{e}}{n-k}\right) = \sigma^2. \quad (6.20)$$

به عبارت دیگر، اگر $\frac{\mathbf{e}'\mathbf{e}}{n-k}$ را تخمینی از σ^2 بدانیم، آنگاه امید ریاضی این تخمین برابر مقدار واقعی پارامتر است؛ یعنی $\frac{\mathbf{e}'\mathbf{e}}{n-k}$ یک تخمین ناریب از σ^2 است. بنابراین

$$\hat{\sigma}_U^2 = \frac{\mathbf{e}'\mathbf{e}}{n-k}. \quad (6.21)$$

اگر تخمین واریانس جمله اختلال را در فرمولهای ۶-۸ و ۶-۹ یا ۶-۱۰ و ۶-۱۱ قرار دهیم، فرمولهای محاسباتی واریانس و کوواریانس $\hat{\beta}$ به دست می آید. برای مدل های رگرسیون با ضریب ثابت که $(\mathbf{X}'\mathbf{X})$ برحسب مشاهدات اصلی محاسبه شده و جمله اختلال کلاسیک است، داریم

$$\text{Var}(\hat{\beta}_i) = \hat{\sigma}_U^2 [(\mathbf{X}'\mathbf{X})^{-1}]_{ii}, \quad (6.22)$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = \hat{\sigma}_U^2 [(\mathbf{X}'\mathbf{X})^{-1}]_{ij}. \quad (6.23)$$

برای مدلهایی که $(X'X)$ برحسب انحراف از میانگین محاسبه شده و فرضهای کلاسیک را نیز شامل است خواهیم داشت

$$\text{Var}(\hat{\beta}_i) = \hat{\sigma}_U^2 [(X'X)^{-1}]_{(i-1)(i-1)}, \quad (6.24)$$

$$\text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = \hat{\sigma}_U^2 [(X'X)^{-1}]_{(i-1)(j-1)}. \quad (6.25)$$

آزمون واریانس جمله اختلال

بحث آزمون σ^2 در رگرسیون چندمتغیره دقیقاً مشابه مطالبی است که در فصل دوم با همین عنوان برای رگرسیون ساده مطرح کردیم. بنابراین می توان در این قسمت، مطالب را فهرست وار مطرح کرد. براساس معادله ۲-۳۸ می توان گفت که

$$\frac{e'e}{\sigma^2} \sim \chi^2(n-k),$$

که در آن $e'e$ مجموع مربعات پسماند در مدلی است که k پارامتر دارد. ملاحظه می شود که $\frac{e'e}{\sigma^2}$ دارای توزیع χ^2 با $(n-k)$ درجه آزادی است. اگر بخواهیم یک فاصله اطمینان ۹۵ درصد برای $\frac{e'e}{\sigma^2}$ به دست آوریم، باید سطحی از منحنی تغییرات χ^2 را تعیین کنیم که ۹۵ درصد سطح زیر منحنی در ناحیه مرکزی و 0.025 در هر یک از دو طرف باشد. اگر نقاط متناظر با ابتدا و انتهای این سطح در روی محور χ^2 را $\chi^2_{0.025}$ و $\chi^2_{0.975}$ بنامیم، خواهیم داشت

$$\text{Pr}(\chi^2_{0.025} < \frac{e'e}{\sigma^2} < \chi^2_{0.975}) = 0.95,$$

یا به عبارت دیگر، در سطح احتمال ۹۵ درصد داریم

$$\chi^2_{0.025} < \frac{e'e}{\sigma^2} < \chi^2_{0.975}.$$

اما در آزمون واریانس جمله اختلال، هدف ما این است که یک فاصله اطمینان برای σ^2 به دست آوریم. با مراجعه به معادله ۶-۲۱ داریم

$$\hat{\sigma}_U^2 = \frac{e'e}{n-k},$$

یا

$$e'e = (n-k) \hat{\sigma}_U^2.$$

معادله فوق را در فاصله اطمینان برای $\frac{e'e}{\sigma^2}$ جایگزین می‌کنیم،

$$\chi_{0.1, 170}^2 < \frac{(n-k) \hat{\sigma}^2}{\sigma^2} < \chi_{0.9, 170}^2.$$

دو طرف رابطه فوق را معکوس کرده، سپس در $(n-k) \hat{\sigma}^2$ ضرب می‌کنیم. خواهیم داشت

$$\frac{(n-k) \hat{\sigma}^2}{\chi_{0.9, 170}^2} < \sigma^2 < \frac{(n-k) \hat{\sigma}^2}{\chi_{0.1, 170}^2}.$$

رابطه فوق فاصله اطمینان ۹۵ درصد برای σ^2 را تعیین می‌کند. فرضیه‌های مختلف برای σ^2 را می‌توان با داشتن این فاصله اطمینان آزمون کرد. اگر H_0 در قلمرو این فاصله اطمینان قرار گیرد، قبول شده، در غیر این صورت رد خواهد شد.

مثال ۹-۱ مدل رگرسیون را که در مثال ۵-۲ مطرح شده است، ملاحظه کنید. می‌خواهیم برای واریانس جمله اختلال این مدل، یک فاصله اطمینان ۹۵ درصد بسازیم.

بر اساس محاسبات انجام شده در مثال ۵-۲ می‌دانیم $\hat{\sigma}_U^2 = 0.75$. همچنین می‌دانیم $n = 5$ و $k = 3$. مقادیر $\chi_{0.1, 170}^2$ را از جدول χ^2 و با درجات آزادی $n - k = 2$ به دست می‌آوریم. خواهیم داشت

$$\chi_{0.9, 170}^2 = 7/38 \quad , \quad \chi_{0.1, 170}^2 = 0/0506.$$

پنابراین

$$\Pr \left(\frac{(5-3) \hat{\sigma}^2}{\chi_{0.9, 170}^2} < \sigma^2 < \frac{(5-3) \hat{\sigma}^2}{\chi_{0.1, 170}^2} \right) = 0/95,$$

یا

$$\frac{(5-3) \cdot 0/75}{7/38} < \sigma^2 < \frac{(5-3) (0/75)}{0/0506},$$

$$0/203 < \sigma^2 < 29/64.$$

ملاحظه می‌شود که فاصله اطمینان بسیار بزرگ است که به معنی پایین بودن دقت در تخمین واریانس جمله اختلال است. ضعف تخمینها در تمام مواردی که در مورد مثال ۵-۲ مطرح کردیم یا خواهیم کرد، ملاحظه می‌شود. این امر به دلیل تعداد بسیار کم مشاهدات ($n=5$) در مقایسه با تعداد پارامترها ($k=3$) است.

۶-۳ آزمون هر یک از پارامترها

در این قسمت، ابتدا آزمون $\beta_1 = a$ که a عدد ثابت و معلومی است، و نیز آزمون $\beta_1 = 0$ بررسی می‌شود. در ادامه این قسمت فرضیه معنی دار بودن زیر مجموعه‌ای از پارامترهای یک مدل رگرسیون چند متغیره، یعنی فرضیه $\beta_1 = \beta_{1+1} = \dots = \beta_j = 0$ مطرح می‌شود.

۱. آزمون فرضیه $\beta_1 = a$

مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i.$$

روش آزمون فرضیه در مورد هر یک از پارامترهای این مدل دقیقاً مانند آزمون فرضیه در رگرسیون ساده است. پارامتر β_1 را در نظر می‌گیریم. فرض کنید می‌خواهیم فرضیه $H_0: \beta_1 = a$ را در مقابل فرضیه $H_1: \beta_1 \neq a$ در سطح معنی دار α درصد آزمون کنیم. ابتدا مدل را تخمین می‌زنیم و $\hat{\beta}_1$ را به دست آورده و آن را استاندارد می‌کنیم. می‌دانیم برای استاندارد کردن $\hat{\beta}_1$ باید آن را از میانگین خود کم کرده و بر انحراف معیارش تقسیم کرد. با توجه به اینکه انحراف معیار همان جذر واریانس است، ابتدا باید واریانس $\hat{\beta}_1$ را تخمین زد. فرمول ۶-۲۲ یا ۶-۲۴ تخمین واریانس $\hat{\beta}_1$ را تعیین می‌کند. با داشتن تخمین واریانس $\hat{\beta}_1$ ، مقدار استاندارد شده $\hat{\beta}$ توزیع t با $(n-k)$ درجه آزادی خواهد داشت؛ بنابراین

$$t = \frac{\hat{\beta}_1 - E(\hat{\beta}_1)}{\sqrt{\text{Var}(\hat{\beta}_1)}} \sim t(n-k).$$

با توجه به $E(\hat{\beta}_i) = \beta_i$ و $SE(\hat{\beta}_i) = \sqrt{\text{Var}(\hat{\beta}_i)}$ خواهیم داشت

$$t = \frac{\hat{\beta}_i - \beta_i}{SE(\hat{\beta}_i)} \sim t(n-k).$$

اگر فرضیه H_0 صحیح باشد، داریم

$$t = \frac{\hat{\beta}_i - a}{SE(\hat{\beta}_i)} \sim t(n-k). \quad (6.26)$$

مقدار t که از معادله فوق به دست می‌آید، یعنی آماره آزمون، را با مقدار جدول t مقایسه می‌کنیم؛ اگر آماره آزمون در ناحیه بحرانی قرار گرفت فرضیه H_0 رد شده و در غیر این صورت رد نخواهد شد. اگر آزمون در سطح معنی دار α درصد باشد، مقدار t را از جدول، باید در سطح معنی دار $\alpha/2$ و با $(n-k)$ درجه آزادی به دست آورد؛ زیرا آزمون ما دو طرفه است.

براساس رابطه ۶۲۶ می‌توان برای β_i فاصله اطمینان نیز ساخت. فرض کنید می‌خواهیم یک فاصله اطمینان $(1 - \alpha)$ درصدی برای β_i بسازیم. با توجه به معادله

۴۱-۲ داریم

$$-t_{\alpha/2} < \frac{\hat{\beta}_i - \beta_i}{SE(\hat{\beta}_i)} < t_{\alpha/2},$$

یا

$$\hat{\beta}_i - t_{\alpha/2} SE(\hat{\beta}_i) < \beta_i < \hat{\beta}_i + t_{\alpha/2} SE(\hat{\beta}_i). \quad (6.27)$$

بدیهی است با داشتن فاصله اطمینان برای β_i ، می‌توان فرضیه $H_0: \beta_i = a$ را نیز آزمون کرد. اگر $\beta_i = a$ در فاصله اطمینان واقع شود، H_0 را قبول می‌کنیم و در غیر این صورت H_0 رد خواهد شد.

آزمون معنی دار بودن β_i

آزمون فرضیه $H_0: \beta_i = a$ را آزمون فرضیه معنی دار بودن β_i می‌گویند. اگر H_0 را نتوان

رد کرد، X_{ii} در توضیح تغییرات Y_i تأثیری نخواهد داشت. کافی است در رابطه ۶-۲۶ به جای a مقدار صفر قرار دهیم. خواهیم داشت

$$t = \frac{\hat{\beta}_i - 0}{SE(\hat{\beta}_i)} = \frac{\hat{\beta}_i}{SE(\hat{\beta}_i)} \quad (6.28)$$

مقدار آماره آزمون را با مقدار جدول t مقایسه می‌کنیم. اگر آماره آزمون در ناحیه بحرانی افتاد، H_0 رد می‌شود، در غیر این صورت H_0 را قبول می‌کنیم. یادآوری می‌شود که اگر $H_0: \beta_i = 0$ ، آزمون ما دوطرفه خواهد بود. اگر آزمون در سطح معنی‌دار α درصد اجرا شود، باید مقدار t را با $(n-k)$ درجه آزادی و در سطح معنی‌دار $\alpha/2$ درصد از جدول t به دست آوریم.

مثال ۶-۲ مدلی را که در مثال ۵-۲ مطرح کردیم، یک بار دیگر در نظر می‌گیریم،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

با استفاده از مشاهدات مندرج در جدول ۵-۲ این مدل را به صورت زیر تخمین زده‌ایم،

$$\hat{Y}_i = 4 + 2/5 X_i - 1/5 Z_i.$$

همچنین در مثال ۵-۲ دیدیم که:

$$(X'X)^{-1} = \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix}, \quad \begin{array}{l} ESS = 26/5, \quad TSS = 28, \\ RSS = 1/5, \quad R^2 = 0.9664. \end{array}$$

با اطلاعات فوق

(الف) واریانس U_i را تخمین بزنید.

(ب) واریانس $\hat{\beta}$ و واریانس $\hat{\gamma}$ را تخمین بزنید.

(ج) فرضیه $H_0: \beta_i = 2$ را در مقابل فرضیه $H_1: \beta_i \neq 2$ در سطح معنی‌دار ۵ درصد آزمون کنید.

(د) فرضیه $H_0: \gamma = 0$ را در مقابل فرضیه $H_1: \gamma \neq 0$ در سطح معنی‌دار ۵ درصد آزمون کنید.

ه) برای β یک فاصله اطمینان ۹۵ درصد بسازید.

برای قسمت الف، یعنی تخمین واریانس U_1 از فرمول ۶-۲۱ استفاده می‌کنیم،

$$\hat{\sigma}_U^2 = \frac{\mathbf{e}'\mathbf{e}}{n-k}.$$

می‌دانیم $\mathbf{e}'\mathbf{e} = 1/5$ و $n = 5$ و $k = 3$. بنابراین

$$\hat{\sigma}_U^2 = \frac{1/5}{5-3} = 0.125.$$

در قسمت ب، تخمین واریانس $\hat{\beta}$ و تخمین واریانس \hat{y} را باید از فرمول ۶-۲۴ به دست آورد؛ زیرا اولاً مدل رگرسیون مفروض شامل جمله ثابت است. ثانیاً ماتریس $(\mathbf{X}'\mathbf{X})$ برحسب مقادیر انحراف از میانگین محاسبه شده است. β_1 در واقع دومین پارامتر در مدل رگرسیون است؛ بنابراین، باید $\hat{\sigma}_U^2$ را در اولین عنصر قطری $(\mathbf{X}'\mathbf{X})^{-1}$ ضرب کنیم. به همین ترتیب برای محاسبه واریانس \hat{y} - که سومین پارامتر است - باید دومین عنصر قطری ماتریس $(\mathbf{X}'\mathbf{X})^{-1}$ در $\hat{\sigma}_U^2$ ضرب شود. خواهیم داشت

$$\text{Var}(\hat{\beta}) = 0.125(1) = 0.125,$$

$$\text{Var}(\hat{y}) = 0.125(2/5) = 0.05.$$

در قسمت ج، رابطه ۶-۲۶ را می‌نویسیم،

$$t = \frac{\hat{\beta}_1 - a}{\text{SE}(\hat{\beta}_1)},$$

که با توجه به $\beta = 2$: $H_0: \beta = 2$ مقدار a برابر ۲ است. مقدار آماره آزمون برابر است با

$$t = \frac{2/5 - 2}{\sqrt{0.125}} = \frac{0.05}{0.3536} = 0.1415.$$

می‌دانیم با دو درجه آزادی و سطح معنی‌دار ۰/۰۲۵ داریم $t_{0.025}(2) = 4.303$. بنابراین آماره آزمون در ناحیه بحرانی قرار نمی‌گیرد، پس معنی‌داریست و فرضیه H_0 رد نمی‌شود.

در قسمت د، برای آزمون معنی دار بودن فرضیه $\gamma = 0$: از رابطه ۶۲۸ استفاده می‌کنیم. آماره آزمون عبارت است از

$$t = \frac{\hat{\gamma}}{SE(\hat{\gamma})} = \frac{-1/5}{\sqrt{1/875}} = \frac{-1/5}{1/369} = -1/1.$$

ملاحظه می‌شود که آماره آزمون معنی دار نیست، زیرا در ناحیه بحرانی قرار نمی‌گیرد. بنابراین $H_0: \gamma = 0$ رد نخواهد شد.

در قسمت ه، برای به دست آوردن فاصله اطمینان ۹۵ درصد برای β ، می‌توان از رابطه ۶۲۷ استفاده کرد،

$$\hat{\beta} - t_{\alpha/2} SE(\hat{\beta}) < \beta < \hat{\beta} + t_{\alpha/2} SE(\hat{\beta}),$$

$$2/5 - 4/30.3(0/1866) < \beta < 2/5 + 4/30.3(0/1866),$$

$$2/5 - 3/7 < \beta < 2/5 + 3/7,$$

$$-1/2 < \beta < 6/2.$$

۴. آزمون معنی دار بودن زیر مجموعه‌ای از پارامترها: «آزمون والد»
مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_2 X_{i1} + \dots + \beta_m X_{im} + V_i. \quad (6.29)$$

فرض می‌کنیم براساس یک نظریه، نه تنها تمام متغیرهای توضیحی فوق ضروری است، بلکه باید r متغیر توضیحی دیگر نیز به مدل اضافه کرد، یعنی

$$Y_i = \beta_1 + \beta_2 X_{i1} + \dots + \beta_m X_{im} + \beta_{m+1} X_{(m+1)i} + \dots + \beta_n X_{ni} + U_i, \quad (6.30)$$

که در آن $r = k - m$. برای آزمون صحت نظریه مفروض باید فرضیه زیر را آزمود.

$$H_1: \beta_{m+1} = \beta_{m+2} = \dots = \beta_k = 0 \quad (6.31)$$

بنابراین موضوع بحث ما در این قسمت تخمین مدل ۶.۳۰ و آزمون فرضیه ۶.۳۱ است. می توان رابطه ۶.۳۱ را مجموعه‌ای از محدودیتهای صفری برای پارامترهای مدل ۶.۳۰ تعبیر کرد. بنابراین، اگر این محدودیتهای را در مدل ۶.۳۰ قرار دهیم، یک مدل رگرسیون مقید^۱ خواهیم داشت که در واقع همان مدل ۶.۲۹ است. در ادامه بحث، هرگاه برای یک متغیر از اندیس r استفاده کنیم، بدین معنی است که آن متغیر متعلق به مدل رگرسیون مقید ۶.۲۹ است. متغیرهای بدون اندیس متعلق به مدل غیر مقید ۶.۳۰ خواهد بود.

در قسمت ۵.۵ و در رابطه ۵.۵۹ دیدیم که اگر در یک رگرسیون چند متغیره، تعداد متغیرهای توضیحی را افزایش دهیم، مجموع مربعات پسماند کاهش یافته، یا حداقل ثابت می ماند. این امر بدین معنی است که R^2 افزایش یافته و یا حداقل ثابت خواهد ماند. بنابراین، انتظار ما این است که مجموع مربعات پسماند در مدل ۶.۳۰، یعنی $\sum e_i^2 = RSS$ ، از مقدار مشابه آن در مدل مقید ۶.۲۹، یعنی $\sum e_i^2 = RSS_r$ ، کمتر بوده یا حداقل مساوی آن باشد. در اینجا کمیت زیر را تشکیل می دهیم،

$$\frac{(RSS_r - RSS)}{RSS}$$

اگر فرضیه H_0 صحیح باشد؛ یعنی متغیرهای توضیحی $X_{(m+1)}$ تا X_{kt} نتوانند تغییرات Y_t را توضیح دهند، دو مدل مقید ۶.۲۹ و غیر مقید ۶.۳۰ در واقع برای جامعه مشاهدات متغیرها، با یکدیگر تفاوت ندارند و یکی هستند. در چنین حالتی خواهیم داشت $RSS_r = RSS$ ؛ بنابراین کسر فوق برابر صفر می شود. البته اگر در عمل و با استفاده از نمونه گیری، این دو مدل را تخمین بزنیم، قطعاً نتایج تخمین پارامترها با یکدیگر متفاوت خواهد بود، که می توان گفت این تفاوت ناشی از خطای نمونه گیری است.

نتیجه می‌گیریم که مقدار کسر ۰.۳۱ یا صفر است یا عددی مثبت و هیچگاه منفی نخواهد بود. اگر صفر شد یا تفاوت مختصری با صفر داشت، فرضیه H_0 را می‌توان قبول کرد. بدیهی است برای پاسخ دادن به این سؤال که این «تفاوت مختصر» در عمل دقیقاً چقدر است، باید از آزمونهای آماری استفاده کرد؛ بنابراین باید یک تفسیر آماری از کسر فوق ارائه دهیم. به ترتیب زیر عمل می‌کنیم.

الف) مدل رگرسیون ۰.۲۹ را تخمین زده و مجموع مربعات پسماند آن، یعنی RSS_p را به دست می‌آوریم.

ب) به همین ترتیب مقدار RSS را از مدل ۰.۳۰ حساب می‌کنیم.

ج) تفاوت $(RSS_p - RSS)$ ، در واقع برابر کاهش است که در مجموع مربعات پسماند

مدل ۰.۲۹ حاصل می‌شود و این کاهش به علت افزایش تعداد r متغیر توضیحی است.

د) بر اساس معادله ۲.۳۸ می‌دانیم که در یک مدل رگرسیون با k پارامتر، آماره $\frac{\sum e_i^2}{\sigma^2}$ دارای توزیع χ^2 با $(n-k)$ درجه آزادی است؛ یعنی

$$\frac{RSS}{\sigma^2} \sim \chi^2 (n-k). \quad (6.32)$$

می‌گوییم که اگر فرضیه H_0 صحیح باشد، یعنی درحقیقت مدل‌های ۰.۲۹ و ۰.۳۰ با یکدیگر تفاوتی نداشته باشد، مقادیر واقعی واریانس جمله اختلال در هر یک از این دو معادله با یکدیگر برابر است. اگر صورت و مخرج کسر ۰.۳۱ را بر واریانس جمله اختلال تقسیم کنیم، هر کدام دارای توزیع χ^2 خواهند بود. درجات آزادی هر یک به صورت زیر است. درجه آزادی صورت r است؛ زیرا درجات آزادی RSS_p و RSS به ترتیب برابر است با $(n-m)$ و $(n-k)$ ؛ بنابراین درجات آزادی تفاوت آنها عبارت خواهد بود از

$$(n-m) - (n-k) = k - m = r.$$

درجات آزادی مخرج، یعنی تغییرات توضیح داده نشده نیز برابر $(n-k)$ است.

ه) می‌دانیم نسبت دو توزیع χ^2 - که هر یک بر درجات آزادی خود تقسیم شده باشد - توزیع F دارد؛ بنابراین اگر صورت و مخرج کسر به دست آمده در مرحله «د» را به ترتیب بر r و $(n-k)$ تقسیم کنیم، کسر به دست آمده توزیع F خواهد داشت. با توجه به

اینکه واریانس جمله اختلال، یعنی σ^2 از صورت و منخرج حذف می شود، در نتیجه

$$\frac{(RSS_r - RSS)/r}{RSS/(n-k)} \sim F(r, n-k). \quad (۶.۳۳)$$

یادآوری می کنیم که رابطه فوق به شرطی برقرار است که H_0 صحیح باشد.

و) با تعیین سطح معنی دار بودن آزمون، به راحتی می توان فرضیه H_0 را آزمون کرد. مقدار F با درجات آزادی r و $(n-k)$ را در سطح معنی دار α از جدول F به دست می آوریم و با آماره آزمون که در رابطه ۶.۳۳ محاسبه شده است مقایسه می کنیم. اگر آماره آزمون از مقدار جدول F بیشتر بود معنی دار است و فرضیه H_0 رد می شود.

مثال ۶.۳ برای تخمین تابع مصرف، دو مدل به شرح زیر پیشنهاد شده است،

$$C_t = \beta_1 + \beta_2 Y_t + U_t,$$

$$C_t = \beta_1 + \beta_2 Y_t + \beta_3 C_{t-1} + \beta_4 L_{t-1} + U_t,$$

که در آن C_t مصرف، Y_t درآمد و L_t داراییهای نقدی است. برای بررسی تفاوت این دو مدل، می خواهیم فرضیه زیر را

$$H_0: \beta_3 = \beta_4 = 0$$

در مقابل فرضیه مخالف صفر و در سطح معنی دار یک درصد آزمون کنیم. همچنین می دانیم

$$n = 31, \quad RSS_r = 13/824, \quad RSS = 3/584.$$

ابتدا آماره آزمون را با استفاده از رابطه ۶.۳۳ می نویسیم،

$$F = \frac{(RSS_r - RSS)/r}{RSS/(n-k)}.$$

در مثال فوق $r = 2$, $k = 4$. بنابراین مقدار عددی آماره آزمون عبارت است از

$$F = \frac{(13/824 - 3/584)/2}{3/584/(31-4)} = 30.57$$

مقدار F ، با درجات آزادی $(2, 27)$ و در سطح معنی دار 0.01 برابر است با

$$F_{0.01}(2, 27) = 5.49.$$

ملاحظه می شود آماره آزمون معنی دار است؛ زیرا در ناحیه بحرانی قرار می گیرد. در نتیجه فرضیه H_0 رد شده و ضرورت وجود C_{t-1} و L_{t-1} به عنوان متغیرهای توضیحی در تابع مصرف اثبات می شود.

۳. آماره آزمون F برحسب R^2

می توان رابطه ۶.۳۳ را فقط برحسب ضرایب تعیین در مدل رگرسیون غیرمقید و مدل رگرسیون مقید نوشت. با توجه به معادله ۵.۴۸ می دانیم

$$R^2 = 1 - \frac{e'e}{y'y},$$

یا

$$RSS = (1 - R^2) TSS.$$

معادله فوق می تواند برای مدل غیرمقید مورد استفاده قرار گیرد. برای مدل مقید خواهیم داشت

$$RSS_p = (1 - R_p^2) TSS.$$

یادآوری می شود که مقدار $TSS = \sum y_t^2$ در هر دو مدل یکسان است. با جایگزینی معادله های فوق در رابطه ۶.۳۳ داریم

$$\frac{[(1 - R_p^2) TSS - (1 - R^2) TSS] / r}{[(1 - R_p^2) TSS] / (n - k)} \sim F(r, n - k),$$

با حذف کل تغییرات (TSS) از صورت و مخرج خواهیم داشت

$$\frac{(n - k)}{r} \cdot \frac{R^2 - R_p^2}{1 - R^2} \sim F(r, n - k). \quad (6.34)$$

مثال ۶.۴ دو مدل مصرف (مثال ۶.۳) را یک بار دیگر ملاحظه کنید،

$$C_t = \beta_1 + \beta_2 Y_t + U_t,$$

$$C_t = \beta_1 + \beta_2 Y_t + \beta_3 C_{t-1} + \beta_4 L_{t-1} + U_t.$$

فرض کنید این دو مدل را تخمین زده و مقادیر ضریب تعیین تعدیل شده زیر را به دست آورده‌ایم،

$$\bar{R}^2 = 0/984 \quad \bar{R}_1^2 = 0/944.$$

با توجه به $n = 31$ ، فرضیه $H_0: \beta_3 = \beta_4 = 0$ را در سطح معنی دار $0/01$ و در مقابل فرضیه مخالف صفر آزمون کنید.

ابتدا باید R^2 را حساب کنیم. معادله 0.64 را نوشته، خواهیم داشت

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k} \right) (1 - R^2),$$

$$0/984 = 1 - \left(\frac{30-1}{30-4} \right) (1 - R^2),$$

$$R^2 = 0/986.$$

برای مدل مقید داریم

$$\bar{R}_r^2 = 1 - \left(\frac{n-1}{n-k} \right) (1 - R_r^2),$$

$$0/944 = 1 - \left(\frac{31-1}{31-2} \right) (1 - R_r^2),$$

$$R_r^2 = 0/946$$

در معادله 0.64 آماره آزمون را بر حسب R^2 و R_r^2 به دست آوریم، یعنی

$$\frac{(n-k)}{r} \cdot \frac{R^2 - R_r^2}{1 - R^2} \sim F(r, n-k),$$

بنابراین

$$F = \frac{(31-4)}{2} \cdot \frac{0/986 - 0/946}{1 - 0/986} = \frac{27}{2} \left(\frac{0/04}{0/014} \right) = 38/07.$$

با توجه به اینکه $F_{0.01}(2, 27) = 5/49$ ؛ بنابراین آماره آزمون معنی دار بوده و فرضیه H_0 رد می شود.

۴. آزمون معنی دار بودن β_1 با آماره F

قبلاً در معادله ۶-۲۸ دیدیم که چگونه می توان آزمون معنی دار بودن β_1 را به کمک آماره t انجام داد. در این قسمت نشان می دهیم که این آزمون، حالت خاصی از آزمون معنی دار بودن زیر مجموعه ای از پارامترهاست. مدل زیر را در نظر می گیریم،

$$Y_t = \beta_1 + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + U_t.$$

می خواهیم فرضیه $H_0: \beta_1 = 0$ را در مقابل فرضیه $H_1: \beta_1 \neq 0$ در سطح معنی دار α درصد آزمون کنیم.

مدل فوق را مدل غیرمقید می نامیم. فرضیه $\beta_1 = 0$ در این مدل قرار می دهیم تا مدل مقید به دست آید. هر دو مدل را تخمین می زنیم و به کمک RSS_0 و RSS_1 ، آماره F را - که در ۶-۳۳ تعریف شده است - می سازیم. با مقایسه آماره آزمون، با مقدار موجود در جدول F نسبت به فرضیه H_0 ، استنباط آماری خواهیم کرد. بنابراین، مراحل آزمون به صورت زیر است.

(الف) متغیر توضیحی X_{1t} را از مدل حذف کرده، مدل مقید به دست آمده را تخمین می زنیم تا RSS_0 به دست آید.

(ب) مدل مفروض را با وجود X_{1t} تخمین می زنیم.

(ج) مقدار کاهش در مجموع مربعات پسماند را که از مرحله (الف) به (ب) ملاحظه می شود حساب می کنیم؛ یعنی $(RSS_0 - RSS_1)$.

(د) آماره آزمون را به شرح زیر می سازیم،

$$\frac{(RSS_0 - RSS_1) / 1}{RSS_1 / (n - k)} \sim F(1, n - k),$$

(ه) از مقایسه آماره آزمون با مقدار جدول F ، می توان در سطح معنی دار

۴۷۰ اقتصادسنجی: تک معادلات با فروض کلاسیک

α نسبت به فرضیه مورد نظر استنباط آماری کرد.

مثال ۶-۵ مدل رگرسیون زیر را که موضوع مثال ۶-۲ بوده است یک بار دیگر ملاحظه کنید،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

در بند «د» مثال مزبور، فرضیه $H_1: \gamma = 0$ را در مقابل فرضیه $H_0: \gamma \neq 0$ در سطح معنی دار ۵ درصد آزمون کرده‌ایم. همین آزمون را با استفاده از آماره F انجام دهید.

برای این منظور، ابتدا متغیر توضیحی Z_i را از مدل حذف کرده و مدل مقید به دست آمده را تخمین می‌زنیم. با استفاده از مثال ۵-۲ و محاسبات جدول ۵-۳ می‌دانیم که $\sum x_i^2 = 10$ و $\sum x_i y_i = 16$.

خواهیم داشت

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{16}{10} = 1.6$$

باید مقدار RSS_r را برای این مدل مقید حساب کنیم. با توجه به معادله ۱-۳۵ داریم

$$R_r^2 = \frac{\hat{\beta} \sum x_i y_i}{\sum x_i^2} = \frac{1.6(16)}{28} = 0.9143.$$

طبق معادله ۱-۴۳ می‌دانیم

$$\sum e_i^2 = (1 - R_r^2) \sum y_i^2,$$

در نتیجه برای این مثال خواهیم داشت

$$\begin{aligned} RSS_r &= (1 - R_r^2) TSS, \\ &= (1 - 0.9143) 28 = 2.4. \end{aligned}$$

حال مدل مفروض را با وجود X_i تخمین زده و تغییرات توضیح داده نشده را به دست می‌آوریم. البته این محاسبات در مثال ۶-۲ انجام شده است، داریم $RSS = 1/5$.

مقدار کاهش در مجموع مربعات پسماند برابر است با $0/9 = 2/4 - 1/5$. علت این بهبود وجود متغیر توضیحی X_t در مدل است که توانسته قدرت توضیحی مدل را افزایش دهد. بدین ترتیب آماره F را به شرح زیر محاسبه می‌کنیم،

$$F = \frac{(RSS_t - RSS) / 1}{RSS / (n - k)}$$

$$= \frac{(2/4 - 1/5) / 1}{1/5 / (5 - 3)} = \frac{0/9}{0/75} = 1/2.$$

مقدار F در سطح معنی‌دار ۵ درصد برابر است با $F_{0.05}(1, 2) = 18/51$ بنابراین آماره آزمون معنی‌دار نیست؛ زیرا در ناحیه بحرانی قرار نگرفته و در نتیجه فرضیه H_0 رد نمی‌شود. قبلاً در مثال ۶-۲ که همین مسأله را با آماره t حل کردیم - مقدار t برابر $1/1$ - به دست آمد، در حالی که در اینجا مقدار F برابر $1/2$ شده است. بدیهی است چنین نتیجه‌ای نشان دهنده درستی محاسبات است؛ زیرا با توجه به معادله $2-51$ می‌دانیم

$$F = t^2,$$

بنابراین، $\sqrt{1/2} = 1/1$.

۶-۴ آنالیز واریانس و آزمون معنی‌دار بودن مدل رگرسیون
مدل رگرسیون زیر را ملاحظه کنید،

$$Y_t = \beta_1 + \beta_2 X_{1t} + \dots + \beta_k X_{kt} + U_t.$$

در این قسمت می‌خواهیم این فرضیه را آزمون کنیم که هیچکدام از متغیرهای توضیحی بر تغییرات Y_t تأثیری ندارد؛ یعنی

$$H_0: \beta_2 = \beta_3 = \dots = \beta_k = 0.$$

فرضیه مخالف، یعنی H_1 این است که $H_0 = 0$ صحیح نیست؛ یعنی حداقل یک متغیر

توضیحی وجود دارد که در تغییرات Y_t مؤثر واقع می‌شود؛ به عبارت دیگر، فرضیه H_1 بر این دلالت دارد که حداقل یک β وجود دارد که صفر نیست. اگر فرضیه H_0 صحیح باشد، تغییرات Y_t صرفاً تصادفی بوده، فقط تابعی از U_t خواهد بود، زیرا در چنین حالتی خواهیم داشت

$$Y_t = \beta_1 + U_t.$$

توجه به این نکته اهمیت فراوان دارد که آزمون «تمام پارامترها» برابر صفر غیر از آزمونهای جداگانه «هریک از پارامترها» برابر صفر است. این نکته را در ادامه مباحث این قسمت مطرح خواهیم کرد.

همان‌گونه که در فصل دوم قسمت ۲-۳ دیدیم، آزمون معنی دار بودن کل رگرسیون را به کمک تجزیه مجموع تغییرات Y_t ، یعنی $\sum y_t^2$ ، به تغییرات توضیح داده شده ($\sum \hat{y}_t^2$) و تغییرات توضیح داده نشده ($\sum e_t^2$) انجام می‌دهند. با توجه به معادله ۵-۴۵ می‌دانیم

$$y = \hat{y} + e.$$

یکی از اهداف مسأله تجزیه واریانس آزمون معنی دار بودن تغییرات توضیح داده شده (قدرت توضیحی مدل) است. در رگرسیون ساده، این آزمون دقیقاً همان $H_0: \beta = 0$ است؛ زیرا در این گونه مدلها فقط یک متغیر توضیحی وجود دارد، اما در رگرسیون چندمتغیره، به علت تعدد متغیرهای توضیحی، این گونه نیست و در ادامه بحث خواهیم دید که باید از آزمون F استفاده شود؛ بنابراین ملاحظه می‌شود که جدول تجزیه واریانس اطلاعاتی دارد که مستقیماً در آزمون معنی دار بودن رگرسیون مورد نیاز است.

مانند جدول ۲-۱ که تجزیه واریانس را برای رگرسیون ساده نشان می‌دهد کل تغییرات متغیر درون‌زا ($y \cdot y$)، تغییرات توضیح داده شده ($\hat{y} \cdot \hat{y}$) و تغییرات توضیح داده نشده یا مجموع مربعات پسماند ($e \cdot e$) را همراه با درجات آزادی و میانگین هرکدام در جدول زیر (۶-۱) تنظیم کرده، آن را جدول تجزیه واریانس برای رگرسیون چندمتغیره می‌نامیم.

جدول ۶-۱ آنالیز واریانس برای مدل رگرسیون چندمتغیره

منبع تغییرات	مجموع مربعات	درجات آزادی	میانگین مربعات با واریانس
تخمین رگرسیون	$ESS = \hat{y}'\hat{y} = \hat{\beta}'X'y$	$(k-1)$	$ESS / (k-1)$
پسماند	$RSS = e'e = y'y - \hat{\beta}'X'y$	$(n-k)$	$RSS / (n-k)$
کل تغییرات	$TSS = y'y$	$(n-1)$	$F = \frac{ESS / (k-1)}{RSS / (n-k)}$

برای تبیین بیشتر جدول تجزیه واریانس به توضیحات زیر اشاره می‌کنیم:
 الف) معادله $ESS = \hat{y}'\hat{y} = \hat{\beta}'X'y$ در واقع همان معادله ۵-۳۹ است که قبلاً ثابت شده است.

ب) برای نشان دادن درستی رابطه $RSS = e'e = y'y - \hat{\beta}'X'y$ ، ابتدا معادله ۵-۴۶ را می‌نویسیم،

$$e'e = y'y - y'X(X'X)^{-1}X'y.$$

با توجه به معادله ۵-۳۴ می‌دانیم

$$y'X(X'X)^{-1} = \hat{\beta}'$$

در نتیجه خواهیم داشت

$$e'e = y'y - \hat{\beta}'X'y \quad (۶-۳۵)$$

ج) برای تبیین این نکته که نسبت $\frac{ESS / (k-1)}{RSS / (n-k)}$ دارای توزیع F با $(k-1)$ و $(n-k)$ درجه آزادی است، ابتدا معادله ۲-۵۱ را یک بار دیگر می‌نویسیم،

$$F = \frac{ESS / 1}{RSS / (n-2)} \sim F(1, n-2).$$

می‌دانیم معادله فوق برای مدل رگرسیونی است که فقط دو پارامتر دارد. می‌توان مشابه این فرمول را برای یک مدل رگرسیون با k پارامتر نوشت. خواهیم داشت

$$F = \frac{\hat{y}'\hat{y} (k-1)}{e'e (n-k)} = \frac{ESS / (k-1)}{RSS / (n-k)} \sim F(k-1, n-k). \quad (۶-۳۶)$$

با وجود این، برای اینکه اشاره‌ای به نحوه استخراج این فرمول داشته باشیم، می‌گوییم آماره F نسبت دو توزیع مستقل χ^2 است که هر یک بر درجات آزادی خود تقسیم شده‌اند. با توجه به تعریف توزیع χ^2 در معادله ۲-۳۷ می‌توان گفت که

$$\frac{(\hat{\beta} - \beta)' (\hat{\beta} - \beta)}{\text{Var}(\hat{\beta})} \sim \chi^2 (k-1). \quad (۲-۳۷)$$

یادآوری می‌کنیم که فرض ما، در این قسمت این است که تمام محاسبات بر حسب انحراف از میانگین انجام شده است؛ بنابراین بردار β در رابطه (۲-۳۶) شامل $(k-1)$ پارامتر است. همچنین مشابه معادله ۲-۳۸، می‌توان نشان داد که

$$\frac{e'e}{\sigma_u^2} \sim \chi^2 (n-k) \quad (۲-۳۸)$$

اگر آماره‌های موجود در رابطه‌های ۲-۳۷ و ۲-۳۸ را بر درجات آزادی خود تقسیم کنیم، نسبت آنها به یکدیگر توزیع F خواهد داشت. با توجه به $\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$ ، داریم

$$\frac{[(\hat{\beta} - \beta)' / \sigma^2 (X'X)^{-1}] / (k-1)}{[e'e / \sigma^2] / (n-k)} \sim F(k-1, n-k).$$

اگر فرضیه H_0 صحیح باشد، یعنی $\beta_1 = \beta_2 = \dots = \beta_k = 0$ ، آنگاه $\beta = 0$ ، در نتیجه خواهیم داشت

$$\frac{\hat{\beta}' (X'X) \hat{\beta} / (k-1)}{e'e / (n-k)} \sim F(k-1, n-k),$$

$$\frac{[\hat{\beta}' (X'X) \hat{\beta}] / (k-1)}{e'e / (n-k)} \sim F(k-1, n-k),$$

با استفاده از معادله ۵-۳۵، یعنی $\hat{y} = X\hat{\beta}$ ، داریم

$$\hat{y}' \hat{y} = \hat{\beta}' X' X \hat{\beta},$$

و با جایگزینی معادله فوق در آماره F خواهیم داشت

$$F = \frac{\hat{\mathbf{y}}\hat{\mathbf{y}}'(k-1)}{\mathbf{e}'\mathbf{e}(n-k)} = \frac{ESS/(k-1)}{RSS/(n-k)} \sim F(k-1, n-k),$$

که دقیقاً همان رابطه ۶-۳۶ است. یادآوری می‌شود که این رابطه، براساس $\beta = 0$ به دست آمده و بنابراین موقعی برقرار است که H_0 صحیح باشد. برای آزمون درستی H_0 ، کافی است آماره آزمون را که از رابطه ۶-۳۶ به دست می‌آید با مقدار جدول F و در سطح معنی دار α درصد مقایسه کنیم. اگر آماره آزمون معنی دار بود؛ یعنی در ناحیه بحرانی قرار گرفت، H_0 رد می‌شود.

مثال ۶-۶ مدل رگرسیون زیر را که موضوع مثال ۶-۲ بوده است دوباره ملاحظه کنید،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

می‌خواهیم فرضیه $\beta = \gamma = 0$ را در سطح معنی دار ۵ درصد آزمون کنیم. آماره آزمون را با استفاده از رابطه ۶-۳۶ می‌نویسیم،

$$\frac{ESS/(k-1)}{RSS/(n-k)}$$

با استفاده از اطلاعات موجود در مثال ۶-۲ می‌دانیم $26/5 =$ تغییرات توضیح داده شده و $1/5 =$ تغییرات توضیح داده نشده، $n=5$ ، $k=3$ ، در نتیجه

$$F = \frac{26/5 / (3-1)}{1/5 / (5-3)} = 17/67$$

مقدار F را با درجات آزادی ۲ و ۲ در سطح معنی دار ۵ درصد از جدول به دست می‌آوریم،

$$F_{\alpha}(2, 2) = 19,$$

در نتیجه آماره آزمون کمتر از مقدار F در جدول است و معنی دار نیست؛ بنابراین فرضیه H_0 رد نمی‌شود.

۱. آماره آزمون F برحسب R^2

می توان رابطه ۶.۳۶ را فقط برحسب R^2 نوشت. در چنین حالتی محاسبات آزمون معنی دار بودن مدل رگرسیون بسیار ساده تر خواهد بود. با توجه به معادله ۵-۴۸ می توان چنین نوشت

$$RSS = e'e = (1 - R^2) y'y, \quad (6.39)$$

$$ESS = y'y = R^2 y'y \quad (6.40)$$

با جایگزینی معادله های ۶.۳۹ و ۶.۴۰ در معادله ۶.۳۶ داریم

$$F = \frac{[R^2 y'y] / (k-1)}{[(1-R^2) y'y] / (n-k)} \sim F(k-1, n-k),$$

بعد از ساده کردن خواهیم داشت

$$F = \left(\frac{n-k}{k-1} \right) \frac{R^2}{(1-R^2)} \sim F(k-1, n-k). \quad (6.41)$$

ملاحظه می شود که برای آزمون معنی دار بودن یک مدل رگرسیون، فقط کافی است مقدار R^2 را بدانیم.

فرمول ۶.۴۱ را می توان از راه دیگری نیز به دست آورد. آزمون صفر بودن زیرمجموعه ای از پارامترهای یک مدل رگرسیون را که در معادله ۶.۳۴ مطرح کردیم یک بار دیگر می نویسیم،

$$F = \frac{(n-k)}{r} \frac{R^2 - R_r^2}{1 - R^2} \sim F(r, n-k).$$

بحث مادر معادله فوق این بود که می خواستیم فرضیه صفر بودن r پارامتر را آزمون کنیم. اگر بخواهیم فرضیه صفر بودن تمام پارامترها، غیر از جمله ثابت را آزمون کنیم، یعنی

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0,$$

r برابر $(k-1)$ خواهد بود. از طرف دیگر، بر فرض درستی فرضیه H_0 ، می توان مدل

رگرسیون مقید را به صورت زیر نوشت

$$Y_i = \beta_1 + U_i.$$

چون هیچ متغیر برون‌زا برای توضیح دادن تغییرات Y_i وجود ندارد؛ بنابراین R^2 برابر صفر می‌شود. مقادیر $t = (k-1)$ و $R^2 = 0$ را در معادله ۶-۳۴ قرار می‌دهیم،

$$F = \left(\frac{n-k}{k-1} \right) \cdot \frac{R^2}{(1-R^2)} \sim F(k-1, n-k),$$

که دقیقاً همان رابطه ۶-۴۱ است.

مثال ۶-۷ مدل رگرسیون زیر را که موضوع مثال ۶-۶ بوده است ملاحظه کنید،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

فرضیه $H_0: \beta = \gamma = 0$ را در سطح معنی‌دار ۵ درصد آزمون کنید. می‌دانیم $R^2 = 0/9464$

با استفاده از مقدار R^2 ، آماره آزمون را از رابطه ۶-۴۱ محاسبه کرده و با مقدار به دست آمده از جدول F مقایسه می‌کنیم. خواهیم داشت

$$F = \left(\frac{n-k}{k-1} \right) \frac{R^2}{1-R^2} = \frac{0-3}{3-1} \left(\frac{0/9464}{1-0/9464} \right) = 17/6,$$

که با توجه به $F_{0.05}(2, 2) = 19$ معنی‌دار نیست و فرضیه H_0 رد نمی‌شود.

۲. آزمون معنی‌دار بودن کل مدل رگرسیون برای مدل‌های فاقد جمله ثابت مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + U_i.$$

می‌خواهیم فرضیه زیر را آزمون کنیم،

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0.$$

فرضیه مخالف این است که حداقل یکی از پارامترها صفر نیست. ملاحظه می‌شود که

مدل رگرسیون فاقد جمله ثابت است. از آزمون والد استفاده می‌کنیم. معادله ۶.۳۳ را دوباره می‌نویسیم،

$$F = \frac{(RSS_r - RSS_0) / r}{RSS_0 / (n - k)} \sim F(r, n - k)$$

با جایگزینی H_0 در مدل رگرسیون مفروض، مدل مقید به شرح زیر به دست می‌آید،

$$Y_i = U_i.$$

با توجه به اینکه در مدل مقید هیچ متغیر توضیحی وجود ندارد، مدل اساساً فاقد قدرت توضیحی بوده در نتیجه مجموع مربعات پسماند برابر مجموع تغییرات متغیر درون‌زا می‌شود؛ یعنی $\sum e_i^2 = \sum Y_i^2$ یا $RSS_r = \sum Y_i^2$. از طرف دیگر، با توجه به معادله ۵.۵۷ داریم

$$\sum Y_i^2 = \sum \hat{Y}_i^2 + \sum e_i^2.$$

بنابراین برای مدل غیرمقید می‌توان چنین نوشت

$$\sum e_i^2 = RSS_0 = \sum Y_i^2 - \sum \hat{Y}_i^2,$$

با استفاده از تساوی $RSS_r = \sum Y_i^2$ خواهیم داشت

$$RSS_0 - RSS_r = \sum \hat{Y}_i^2.$$

با جایگزینی رابطه فوق در صورت کسر آماره F داریم

$$F = \frac{\sum \hat{Y}_i^2 / (k - 1)}{RSS_0 / (n - k + 1)} = \frac{\sum \hat{Y}_i^2 / (k - 1)}{\sum e_i^2 / (n - k + 1)} \sim F(k - 1, n - k + 1),$$

اگر آماره مزبور معنی‌دار باشد، فرضیه H_0 رد می‌شود.

۳. مقایسه آزمون معنی دار بودن کل رگرسیون و معنی دار بودن هر یک از پارامترها*
مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + U_i.$$

دو فرضیه زیر را به ترتیب با آماره F و آماره t آزمون می کنیم،

$$H_0: \beta_2 = \beta_3 = \dots = \beta_k = 0$$

$$H_0: \beta_i = 0, \quad i = 2, 3, \dots, k.$$

آیا اگر یکی از این دو فرضیه قبول یا رد شود، فرضیه دیگر نیز ضرورتاً قبول یا رد خواهد شد؟

ابتدا پاسخ را مطرح کرده، سپس به اثبات آن می پردازیم. در جواب باید گفت که ضرورتاً این گونه نیست؛ زیرا آزمونهای F و t، از نظر استنباط آماری، قابل مقایسه با یکدیگر نیستند، چون بر فرضهای متفاوتی مبتنی هستند. برای مثال، وقتی بخواهیم $\beta_1 = 0$ را با آماره t آزمون کنیم، وجود هیچ فرضی برای β_2 ضروری نیست. در واقع با آماره t، آزمون هر یک از پارامترها مستقل از دیگری انجام می شود، در حالی که آزمون F مستلزم ملاحظه همزمان تمام پارامترهاست. برای اثبات این قضیه کافی است آماره F را بر حسب آماره t بنویسیم. بحث را در قالب یک مدل رگرسیون با دو متغیر توضیحی ادامه می دهیم.

مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + U_i.$$

سؤال این است که آیا فرضیه $H_0: \beta_2 = \beta_3 = 0$ و فرضیه $H_0: \beta_2 = 0$ یا $H_0: \beta_3 = 0$ ، که اولی با آماره F و دومی با آماره t آزمون می شود، ضرورتاً به یک نتیجه منتهی می گردد؟ برای این منظور ابتدا مقادیر تغییرات توضیح داده شده و تغییرات توضیح داده نشده را محاسبه کرده و سپس آنها را در معادله ۶-۳۶ قرار می دهیم تا رابطه F با t به دست آید. با توجه به

تعریف تغییرات توضیح داده شده و نیز معادله‌های ۱-۳۰ و ۴-۶ داریم

$$\begin{aligned} ESS &= \sum (\hat{Y}_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2, \\ &= \sum (\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} - \bar{Y})^2, \\ &= \sum [(\bar{Y} - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2) + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} - \bar{Y}]^2, \\ &= \sum (\hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i})^2. \end{aligned}$$

بدین ترتیب خواهیم داشت

$$ESS = \hat{\beta}_1^2 \sum x_{1i}^2 + \hat{\beta}_2^2 \sum x_{2i}^2 + 2\hat{\beta}_2 \hat{\beta}_3 \sum x_{2i} x_{3i}. \quad (6.42)$$

از معادله ۶-۲۱ استفاده کرده واریانس U_i را برای این حالت می‌نویسیم،

$$\hat{\sigma}_u^2 = \frac{e'e}{n-3} = \frac{RSS}{n-3}. \quad (6.43)$$

معادله‌های ۶-۴۲ و ۶-۴۳ را در آماره F در معادله ۶-۳۶ قرار می‌دهیم. خواهیم داشت

$$F = \frac{\hat{\beta}_1^2 \sum x_{1i}^2 + \hat{\beta}_2^2 \sum x_{2i}^2 + 2\hat{\beta}_2 \hat{\beta}_3 \sum x_{2i} x_{3i}}{2\hat{\sigma}^2}. \quad (6.44)$$

حال به بررسی آزمون t می‌پردازیم. آماره t برای آزمون $\beta_2 = 0$ و آماره t برای

آزمون $\beta_3 = 0$ به ترتیب برابر است با

$$t_2 = \frac{\hat{\beta}_2}{SE(\hat{\beta}_2)}, \quad t_3 = \frac{\hat{\beta}_3}{SE(\hat{\beta}_3)}.$$

مقادیر فوق را مجدور کرده داریم

$$t_2^2 = \frac{\hat{\beta}_2^2}{\text{Var}(\hat{\beta}_2)}, \quad t_3^2 = \frac{\hat{\beta}_3^2}{\text{Var}(\hat{\beta}_3)}. \quad (6.45)$$

با توجه به معادله‌های ۴-۲۴ و ۴-۲۵ می‌دانیم که

$$\text{Var}(\hat{\beta}_y) = \frac{\hat{\sigma}^2}{(1-r_{yy}^2) \sum x_{yt}^2}, \quad \text{Var}(\hat{\beta}_x) = \frac{\hat{\sigma}^2}{(1-r_{xx}^2) \sum x_{xt}^2}.$$

معادله‌های فوق را در معادله ۶-۴۵ قرار می‌دهیم، خواهیم داشت

$$t_y^2 = \frac{\hat{\beta}_y^2 (1-r_{yy}^2) \sum x_{yt}^2}{\hat{\sigma}^2}, \quad t_x^2 = \frac{\hat{\beta}_x^2 (1-r_{xx}^2) \sum x_{xt}^2}{\hat{\sigma}^2}.$$

مقادیر $\hat{\beta}_y$ و $\hat{\beta}_x$ را از معادله‌های فوق به دست می‌آوریم،

$$\hat{\beta}_y = \frac{t_y^2 \hat{\sigma}^2}{(1-r_{yy}^2) \sum x_{yt}^2}, \quad \hat{\beta}_x = \frac{t_x^2 \hat{\sigma}^2}{(1-r_{xx}^2) \sum x_{xt}^2}.$$

با استفاده از دو معادله فوق می‌توان $\hat{\beta}_y$ و $\hat{\beta}_x$ را به دست آورد،

$$\hat{\beta}_y \hat{\beta}_x = \frac{t_y t_x \hat{\sigma}^2}{(1-r_{yy}^2) \sqrt{\sum x_{yt}^2 \sum x_{xt}^2}}.$$

مقادیر $\hat{\beta}_y$ ، $\hat{\beta}_x$ و $\hat{\beta}_y \hat{\beta}_x$ را در معادله ۶-۴۴ جایگزین می‌کنیم. در نتیجه

$$F = \frac{t_y^2 + t_x^2 + 2(t_y t_x \sum x_{yt} x_{xt}) / \sqrt{\sum x_{yt}^2 \sum x_{xt}^2}}{2(1-r_{yy}^2)},$$

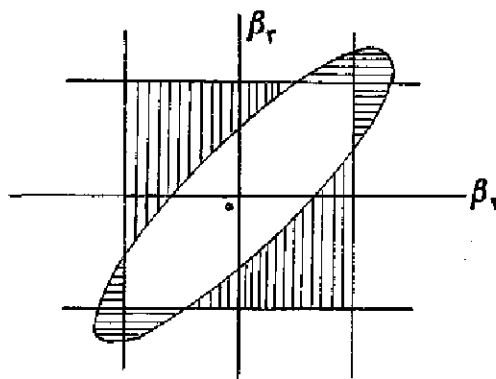
و با توجه به معادله ۱-۳۲ داریم

$$F = \frac{t_y^2 + t_x^2 + 2t_y t_x r_{yx}}{2(1-r_{yy}^2)} \sim F(2, n-3) \quad (6.46)$$

فرض می‌کنیم که مقادیر t_y و t_x کوچک هستند، به گونه‌ای که از مقادیر t در جدول کمتر شده، معنی دار نمی‌شوند. در نتیجه فرضیه‌های $H_0: \beta_y = 0$ و $H_0: \beta_x = 0$ هر دو قبول خواهد شد. اما اگر متغیرهای توضیحی مدل، یعنی X_{yt} و X_{xt} ، از شدت همبستگی تقریباً زیادی برخوردار باشد، یعنی مقدار r_{yx} به یک نزدیک شود، آنگاه $(1-r_{yy}^2)$ عدد بسیار

کوچکی شده، بنابراین کسر F بسیار بزرگ خواهد شد، به گونه‌ای که به راحتی از مقدار جدول F تجاوز کرده، فرضیه $H_0: \beta_1 = \beta_2 = 0$ قبول می‌شود. به همین ترتیب می‌توان گفت که در مواردی ممکن است مقادیر $\hat{\beta}_1$ یا $\hat{\beta}_2$ یا هر دو به طور معنی‌داری با صفر متفاوت باشد، در حالی که آماره F ، معنی‌دار نیست. علت این امر تفاوت ماهوی آزمونهای F و t است. وقتی بخواهیم فرضیه $H_0: \beta_1 = \beta_2 = 0$ را با آماره t آزمون کنیم، هیچ ضرورتی ندارد که فرضیه‌ای درباره β_1 مطرح کنیم. اساساً آزمونهای t برای پارامترهای مختلف مستقل از یکدیگر انجام می‌شود، در حالی که آزمون F براساس فرضیه‌ای شکل می‌گیرد که شامل تمام پارامترهای متغیرهای توضیحی است.

تفاوت در استنباط آماری حاصل از آزمونهای F و t را می‌توان در یک نمودار نیز نشان داد. وقتی بخواهیم فرضیه $H_0: \beta_1 = \beta_2 = 0$ را آزمون کنیم در واقع به جای یک فاصله اطمینان، باید از یک «ناحیه اطمینان» صحبت کرد؛ یعنی ناحیه‌ای که اگر آماره آزمون F در آن واقع شود، معنی‌دار نیست و H_0 قبول می‌شود. در آزمون t نیز دو فاصله اطمینان برای β_1 و β_2 به دست می‌آوریم که در نتیجه یک ناحیه اطمینان به وجود می‌آورند. همان گونه که در نمودار ۶-۱ ملاحظه می‌شود قسمتهایی وجود دارد که این دو ناحیه اطمینان بر هم منطبق نیست. در بند (۶) همین قسمت نشان خواهیم داد که اگر



نمودار ۶-۱ نواحی اطمینان برای آزمونهای F و t

بخواهیم با آماره F ، دو فرضیه را همزمان آزمون کنیم، ناحیه اطمینان، یک بیضی خواهد بود. همچنین می توان نشان داد که هر قدر همبستگی دو متغیر توضیحی در مدل رگرسیون بیشتر باشد، این بیضی باریکتر خواهد بود و برعکس. در نمودار فوق، قسمتهایی که با هاشور عمودی مشخص شده است، ناحیه اطمینان برای $\beta_1 = 0$ و $\beta_2 = 0$ است که از آماره t به دست می آید. ملاحظه می شود که فقط ناحیه هاشور نخورده داخل بیضی، فصل مشترک نواحی اطمینان آماره های F و t است. اگر آماره F در قسمتهایی قرار گیرد که با هاشورهای افقی مشخص شده است، فرضیه H_0 قبول می شود، در حالی که دقیقاً در همین نواحی فرضیه H_1 با آزمون t رد خواهد شد. به همین ترتیب در نواحی مشخص شده با هاشورهای عمودی، فرضیه H_0 با آزمون t قبول شده و با آزمون F رد می شود؛ بنابراین، با افزایش همبستگی بین متغیرهای توضیحی، نواحی هاشور زده شده، بزرگتر می شود در نتیجه تفاوت نتایج حاصل از آزمونهای F و t بزرگتر خواهد شد.

مثال ۶-۸ مدل رگرسیون زیر مفروض است،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + U_i$$

با استفاده از ۲۰ مشاهده، کمتهای زیر را بر حسب انحراف از میانگین محاسبه کرده ایم.

$$\begin{aligned} \sum y_i^2 &= 100 & \sum x_{2i}^2 &= 100 \\ \sum y_i x_{2i} &= 90 & \sum x_{2i}^2 &= 100 \\ \sum y_i x_{3i} &= 90 & \sum x_{2i} x_{3i} &= 90 \end{aligned}$$

اولاً، پارامترهای β_2 و β_3 را تخمین بزنید.

ثانیاً، با استفاده از آزمون t نشان دهید که هر یک از دو فرضیه $\beta_1 = 0$ و $\beta_2 = 0$ رد نمی شود.

ثالثاً، با استفاده از آزمون F نشان دهید دو فرضیه $\beta_1 = 0$ و $\beta_2 = 0$ به طور همزمان رد می شود؛ یعنی: $H_0: \beta_1 = \beta_2 = 0$.

۱. از معادله‌های ۴-۱۰ داریم

$$\hat{\beta}_1 = \frac{\sum x_{rt}^2 \sum x_{rt} y_t - \sum x_{rt} x_{rt} \sum x_{rt} y_t}{\sum x_{rt}^2 \sum x_{rt}^2 - (\sum x_{rt} \sum x_{rt})^2}$$

$$\hat{\beta}_2 = \frac{\sum x_{rt}^2 \sum x_{rt} y_t - \sum x_{rt} x_{rt} \sum x_{rt} y_t}{\sum x_{rt}^2 \sum x_{rt}^2 - (\sum x_{rt} \sum x_{rt})^2}$$

با استفاده از محاسبات انجام شده، خواهیم داشت

$$\hat{\beta}_1 = \frac{100(90) - 90(90)}{100(100) - (90)^2} = \frac{100}{970} = 0/1031,$$

$$\hat{\beta}_2 = \frac{100(90) - 90(90)}{100(100) - (90)^2} = \frac{100}{970} = 0/1031.$$

۲. برای آزمون t باید واریانس $\hat{\beta}_1$ و $\hat{\beta}_2$ را به دست آوریم. ابتدا باید تغییرات توضیح داده نشده را به دست آورد، که خود مستلزم محاسبه تغییرات توضیح داده شده است. از معادله ۵-۳۹ داریم

$$\begin{aligned} ESS &= \hat{y}' \hat{y} = \hat{\beta}' X' y, \\ &= [\hat{\beta}_1 \quad \hat{\beta}_2] \begin{bmatrix} \sum x_{rt} y_t \\ \sum x_{rt} y_t \end{bmatrix} = \hat{\beta}_1 \sum x_{rt} y_t + \hat{\beta}_2 \sum x_{rt} y_t, \end{aligned}$$

در نتیجه

$$ESS = 0/1031(90) + 0/1031(90) = 83$$

از صورت مسأله می‌دانیم که: $\sum y_t^2 = y'y = TSS = 100$

بنابراین

$$RSS = TSS - ESS = 100 - 83 = 17.$$

از معادله ۶-۲۱ مقدار واریانس U_1 را حساب می‌کنیم، خواهیم داشت

$$\hat{\sigma}_{U_1}^2 = \frac{e'e}{n-k} = \frac{17}{20-2} = 1.$$

از معادله‌های ۴-۲۴ و ۴-۲۵ داریم

$$\text{Var}(\hat{\beta}_r) = \frac{\hat{\sigma}_u^2}{(1-r_{rr}^2) \sum x_{rt}^2}, \quad \text{Var}(\hat{\beta}_r) = \frac{\hat{\sigma}_u^2}{(1-r_{rr}^2) \sum x_{rt}^2}$$

ابتدا باید r_{rr}^2 را محاسبه کرد

$$r_{rr}^2 = \frac{(\sum x_{rt} x_{rt})^2}{\sum x_{rt}^2 \sum x_{rt}^2} = \frac{(90)^2}{100(100)} = 0/9025$$

بنابراین

$$\text{Var}(\hat{\beta}_r) = \frac{1(100)^2}{(100^2 - 90^2) 100} = \frac{100}{970} = 0/102064$$

$$\text{Var}(\hat{\beta}_r) = \frac{1(100)^2}{(100^2 - 90^2) 100} = \frac{100}{970} = 0/102064$$

آماره آزمون t عبارت است از

$$t_r = \frac{\hat{\beta}_r}{\text{SE}(\hat{\beta}_r)} = \frac{0/461}{\sqrt{0/102064}} = 1/444$$

$$t_r = \frac{\hat{\beta}_r}{\text{SE}(\hat{\beta}_r)} = \frac{0/461}{\sqrt{0/102064}} = 1/444$$

مقدار به دست آمده از جدول t با ۱۷ درجه آزادی و در سطح معنی‌دار ۵ درصد برابر است با ۲/۱۱۰، در نتیجه آماره آزمون معنی‌دار نیست و فرضیه‌های $\beta_r = 0$ و $\beta_r = 0$ رد نمی‌شود.

۳. برای آزمون F از آماره ۶-۳۶ استفاده می‌کنیم،

$$F = \frac{\text{ESS} / (k-1)}{\text{RSS} / (n-k)}$$

$$F = \frac{13 / (3-1)}{17(20-3)} = 41/5$$

مقدار به دست آمده از جدول F با درجات آزادی ۲ و ۱۷ و در سطح معنی دار ۵ درصد برابر است با

$$F_{\%5}(2, 17) = 3/59,$$

بنابراین آماره آزمون از مقدار F در جدول به مراتب بیشتر شده، در ناحیه بحرانی قرار می‌گیرد، نتیجه می‌گیریم که فرضیه $H_0: \beta_1 = \beta_2 = 0$ رد می‌شود.

بنابراین ملاحظه می‌شود که با آزمون t هر دو فرضیه قبول می‌شود. در حالی که با آزمون F هر دو فرضیه رد خواهد شد. علت این امر را می‌توان این گونه توضیح داد که تأثیر مستقل X_{21} و X_{22} روی Y_1 بسیار ضعیف است، در حالی که تأثیر مشترک آنها بسیار قوی است. البته دلیل اشتراک قدرت توضیحی بسیار X_{21} و X_{22} را می‌توان در ضریب همبستگی بسیار قوی آنها جستجو کرد. می‌دانیم $r_{23}^2 = 0/9025$ در نتیجه $r_{23} = 0/95$. این نکته را در بحث از فرمول ۶-۴۶ بیان کردیم. برای تبیین بیشتر این نکته تأثیر مستقل X_{21} و X_{22} و تأثیر مشترک آنها را در تغییرات توضیح داده شده Y_1 محاسبه می‌کنیم. معادله ۶-۴۲ جواب این سؤال است. این معادله را یک بار دیگر می‌نویسیم،

$$ESS = \hat{\beta}_1^2 \sum x_{21}^2 + \hat{\beta}_2^2 \sum x_{22}^2 + 2\hat{\beta}_1 \hat{\beta}_2 \sum x_{21} x_{22}.$$

دو جمله اول سمت راست، تأثیرات مستقیم و مستقل X_{21} و X_{22} روی Y_1 را اندازه‌گیری می‌کند و جمله سوم تأثیر مشترک را منعکس می‌سازد. با اطلاعاتی که داریم، هر یک را محاسبه می‌کنیم،

$$\hat{\beta}_1^2 \sum x_{21}^2 = (0/461)^2 (100) = 21/3,$$

$$\hat{\beta}_2^2 \sum x_{22}^2 = (0/461)^2 (100) = 21/3,$$

$$\frac{2\hat{\beta}_1 \hat{\beta}_2 \sum x_{21} x_{22}}{\sum \hat{y}_i^2} = \frac{2(0/461)(0/461)(95)}{83} = 40/4$$

ملاحظه می‌شود که تأثیر مشترک X_{21} و X_{22} روی Y_1 تقریباً دو برابر تأثیر مستقل هر یک

از آنهاست. همان گونه که قبلاً اشاره شد این امر می تواند به علت همبستگی بسیار شدید X_{1i} و X_{2i} در مدل رگرسیون مفروض باشد.

۶.۵ آزمون یک ترکیب خطی از پارامترها

آزمون محدودیتهای خطی بین پارامترها، کاربرد وسیعی در تحلیلهای اقتصادسنجی دارد. فرض کنید در مدل

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i,$$

می خواهیم ترکیب خطی زیر را بین پارامترها آزمون کنیم،

$$H_0: p\beta_2 + q\beta_3 = d.$$

روش کار دقیقاً مانند آزمون $\beta_1 = a$ است. همان گونه که $\hat{\beta}_1$ را استاندارد کرده، آماره آزمون را محاسبه و با مقدار t از جدول مقایسه می کردیم، در اینجا نیز باید $(p\hat{\beta}_2 + q\hat{\beta}_3)$ را استاندارد کرده، آماره آزمون را به دست آوریم. کافی است $(p\hat{\beta}_2 + q\hat{\beta}_3)$ را از میانگین خود کم کرده، بر انحراف معیارش تقسیم کنیم، خواهیم داشت

$$t = \frac{(p\hat{\beta}_2 + q\hat{\beta}_3) - E(p\hat{\beta}_2 + q\hat{\beta}_3)}{\sqrt{\text{Var}(p\hat{\beta}_2 + q\hat{\beta}_3)}} \sim t(n-k),$$

رابطه فوق را می توان چنین نوشت،

$$t = \frac{(p\hat{\beta}_2 + q\hat{\beta}_3) - (p\beta_2 + q\beta_3)}{p^2 \text{Var} \hat{\beta}_2 + q^2 \text{Var} \hat{\beta}_3 + 2pq \text{Cov}(\hat{\beta}_2, \hat{\beta}_3)} \sim t(n-k).$$

با داشتن تخمینهای $\hat{\beta}_2$ و $\hat{\beta}_3$ و یا توجه به اینکه $(p\beta_2 + q\beta_3) = d$ و نیز محاسبه واریانس $\hat{\beta}_2$ و واریانس $\hat{\beta}_3$ و کوواریانس $\hat{\beta}_2$ و $\hat{\beta}_3$ ، آماره آزمون به سهولت محاسبه می شود. اگر آماره آزمون معنی دار باشد، یعنی از مقدار جدول t بزرگتر شود و در ناحیه بحرانی قرار بگیرد، فرضیه H_0 رد می شود، در غیر این صورت قبول خواهد شد.

مثال ۶-۹ مدل رگرسیون زیر را که در مثال ۶-۲ مطرح کردیم، یک بار دیگر ملاحظه کنید،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

می دانیم

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix},$$

و نیز $\hat{\sigma}^2 = 0/75$ ، $\text{Var}(\hat{\beta}) = 0/75$ و $\text{Var}(\hat{\gamma}) = 1/875$ ، $n = 5$. همچنین می دانیم

$$\hat{Y}_i = \xi + 2/5 X_i - 1/5 Z_i.$$

می خواهیم فرضیه $H_0: 2\beta - 3\gamma = 10$ را در مقابل فرضیه $H_1: \neq 0$ در سطح معنی دار α درصد آزمون کنیم.

برای این منظور آماره آزمون را محاسبه می کنیم. می دانیم

$$\text{Var}(2\hat{\beta} - 3\hat{\gamma}) = 4 \text{Var}(\hat{\beta}) + 9 \text{Var}(\hat{\gamma}) - 12 \text{Cov}(\hat{\beta}, \hat{\gamma}).$$

با توجه به $\text{Cov}(\hat{\beta}, \hat{\gamma}) = 0/75(-1/5) = -1/125$ ، خواهیم داشت

$$\text{Var}(2\hat{\beta} - 3\hat{\gamma}) = 4(0/75) + 9(-1/875) + 12(1/125) = 33/375,$$

در نتیجه

$$t = \frac{(2\hat{\beta} - 3\hat{\gamma}) - (2\beta - 3\gamma)}{\text{SE}(2\hat{\beta} - 3\hat{\gamma})} = \frac{2(2/5) - 3(-1/5) - 10}{\sqrt{33/375}} = \frac{-0/5}{0/771} = -0/086.$$

با توجه به مقدار جدول t در سطح معنی دار $0/025$ ، یعنی $t = 4/303$ ، ملاحظه می شود که آماره آزمون معنی دار نیست و H_0 رد نمی شود.

مثال ۶-۱۰ برای مثال ۶-۹ فرضیه زیر را آزمون کنید،

$$H_0: \beta + \gamma = 1$$

یادآوری می‌کنیم که اگر مدل مفروض در مثال ۶-۹ بیان لگاریتمی از یک تابع تولید گاب-داگلاس به صورت $Q_i = AK^\alpha L^\beta \varepsilon_i$ باشد، فرضیه H_0 همان بازده ثابت نسبت به مقیاس است. آماره آزمون در این حالت برابر است با

$$t = \frac{(\hat{\beta} + \hat{\gamma}) - (\beta + \gamma)}{\sqrt{\text{Var}(\hat{\beta}) + \text{Var}(\hat{\gamma}) + 2 \text{Cov}(\hat{\beta}, \hat{\gamma})}}$$

$$t = \frac{2/5 + (-1/5) - 1}{\sqrt{0/75 + 1/875 - 2(0/75)(-1/5)}} = 0,$$

که با توجه به مقدار جدول t ، یعنی $4/303$ معنی‌دار نیست و H_0 قبول می‌شود.

۱. آزمون $C'\beta = r$

حالت عمومی آزمون یک محدودیت خطی در مدل رگرسیون $Y = X\beta + u$ را می‌توان به صورت زیر نشان داد،

$$H_0: C'\beta = r \quad (6-47)$$

ابتدا بردارهای $C \rightarrow (1 \times k)$ و $r \rightarrow (1 \times 1)$ را تعریف می‌کنیم. برای این منظور، $C'\beta = r$ را برای مثالهای ۶-۹ و ۶-۱۰ می‌نویسیم. فرضیه $H_0: 2\beta - 3\gamma = 10$ در مثال ۶-۱۰ به زبان معادله ۶-۴۷ عبارت است از

$$[2 \quad -3] \begin{bmatrix} \beta \\ \gamma \end{bmatrix} = [10] , \quad (6-48)$$

بنابراین بردار C' برابر است با $[2 \quad -3]$ و نیز یک اسکالر یک در یک بوده، برابر $[10]$ است. برای مثال ۶-۱۰، فرضیه $H_0: \beta + \gamma = 1$ به زبان ماتریسی برابر است با

$$[1 \quad 1] \begin{bmatrix} \beta \\ \gamma \end{bmatrix} = [1] ,$$

در نتیجه $C' = [1 \quad 1]$ و $r = [1]$.

در نوشتن بردار C' باید دقت شود که علاوه بر ضرایب پارامترهایی که در H_0 وجود دارد، باید برای سایر پارامترهای مدل نیز ضریب صفر را منظور کرد تا بدین ترتیب C' همواره یک بردار $(1 \times k)$ شود. برای مثال، اگر بخواهیم در مدل رگرسیون زیر

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + \beta_0 X_{0i} + U_i$$

فرضیه $10: 2\beta_2 - 3\beta_3 = 10$ را H_1 را آزمون کنیم، معادله $C'\beta = r$ عبارت است از

$$\begin{bmatrix} 2 & -3 & 0 & 0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_0 \end{bmatrix} = [10]$$

یادآوری می‌کنیم که پارامتر β_1 در C' وارد نشده؛ زیرا فرض بر این است که $(X'X)$ براساس مقادیر انحراف از میانگین محاسبه شده است. اما اگر در محاسبه $(X'X)$ از مشاهدات اصلی استفاده کنیم، ضرورتاً پارامتر β_1 نیز در C' ظاهر خواهد شد. در حالت کلی، اگر ماتریس $(X'X)$ براساس مقادیر اصلی محاسبه شده باشد، C' یک بردار $(1 \times k)$ است زیرا $(X'X) \rightarrow (n \times k)$ ؛ در غیر این صورت هرگاه $(X'X) \rightarrow (k-1) \times (k-1)$ ، آنگاه $(1 \times k-1) \rightarrow C'$ ، یعنی C' فاقد جمله ثابت مدل رگرسیون خواهد بود.

بعد از آشنایی با $C'\beta = r$ ، به این نکته اشاره می‌کنیم که برای آزمون H_0 در معادله ۶-۴۷ باید $C'\beta = r$ را استاندارد کرده، آماره آزمون به دست آمده را با مقدار جدول t مقایسه کنیم؛ بنابراین ابتدا میانگین و واریانس $C'\hat{\beta}$ باید محاسبه شود،

$$E(C'\hat{\beta}) = C'E(\hat{\beta}) = C'\beta,$$

$$\text{Var}(C'\hat{\beta}) = E[C'(\hat{\beta}-\beta)(\hat{\beta}-\beta)'C],$$

$$= C'E[(\hat{\beta}-\beta)(\hat{\beta}-\beta)']C.$$

با توجه به معادله ۶.۷ می توان نوشت،

$$\text{Var}(\mathbf{C}'\hat{\beta}) = \sigma^2 \mathbf{C}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C} \quad (6.49)$$

بدین ترتیب آماره آزمون عبارت خواهد بود از

$$t = \frac{\mathbf{C}'\hat{\beta} - \mathbf{C}'\beta}{\sqrt{\sigma^2 \mathbf{C}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}}} \sim t(n-k),$$

$$t = \frac{\mathbf{C}'\hat{\beta} - \tau}{\sigma \sqrt{\mathbf{C}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}}} \sim t(n-k). \quad (6.50)$$

اگر آماره آزمون فوق معنی دار بود؛ یعنی از مقدار موجود در جدول t بزرگتر شد، فرضیه H_0 رد می شود.

مثال ۶.۱۱ فرضیه $H_0: 2\beta - 3\beta = 10$ را که موضوع آزمون مثال ۶.۹ بود به زبان ماتریسی و با استفاده از رابطه ۶.۵۰ آزمون کنید.

فرضیه H_1 را قبلاً در معادله ۶.۴۸ به زبان ماتریسی و به صورت $\mathbf{C}'\beta$ نوشته ایم. $\mathbf{C}'\hat{\beta}$ نیز به همین ترتیب به دست می آید،

$$\mathbf{C}'\hat{\beta} = [2 \quad -3] \begin{bmatrix} 2/5 \\ -1/5 \end{bmatrix} = 9/5.$$

برای محاسبه واریانس $\mathbf{C}'\hat{\beta}$ از معادله ۶.۴۹ استفاده می کنیم،

$$\text{Var}(\mathbf{C}'\hat{\beta}) = \hat{\sigma}^2 \mathbf{C}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C},$$

$$= 0.75 [2 \quad -3] \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \end{bmatrix} = 33/375.$$

آماره آزمون از معادله ۶.۵۰ به دست می آید،

$$t = \frac{9/5 - 10}{\sqrt{33/375}} = \frac{-0.5}{0.771} = -0.086.$$

با توجه به مقدار جدول t در سطح معنی دار 0.025 ، یعنی $t = 4/3.03$ ، ملاحظه می شود که آماره آزمون معنی دار نیست و فرضیه H_0 رد نمی شود.

سؤال این است که اگر در یک مدل رگرسیون چندمتغیره، دو یا چند پارامتر را در نظر بگیریم و نسبت به هر پارامتر فرضیه ای را مطرح نماییم و سپس بخواهیم این فرضیه ها را همزمان آزمون کنیم از چه روشی می توان استفاده کرد. با توجه به مطالبی که مطرح شده است به نظر می رسد که نمی توان از آزمون t استفاده کرد. در ادامه بحث خواهیم دید که آزمون F جواب این مسأله است.

۶-۶ آزمون همزمان پارامترها و تعیین نواحی اطمینان*

مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_2 X_{i1} + \dots + \beta_k X_{ik} + U_i.$$

فرض کنید می خواهیم فرضیه های زیر را به طور همزمان آزمون کنیم،

$$H_1: \beta_2 = a_1,$$

$$H_2: \beta_3 = a_2,$$

$$\vdots$$

$$H_r: \beta_r = a_r.$$

برای این منظور ابتدا فرض می کنیم که فرضیه های H_1 همه صحیح است. مقادیر β در فرضیه های H_1 را در مدل رگرسیون مفروض جایگزین می کنیم تا مدل رگرسیون مقید به دست آید. بدیهی است تعداد پارامترهایی که باید در مدل رگرسیون مقید تخمین زده شود، از تعداد پارامترهای مدل اصلی کمتر است. تفاوت این دو دقیقاً به اندازه تعداد فرضیه هایی است که می خواهیم آزمون کنیم. اگر مجموع مربعات پسماند را در مدل اصلی با RSS و در مدل مقید با RSS_p مشخص کنیم؛ مسأله مفروض دقیقاً «آزمون معنی دار بودن زیرمجموعه ای از پارامترها» می شود که در قسمت ۶-۳ مطرح شد. اگر تعداد فرضیه ها را r بنامیم، آماره آزمون عبارت خواهد بود از

$$\frac{(RSS_r - RSS) / r}{RSS / (n - k)} \sim F(r, n - k) \quad (6.51)$$

که همان فرمول ۶.۳۳ است.

مثال ۶-۱۲ تابع تولید زیر را در نظر می‌گیریم،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i,$$

که در آن Y_i ، X_{1i} و X_{2i} به ترتیب لگاریتم تولید، لگاریتم کار و لگاریتم سرمایه است. براساس یک نمونه شامل ۲۳ شرکت مختلف تولیدی و محاسبات انجام شده در مثال ۴-۲ تخمین زیر را به دست آورده‌ایم،

$$\hat{Y}_i = 4 + 0.7X_{1i} + 0.2X_{2i}.$$

همچنین از مثال ۴-۲ می‌دانیم که $\sum y_i^2 = 10$ ، $\sum x_{1i}^2 = 12$ و $\sum x_{2i}^2 = 10$ ، $\sum x_{1i}y_i = 10$ ، $RSS = 1/4$ می‌خواهیم دو فرضیه زیر را به طور همزمان در سطح معنی‌دار ۵ درصد آزمون کنیم،

$$H_0: \beta_1 = 1, \quad H_0: \beta_2 = 0,$$

$$H_1: \beta_1 \neq 1, \quad H_1: \beta_2 \neq 0.$$

برای اینکه بتوان آماره آزمون را از معادله ۶.۳۳ به دست آورد کافی است مجموع مربعات پسماند را برای مدل مقید، یعنی RSS_r محاسبه کنیم. برای این منظور می‌گوییم که اگر فرضیه‌های فوق صحیح باشد، مدل رگرسیون مقید به صورت زیر خواهد بود،

$$Y_i = \alpha + (1) X_{1i} + 0 (X_{2i}) + U_i,$$

$$Y_i = \alpha + X_{1i} + U_i.$$

برای به دست آوردن $\hat{\alpha}$ باید $\sum e_i^2$ را حداقل کنیم،

$$\begin{aligned} \frac{d \sum e_i^2}{d \hat{\alpha}} &= \frac{d \sum (Y_i - \hat{\alpha} - X_{1i})^2}{d \hat{\alpha}}, \\ &= 2(-1) \sum (Y_i - \hat{\alpha} - X_{1i}) = 0, \end{aligned}$$

در نتیجه خواهیم داشت

$$\hat{\alpha} = \bar{Y} - \bar{X}_1$$

با استفاده از تعریف RSS_r داریم

$$\begin{aligned} RSS_r &= \sum (Y_i - \hat{Y}_i)^2, \\ &= \sum [(Y_i - (\hat{\alpha} + X_{1i}))]^2, \\ &= \sum [(Y_i - (\bar{Y} - \bar{X}_1 + X_{1i}))] = \sum [(Y_i - \bar{Y}) - (X_{1i} - \bar{X}_1)]^2. \end{aligned}$$

بدین ترتیب نتیجه می‌گیریم که

$$\begin{aligned} RSS_r &= \sum y_i^2 + \sum x_{1i}^2 - 2 \sum x_{1i} y_i, \\ &= 10 + 12 - 20 = 2. \end{aligned}$$

معادله ۶.۵۱ را می‌نویسیم،

$$F = \frac{(RSS_r - RSS) / r}{RSS / (n - k)}$$

می‌دانیم تعداد محدودیتها که در واقع همان تعداد فرضیه‌هاست برابر است با $r = 2$. در نتیجه:

$$F = \frac{(2 - 1/4) / 2}{1/4 / (23 - 2)} = \frac{0.3}{0.07} = 4/3.$$

آماره آزمون فوق را باید با مقدار جدول F با درجات آزادی ۲۰ و ۲ و سطح معنی‌دار ۵ درصد آزمون کرد. می‌دانیم

$$F_{\%0} (2, 20) = 3/49,$$

در نتیجه آماره آزمون در ناحیه بحرانی قرار گرفته و معنی‌دار است؛ بنابراین، فرضیه‌های H_0 به طور همزمان رد می‌شوند.

ناحیه اطمینان برای چند پارامتر

همان گونه که در قسمت ۶-۴ و در توضیحات نمودار ۶-۱ اشاره کردیم، در رگرسیون چندمتغیره می توان از نواحی اطمینان نیز صحبت کرد. حالت عمومی در به دست آوردن این نواحی اطمینان در قسمت ۶-۸ بررسی می شود. در اینجا برای تبیین این مسأله، قلمرو بحث را به یک مدل رگرسیون با دو متغیر توضیحی محدود می کنیم. مدل مثال ۶-۱۲ را یک بار دیگر می نویسیم،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i.$$

در مثال ۴-۲ برای β_1 و β_2 فواصل اطمینان را به طور جداگانه محاسبه کردیم. در این قسمت می خواهیم ناحیه اطمینان برای هر دو پارامتر را در سطح معنی دار ۵ درصد به دست آوریم.

بدون اثبات، می گوئیم که در رگرسیونهایی که دو متغیر توضیحی به شرح فوق دارند، آماره زیر دارای توزیع F با درجات آزادی $(2, n-3)$ است.

$$F = \frac{1}{2\hat{\sigma}_U^2} [\sum x_{1i}^2 (\hat{\beta}_1 - \beta_1)^2 + 2 \sum x_{1i} x_{2i} (\hat{\beta}_1 - \beta_1)(\hat{\beta}_2 - \beta_2) + \sum x_{2i}^2 (\hat{\beta}_2 - \beta_2)^2]. \quad (6.52)$$

به کمک معادله ۶-۵۲ می توان اولاً، نواحی اطمینان برای پارامترهای β_1 و β_2 را به دست آورد؛ ثانیاً، فرضیه های مختلف برای β_1 و β_2 را به طور همزمان آزمون کرد. البته در همین قسمت دیدیم که به کمک معادله ۶-۵۱ می توان آزمونهای همزمان پارامترها را نیز انجام داد. در این قسمت نشان می دهیم که استفاده از معادله های ۶-۵۱ و ۶-۵۲ به یک جواب منتهی می شوند. بنابراین می توان گفت که اهمیت معادله ۶-۵۲ بیشتر در به دست آوردن نواحی اطمینان است.

مثال ۶-۱۳ برای به دست آوردن ناحیه اطمینان ۹۵ درصدی برای β_1 و β_2 در مثال ۶-۱۲، معادله ۶-۵۲ را به صورت زیر می نویسیم. توجه داریم که $F_{\%95}(2, 20) = 3/49$. همچنین می دانیم در ناحیه اطمینان، آماره آزمون باید کمتر یا مساوی مقدار F از جدول باشد.

$$\sum x_{1i}^2 (\hat{\beta}_1 - \beta_1)^2 + 2 \sum x_{1i} x_{2i} (\hat{\beta}_1 - \beta_1) (\hat{\beta}_2 - \beta_2) + \sum x_{2i}^2 (\hat{\beta}_2 - \beta_2)^2 \leq 3/49 (2\hat{\sigma}^2). \quad (6.53)$$

با استفاده از محاسبات مثال ۴.۲ می‌دانیم

$$\sum y_i^2 = 10, \quad \sum x_{1i}^2 = 12, \quad \sum x_{2i}^2 = 12, \quad \sum x_{1i} x_{2i} = 8, \quad \hat{\sigma}^2 = 0.07$$

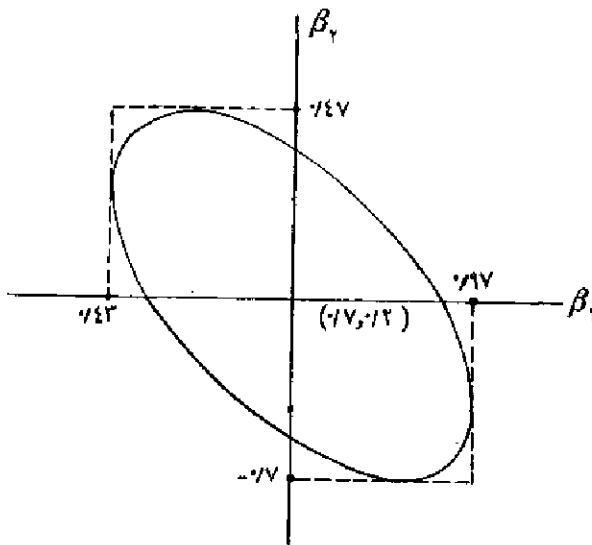
بنابراین معادله ۶.۵۳ به صورت زیر نوشته می‌شود،

$$12(0.07 - \beta_1)^2 + 2(8)(0.07 - \beta_1)(\beta_2 - \beta_2) + 12(0.07 - \beta_2)^2 \leq 3/49(2)(0.07).$$

دو طرف نامساوی فوق را بر ۱۲ تقسیم کرده، با تبدیل $(\hat{\beta}_1 - \beta_1)$ به $(\beta_1 - \hat{\beta}_1)$ و نیز $(\hat{\beta}_2 - \beta_2)$ به $(\beta_2 - \hat{\beta}_2)$ خواهیم داشت

$$(\beta_1 - 0.07)^2 + \frac{4}{3}(\beta_1 - 0.07)(\beta_2 - 0.07) + (\beta_2 - 0.07)^2 \leq 0.41. \quad (6.54)$$

ملاحظه می‌شود که منحنی نمودار تغییرات رابطه ۶.۵۴ یک بیضی است که مرکز آن نقطه $(0.07, 0.07)$ است. نمایش هندسی این بیضی در نمودار ۶.۲ منعکس است. ذکر دو نکته



نمودار ۶.۲ ناحیه اطمینان برای آزمونهای همزمان

در مورد این نمودار ضروری است. نکته اول این است که اگر $Cov(\hat{\beta}_1, \hat{\beta}_2) < 0$ ، آنگاه این بیضی به سمت چپ متمایل است. در مثال ۴-۲ نیز دیدیم که $Cov(\hat{\beta}_1, \hat{\beta}_2) = -0.007$ ؛ بنابراین در نمودار ۶-۲ جهت بیضی به سمت چپ است. اگر $Cov(\hat{\beta}_1, \hat{\beta}_2) > 0$ ، بیضی به سمت راست متمایل خواهد بود. با توجه به معادله ۴-۲۶ می توان گفت که علامت کوواریانس $\hat{\beta}_1, \hat{\beta}_2$ تابع معکوسی از علامت $\sum x_{1i} x_{2i}$ است. نتیجه می گیریم که در مدل های رگرسیون با دو متغیر توضیحی، هرگاه $\sum x_{1i} x_{2i}$ مثبت باشد، ناحیه اطمینان برای آزمون همزمان دو پارامتر به سمت چپ متمایل است و برعکس. دومین نکته در مورد این نمودار این است که ناحیه اطمینان ۹۵ درصد برای β_1, β_2 در واقع سطح محصور در بیضی است. این سطح با فواصل اطمینان ۹۵ درصد که با آزمون ۴ به طور جدا گانه برای β_1, β_2 به دست می آید متفاوت است.

برای تبیین این نکته مفید است که فاصله اطمینان ۹۵ درصد برای β_1, β_2 را که در مثال ۴-۲ به دست آورده ایم یک بار دیگر می نویسیم،

$$\beta_1 : (0.49, 0.91) ,$$

$$\beta_2 : (-0.01, 0.41) ,$$

این فاصله های اطمینان را می توان با قلمرو تغییرات β_1 و β_2 که در آزمون همزمان تعیین شده است و به راحتی از طریق مماسهای بیضی به دست می آید، مقایسه نمود،

$$\beta_1 : (0.43, 0.97) ,$$

$$\beta_2 : (-0.07, 0.47) ,$$

استفاده دیگری که می توان از معادله ۶-۵۲ کرد، آزمون همزمان فرضیه های مختلف در مورد پارامترهای β_1 و β_2 است. بدیهی است با داشتن نواحی اطمینان، به راحتی میتوان فرضیه های مختلف را به طور همزمان آزمون کرد. اگر فرضیه های مورد نظر در داخل ناحیه اطمینان قرار بگیرد، قبول و در غیر این صورت رد می شود. اما معمولاً تعیین نواحی اطمینان تا حدی مشکل است و به همین دلیل در این موارد

مستقیماً و به کمک معادله ۶-۵۲ از آزمون F استفاده می‌شود. می‌دانیم آماره آزمون در رابطه ۶-۵۲ دارای توزیع F با درجات آزادی $(3 - n, 2)$ است. مقدار عددی آماره آزمون را حساب می‌کنیم. مقادیر $H_0: \beta_1 = 0$ و $H_1: \beta_1 = 1$ را در ۶-۵۲ قرار داده خواهیم داشت

$$F = \frac{1}{2(0.07)} [12(0.07 - 1)^2 + 2(8)(0.07 - 1)(0.02 - 0) + 12(0.02 - 0)^2] = 4/3.$$

مقدار موجود در جدول F با درجات آزادی $2, 20$ و در سطح معنی‌دار 5 درصد برابر است با

$$F_{\%5}(2, 20) = 3/49.$$

با توجه به اینکه آماره آزمون در ناحیه بحرانی قرار می‌گیرد؛ معنی‌دار است و فرضیه‌های H_0 هر دو همزمان رد می‌شوند. یادآوری می‌شود که مقدار آماره آزمون به دست آمده از ۶-۵۲ دقیقاً برابر مقداری است که از معادله ۶-۵۱ به دست آوردیم. با توجه به سهولت معادله ۶-۵۱ ملاحظه می‌شود که آزمون همزمان پارامترها بهتر است به کمک معادله ۶-۵۱ صورت پذیرد. تعمیم نتایج به حالت کلی نیز با استفاده از این معادله بهتر انجام می‌شود. این نکته را در قسمت ۶-۸ نشان خواهیم داد.

۶-۷ آزمون تساوی چند پارامتر با یکدیگر

آزمون دیگری که می‌توان در مدل‌های رگرسیون چند متغیره مطرح کرد، آزمون تساوی دو یا بیشتر از دو پارامتر با یکدیگر است. برای مثال، مدل مصرف زیر را در نظر می‌گیریم،

$$C_t = \beta_1 + \beta_2 W_t + \beta_3 NW_t + \beta_4 C_{t-1} + U_t,$$

که در آن C_t مصرف، W_t درآمدهای حاصل از دستمزد و NW_t درآمدهای حاصل از منابع غیر دستمزد است. می‌توان این فرضیه را مطرح کرد که آیا میل نهایی به مصرف در افرادی که حقوق یا دستمزد می‌گیرند با میل نهایی به مصرف در صاحبان مشاغل آزاد که درآمدها از طریق حقوق و دستمزد نیست، یکسان است. در

این صورت باید فرضیه زیر را آزمون کرد،

$$H_0: \beta_r = \beta_s,$$

$$H_1: \beta_r \neq \beta_s.$$

در حالت کلی می توان یک مدل رگرسیون با k پارامتر را در نظر گرفت و فرضیه زیر را مطرح کرد،

$$H_0: \beta_i = \beta_j,$$

$$H_1: \beta_i \neq \beta_j.$$

دوره حل می توان در نظر گرفت. در راه حل اول کافی است که توزیع تفاوت $\hat{\beta}_i$ و $\hat{\beta}_j$ را ملاحظه کرده، فرضیه $H_0: (\hat{\beta}_i - \hat{\beta}_j) = 0$ را در مقابل فرضیه مخالف صفر آزمون کرد. بنابراین، باید ابتدا $\hat{\beta}_i - \hat{\beta}_j$ را استاندارد کرد تا آماره آزمون به دست آید. می دانیم

$$E(\hat{\beta}_i - \hat{\beta}_j) = \beta_i - \beta_j,$$

$$\text{Var}(\hat{\beta}_i - \hat{\beta}_j) = \text{Var}(\hat{\beta}_i) + \text{Var}(\hat{\beta}_j) - 2 \text{Cov}(\hat{\beta}_i, \hat{\beta}_j).$$

با تخمین واریانس جمله اختلال از معادله ۶۲۱ و تخمین واریانسهای $\hat{\beta}_i$ و $\hat{\beta}_j$ از معادله ۶۲۴ و نیز کوواریانس $\hat{\beta}_i, \hat{\beta}_j$ از معادله ۶۲۵، مقدار تخمین $\text{Var}(\hat{\beta}_i - \hat{\beta}_j)$ به دست می آید. بدین ترتیب آماره آزمون از رابطه زیر نتیجه خواهد شد،

$$t = \frac{(\hat{\beta}_i - \hat{\beta}_j) - E(\hat{\beta}_i - \hat{\beta}_j)}{\sqrt{\text{Var}(\hat{\beta}_i - \hat{\beta}_j)}} \\ = \frac{\hat{\beta}_i - \hat{\beta}_j - (\beta_i - \beta_j)}{\text{SE}(\hat{\beta}_i - \hat{\beta}_j)}$$

اگر فرضیه H_0 صحیح باشد، یعنی اگر $(\beta_i - \beta_j) = 0$ ، آنگاه آماره آزمون زیر دارای

توزیع t با $(n-k)$ درجه آزادی خواهد داشت،

$$t = \frac{\hat{\beta}_i - \beta_i}{SE(\hat{\beta}_i - \beta_i)} \sim t(n-k). \quad (6.55)$$

مقدار جدول t را با $(n-k)$ درجه آزادی و در سطح معنی دار $\alpha/2$ درصدیه دست می آوریم. اگر آماره آزمون (6.55) از این مقدار بیشتر شد، در آن صورت معنی دار بوده، فرضیه H_0 رد می شود، در غیر این صورت رد نخواهد شد.

مثال 6-14 مدل رگرسیون زیر، موضوع مثال 6-2 را دوباره می نویسیم،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

با استفاده از محاسبات انجام شده در مثال 6-1 فرضیه زیر را در سطح معنی دار 5 درصد آزمون کنید،

$$H_0: \beta = \gamma,$$

$$H_1: \beta \neq \gamma.$$

می دانیم $\hat{\beta} = 2/5$ ، $\hat{\gamma} = -1/5$ ، $\text{Var}(\hat{\beta}) = 0/75$ ، $\text{Var}(\hat{\gamma}) = 1/875$ و $\text{Cov}(\hat{\beta}, \hat{\gamma}) = 0/75(-1/5) = -1/125$ ، با این اطلاعات، آماره آزمون قابل محاسبه است.

$$(\hat{\beta} - \hat{\gamma}) = 2/5 - (-1/5) = 3/5,$$

$$\begin{aligned} \text{Var}(\hat{\beta} - \hat{\gamma}) &= \text{Var}(\hat{\beta}) + \text{Var}(\hat{\gamma}) - 2 \text{Cov}(\hat{\beta}, \hat{\gamma}), \\ &= 0/75 + 1/875 - 2(-1/125) = 4/875. \end{aligned}$$

از معادله 6.55 داریم

$$t = \frac{\hat{\beta} - \hat{\gamma}}{\sqrt{\text{Var}(\hat{\beta} - \hat{\gamma})}} \sim t(n-3),$$

$$t = \frac{3/5}{\sqrt{4/875}} = \frac{3}{2/2.8} = 1/811.$$

با توجه به اینکه آزمون دو طرفه است، مقدار جدول t با $2 = (n - 3)$ درجه آزادی و در سطح معنی دار ۲۵ درصد برابر است با $4/303$ ؛ بنابراین، آماره آزمون در ناحیه بحرانی قرار نگرفته، معنی دار نیست. بدین ترتیب فرضیه H_0 رد نمی شود.

آزمون تساوی بیش از دو پارامتر با یکدیگر

مدل رگرسیون چندمتغیره زیر مفروض است،

$$Y_t = \beta_1 + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + U_t.$$

می خواهیم فرضیه زیر را آزمون کنیم،

$$H_1: \beta_2 = \beta_3 = \beta_4.$$

فرضیه H_1 این است که حداقل یکی از تساویهای فوق برقرار نباشد. همان گونه که می دانیم در روال عمومی آزمون F ، فرض می شود، H_0 صحیح است؛ سپس مدل رگرسیون را مشروط به درستی H_0 می نویسیم و آن را مدل رگرسیون مقید می نامیم. بنابراین

$$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_2 X_{3t} + \beta_2 X_{4t} + \dots + \beta_k X_{kt} + U_t.$$

از β_2 فاکتور می گیریم،

$$Y_t = \beta_1 + \beta_2 (X_{2t} + X_{3t} + X_{4t}) + \dots + \beta_k X_{kt} + U_t,$$

یا

$$Y_t = \beta_1 + \beta_2 X_{2t}^* + \beta_3 X_{3t} + \dots + \beta_k X_{kt} + U_t,$$

که در آن

$$X_{2t}^* = (X_{2t} + X_{3t} + X_{4t}).$$

مجموع مربعات پسماند در مدل آخر را RSS_r می نامیم و با توجه به معادله های ۶.۳۳ یا ۶.۵۱ خواهیم داشت

$$F = \frac{(RSS_r - RSS) / r}{RSS / (n - k)} \sim F(r, n - k),$$

که در آن، RSS مجموع مربعات پسماند در مدل اولیه و r تعداد محدودیتهاست. باید توجه کنیم که در این مثال، $r=2$ است، نه ۳؛ زیرا اگر دو محدودیت برقرار باشد، مثلاً $\beta_1 = \beta_2$ و $\beta_3 = \beta_4$ ، محدودیت سوم، یعنی $\beta_1 = \beta_2 = \beta_3 = \beta_4$ نیز برقرار خواهد بود. به عبارت دیگر فقط می توان دو محدودیت مستقل در H_0 تشخیص داد. بنابراین، r تعداد محدودیتهای مشاهده شده نبوده، بلکه تعداد محدودیتهای مستقل است. اگر آماره آزمون فوق معنی دار باشد، فرضیه H_0 رد می شود، در غیر این صورت قبول خواهد شد.

مثال ۶-۱۵ در مثال ۶-۱۴ فرضیه $\beta = \gamma$ را در مقابل فرضیه $\beta \neq \gamma$ در سطح معنی دار ۵ درصد آزمون کنید.

گفتیم که فرضیه های تساوی چند پارامتر را باید با آماره F آزمون کرد؛ زیرا فرضیه های تساوی دو پارامتر را همواره می توان با آماره t آزمون کرد؛ با وجود این، با توجه به اینکه محاسبات مثال ۶-۱۴ کامل است، بنابراین، خیلی ساده تر خواهد بود اگر آزمون F را در مورد همین مثال بررسی کنیم. روش کار مهم است وگرنه تساوی دو یا بیشتر از دو پارامتر، از نظر روش آزمون تفاوتی ندارد. مدل رگرسیون را یک بار دیگر می نویسیم.

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i,$$

در مثال ۶-۲ دیدیم که مجموع مربعات پسماند برای این مدل برابر است با $RSS = 1/5$. حال باید مدل مقید را بسازیم. با جایگزینی فرضیه H_0 در این مدل خواهیم داشت

$$\begin{aligned} Y_i &= \alpha + \beta X_i + \beta Z_i + U_i \\ &= \alpha + \beta (X_i + Z_i) + U_i, \end{aligned}$$

یا با تعریف $X_i^* = (X_i + Z_i)$ ، داریم

$$Y_i = \alpha + \beta X_i^* + U_i.$$

باید این مدل مقید را تخمین زده، مجموع مربعات پسماند را به دست آوریم. با استفاده از

مثال ۵.۲ و جدول ۵.۳ محاسبات زیر را انجام می‌دهیم.

جدول ۶.۱

$x_i^* = x_i + z_i$	x_i^*	$x_i^* y_i$	x_i^{*2}	y_i^2
۸	۰	۰	۰	۱
۵	-۳	۹	۹	۹
۱۱	۳	۱۲	۹	۱۶
۶	-۲	۲	۴	۱
۱۰	۲	۲	۴	۱
Σ ۴۰	۰	۲۵	۲۶	۲۸

با مراجعه به معادله ۲.۴۷ می‌دانیم

$$\begin{aligned}
 RSS_r &= \sum e_{t(r)}^2 = \sum y_i^2 - \frac{(\sum x_i^* y_i)^2}{\sum x_i^{*2}} \\
 &= 28 - \frac{(25)^2}{26} = \frac{103}{26} = 3/96.
 \end{aligned}$$

معادله ۶.۵۱ را می‌نویسیم،

$$\begin{aligned}
 F &= \frac{(RSS_r - RSS) / r}{RSS / (n - k)} \sim F(r, n - k) \\
 &= \frac{(3/96 - 1/5) / 1}{(1/5)(5 - 2)} \sim F(1, 2).
 \end{aligned}$$

بنابراین مقدار آماره آزمون برابر است با $F = 3/28$. مقدار جدول F با درجات آزادی $(1, 2)$ و در سطح معنی‌دار ۵ درصد برابر است یا $18/51$ ؛ بنابراین، آماره آزمون معنی‌دار نبوده و H_0 رد نمی‌شود. همان‌گونه که انتظار می‌رفت دو آزمون F و t در این مورد خاص، که تساوی دو پارامتر را می‌خواهیم آزمون کنیم، دقیقاً به یک جواب می‌رسند؛ زیرا بنابر معادله ۲.۵۱ می‌دانیم $F = t^2$ و با توجه به $F = 3/28$ در مثال ۶.۱۴ و $t = 1/811$

در مثال ۶-۱۳ نتیجه می‌گیریم

$$\sqrt{3/28} = 1/811$$

۶-۸ آزمون همزمان چند ترکیب خطی از پارامترها: آماره عمومی آزمون و کاربرد آن*

عمومی‌ترین حالت در آزمون فرضیه این است که مجموعه‌ای از ترکیبهای خطی بین پارامترها را آزمون کنیم. ابتدا این آزمون را به زبان ماتریسی ارائه داده، سپس نشان می‌دهیم که چگونه بسیاری از آزمونهایی که تا به حال مطرح کرده‌ایم، می‌تواند حالت خاصی از این آزمون باشد.

مدل رگرسیون $y = X\beta + u$ مفروض است که در آن $y \rightarrow n \times 1$, $X \rightarrow n \times k$ و $u \rightarrow n \times 1$ می‌خواهیم مجموعه r فرضیه‌زیر را به‌طور همزمان آزمون کنیم،

$$H_0: R\beta = r, \quad (6.56)$$

که در آن $R \rightarrow r \times k$ و $r \leq k$ و $r \rightarrow r \times 1$. عناصر ماتریس R و بردار r مقادیر معلومی است. تعداد ردیفهای ماتریس R ، یعنی r در واقع تعداد محدودیتهای خطی است که می‌خواهیم آزمون کنیم؛ بنابراین r فرضیه برای آزمون داریم و بدیهی است که باید r از k کوچکتر یا مساوی آن باشد. توجه به این نکته نیز اهمیت دارد که ماتریس R باید دارای «رتبه کامل سطری»^۱ باشد؛ یعنی سطرهای آن از یکدیگر استقلال کامل داشته باشد که در واقع به معنی عدم همبستگی خطی بین فرضیه‌هاست. معادله ۶-۵۶ دارای اهمیت خاصی است، زیرا می‌توان فرضیه‌های مختلفی را در قالب آن بیان نمود. برای مثال، فرض کنید می‌خواهیم دو فرضیه زیر را به‌طور همزمان آزمون کنیم،

1. Full Row Rank

برای توضیح بیشتر به پیوست «۵-الف» مراجعه شود.

$$\begin{cases} H_2: 2\beta_1 - 3\beta_2 = 10, \\ H_1: 2\beta_1 - 3\beta_2 \neq 10, \end{cases} \quad \text{و} \quad \begin{cases} H_2: 2\beta_1 - 4\beta_0 = 2, \\ H_1: 2\beta_1 - 4\beta_0 \neq 2. \end{cases}$$

معادله ۶۵۶ به صورت زیر خواهد بود،

$$\begin{bmatrix} 0 & 2 & -3 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & -4 & \dots & 0 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \vdots \\ \beta_k \end{bmatrix} = \begin{bmatrix} 10 \\ 2 \\ \vdots \\ r \end{bmatrix}$$

R β

ملاحظه می شود که $2 \times k \rightarrow R \rightarrow k \times 1$ و $\beta \rightarrow k \times 1$ و $2 \times 1 \rightarrow r$ ؛ بنابراین، فرض بر این است که مدل بر حسب مشاهدات اصلی محاسبه شده است.

برای آزمون فرضیه $R\beta = r$ ابتدا بردار $\hat{\beta}$ را به جای β قرار می دهیم. هر قدر اختلاف $R\hat{\beta}$ با r بیشتر باشد، احتمال اینکه فرضیه H_1 رد شود، بیشتر خواهد بود. بنابراین، مسأله به بررسی توزیع $R\hat{\beta}$ و یافتن روشی که بتوان آماره آزمون را ساخت، تبدیل می شود. اگر آماره آزمون معنی دار باشد، فرضیه H_1 رد می شود و برعکس. ابتدا میانگین $R\hat{\beta}$ را به دست می آوریم:

$$E(R\hat{\beta}) = RE(\hat{\beta}) = R\beta. \quad (6.57)$$

واریانس $R\hat{\beta}$ برابر است با

$$\begin{aligned} \text{Var}(R\hat{\beta}) &= E[R(\hat{\beta} - \beta)(\hat{\beta} - \beta)'R'] \\ &= RE[(\hat{\beta} - \beta)(\hat{\beta} - r)']R'. \end{aligned}$$

با مراجعه به معادله ۶.۷ خواهیم داشت

$$\text{Var}(\mathbf{R}\hat{\beta}) = \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}' \quad (6.58)$$

با توجه به معادله ۶.۱۲ و با توجه به اینکه $\hat{\beta}$ دارای توزیع نرمال است، با استفاده از معادله‌های ۶.۵۷ و ۶.۵۸ داریم

$$\mathbf{R}\hat{\beta} \sim N[\mathbf{R}\beta, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'] \quad (6.59)$$

حال به بررسی توزیع $\mathbf{R}(\hat{\beta} - \beta)$ توجه می‌کنیم،

$$\begin{aligned} E[\mathbf{R}(\hat{\beta} - \beta)] &= \mathbf{R}E(\hat{\beta}) - \mathbf{R}E(\beta), \\ &= \mathbf{R}\beta - \mathbf{R}\beta = \mathbf{0}, \end{aligned}$$

همچنین

$$\text{Var}[\mathbf{R}(\hat{\beta} - \beta)] = \text{Var}(\mathbf{R}\hat{\beta}),$$

زیرا β و \mathbf{R} مقادیر ثابتی هستند. با جایگزینی معادله ۶.۵۸ در معادله فوق داریم

$$\text{Var}[\mathbf{R}(\hat{\beta} - \beta)] = \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'.$$

بدین ترتیب، با توجه به معادله ۶.۵۹ خواهیم داشت

$$\mathbf{R}(\hat{\beta} - \beta) \sim N[\mathbf{0}, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'] \quad (6.60)$$

اگر فرضیه $\mathbf{R}\beta = \mathbf{r}$: H_0 صحیح باشد، با جایگزینی \mathbf{r} به جای $\mathbf{R}\beta$ در رابطه ۶.۶۰ داریم

$$(\mathbf{R}\hat{\beta} - \mathbf{r}) \sim N[\mathbf{0}, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'] \quad (6.61)$$

از جبر ماتریسی می‌دانیم که اگر $\mathbf{x} \sim N(\mathbf{0}, \Sigma)$ آنگاه

$$\mathbf{x}'\Sigma^{-1}\mathbf{x} \sim \chi^2(n), \quad (6.62)$$

که در آن $(n \times 1) \rightarrow \chi$ رابطه ۶-۶۱ را با استفاده از معادله ۶-۶۲ به صورت زیر می نویسیم،

$$(R\hat{\beta} - r)' [\sigma^2 R(X'X)^{-1} R']^{-1} (R\hat{\beta} - r) - \chi^2(r), \quad (6-63)$$

که در آن، r درجات آزادی توزیع χ^2 است و برابر با مقدار ردیفهای بردار $R\hat{\beta}$ ، یعنی مقدار فرضیه‌هایی است که می‌خواهیم آزمون کنیم. همچنین می‌توان ثابت کرد که معکوس ماتریس $[R(X'X)^{-1} R']$ وجود دارد.

تمام اجزای رابطه ۶-۶۳ بغیر از σ^2 ، مقادیر عددی معلوم و مشخصی است. از معادله ۶-۳۸ داریم

$$\frac{e'e}{\sigma^2} \sim \chi^2(n-k).$$

می‌دانیم، اگر دو توزیع مستقل χ^2 را بر درجات آزادی خود تقسیم کرده، سپس نسبت آنها را تشکیل دهیم، یک توزیع F خواهیم داشت. با تقسیم رابطه ۶-۶۳ بر r و رابطه ۶-۳۸ بر $(n-k)$ و تشکیل نسبت آنها اولاً σ^2 حذف می‌شود ثانیاً یک توزیع F با درجات آزادی $[r, (n-k)]$ خواهیم داشت. نتیجه کلی اینکه اگر $H_0: R\beta = r$ صحیح باشد آنگاه

$$\frac{\{(R\hat{\beta} - r)' [R(X'X)^{-1} R']^{-1} (R\hat{\beta} - r)\} / r}{e'e / (n-k)} \sim F(r, n-k). \quad (6-64)$$

رابطه فوق عمومی‌ترین آماره آزمون در مدل‌های رگرسیون چندمتغیره است. اگر مقدار عددی آماره آزمون از مقدار جدول F بیشتر شود، یعنی آماره آزمون در ناحیه بحرانی قرار بگیرد، فرضیه H_0 رد می‌شود.

در اینجا باید نشان داد که چگونه در محاسبات عددی می‌توان از آماره عمومی آزمون استفاده کرد. در اینجا به چند مورد اشاره می‌کنیم تا قدرت و ظرفیت بسیار فرمول فوق در اجرای آزمون فرضیه‌های مختلف روشن شود.

۱. آزمون هر یک از پارامترها

فرض کنید می‌خواهیم در مدل رگرسیون

$$Y_i = \beta_1 + \beta_1 X_{it} + \dots + \beta_k X_{kt} + U_i$$

فرضیه $H_0: \beta_i = a$ را در مقابل $H_1: \beta_i \neq a$ آزمون کنیم. برای استفاده از آماره آزمون ۶-۶۴ باید R و r را تعریف کنیم:

$$R = \begin{bmatrix} \cdot & \cdot & \dots & \cdot & 1 & \dots & \cdot \end{bmatrix}, \quad r = a.$$

↓
i امین عنصر

بنابراین

$$(R\hat{\beta} - r) = \beta_i - a.$$

باید $R(X'X)^{-1}R'$ را برای حالت $\beta_i = a$ به دست آوریم. با توجه به ساختار بردار R ملاحظه می شود که $R(X'X)^{-1}R'$ یک اسکالر بوده، مقدار آن برابر با i امین عنصر قطری ماتریس $(X'X)^{-1}$ است، یعنی

$$R(X'X)^{-1}R' = a_{ii},$$

که a_{ii} در واقع i امین عنصر قطری ماتریس $(X'X)^{-1}$ است. بدین ترتیب آماره عمومی آزمون، یعنی رابطه ۶-۶۴ و با توجه به $r = a$ به قرار زیر خواهد بود،

$$F = \frac{[(\beta_i - a)' (a_{ii})^{-1} (\beta_i - a)] / 1}{\sum e_i^2 / (n - k)} \sim F(1, n - k).$$

معادله فوق را می توان به صورت زیر نوشت:

$$F = \frac{(\beta_i - a)^2}{a_{ii} \hat{\sigma}_u^2} \sim F(1, n - k),$$

که با توجه به معادله ۶-۲۴ داریم

$$F = \frac{(\beta_i - a)^2}{\text{Var}(\hat{\beta}_i)} \sim F(1, n - k). \quad (6.65)$$

ملاحظه می شود که آماره آزمون در رابطه ۶-۶۵ دقیقاً مجذور آماره t است که در رابطه

۶۲۶ به دست آوردیم. انتظار چنین نتیجه‌ای را نیز داشتیم زیرا طبق معادله ۲-۵۱ می‌دانیم $F = t^2$.

آزمون معنی‌دار بودن β_i

اگر بخواهیم فرضیه $H_0: \beta_i = 0$ را در مقابل $H_1: \beta_i \neq 0$ برای یک مدل رگرسیون چندمتغیره آزمون کنیم، بردارهای R و r به صورت زیر تعریف می‌شود،

$$[0 \quad 0 \quad \dots \quad 1 \quad 0 \quad \dots \quad 0], \quad r = 0.$$

↓
i امین عنصر

بنابراین

$$R\hat{\beta} - r = \beta_i.$$

اسکالر $R'(X'X)^{-1}R$ در این حالت برابر است با i امین عنصر قطری ماتریس $(X'X)$ ، یعنی

$$R'(X'X)^{-1}R = a_{ii}.$$

آماره عمومی آزمون برای این حالت و از رابطه ۶۶۴ و با توجه به $r = 0$ عبارت است از

$$F = \frac{[(\beta_i)' (a_{ii})^{-1} (\beta_i)] / 1}{\sum e_i^2 / (n - k)} \sim F(1, n - k).$$

در نتیجه خواهیم داشت

$$F = \frac{\beta_i^2}{\text{Var}(\hat{\beta}_i)} \sim F(1, n - k), \quad (6.66)$$

که دقیقاً مجذور آماره t است که در رابطه ۶۲۸ به دست آوردیم.

۲. آزمون معنی‌دار بودن زیر مجموعه‌ای از پارامترها

مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + U_i.$$

می‌خواهیم، فرضیه زیر را آزمون کنیم،

$$H_1: \beta_{k-r+1} = \beta_{k-r+2} = \dots = \beta_k = 0.$$

فرضیه فوق دلالت بر این می‌کند که باید r پارامتر آخر مدل رگرسیون آزمون شود. ممکن است در بعضی موارد تمام پارامترهایی که می‌خواهیم تساوی آنها را با سفر آزمون کنیم، به طور منظم پارامترهای مدل را تشکیل ندهند. در این صورت می‌توان با یک جابجایی ساده در ترتیب نوشتن متغیرهای توضیحی، پارامترهای مورد نظر را در انتهای مدل قرار داد.

برای آزمون تساوی r پارامتر با صفر، مطابق معمول باید R و r را برای این حالت تعریف کنیم،

$$R = [0 \quad I_r], \quad r = 0.$$

توجه داریم که در اینجا R را به دو ماتریس 0 و I_r افراز کرده‌ایم. همچنین می‌دانیم $r \rightarrow r \times 1$ و $I \rightarrow k \times r$ ، $0 \rightarrow k \times (k-r)$ ، $R \rightarrow r \times k$ رگرسیون $y = X\hat{\beta} + e$ را نیز به ترتیبی مشابه افراز می‌کنیم،

$$y = [X_0 \quad X_r] \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_r \end{bmatrix} + e = X_0 \hat{\beta}_0 + X_r \hat{\beta}_r + e, \quad (6.67)$$

که در آن $X_0 \rightarrow n \times (k-r)$ و $X_r \rightarrow n \times r$. همچنین X_0 و X_r به ترتیب $(k-r)$ ستون اول و r ستون آخر ماتریس X است.

برای محاسبه آماره آزمون، یعنی رابطه ۶-۶۴، مطابق معمول باید $(R\hat{\beta} - r)$ و $R^{-1}(X'X)^{-1}R$ را به دست آوریم. با توجه به تعاریف r و R داریم

$$R\hat{\beta} - r = \hat{\beta}_r, \quad (6.68)$$

و نیز با دقت در ماتریس X و بردار R ، به سهولت می‌توان گفت که $R^{-1}(X'X)^{-1}R$ برابر است با ماتریس مربع مرتبه r که در سمت راست پایین ماتریس $(X'X)^{-1}$ وجود دارد.

اگر این ماتریس را با C_{rr} نشان دهیم، آنگاه با توجه به ساختار ماتریس $(X'X)$ ، یعنی

$$X'X = \begin{bmatrix} X_1'X_1 & X_1'X_r \\ X_r'X_1 & X_r'X_r \end{bmatrix},$$

می توان نشان داد که^۱

$$R(X'X)^{-1}R' = C_{rr} = (X_r'M_1X_r)^{-1}, \quad (6.69)$$

که در آن:

$$M_1 = I - X_1(X_1'X_1)^{-1}X_1'$$

با مراجعه به معادله‌های ۶.۶۸ و ۶.۶۹ می توان آماره عمومی آزمون، یعنی معادله ۶.۶۴ را برای این حالت به قرار زیر نوشت،

$$F = \frac{[\tilde{\beta}'_r(X_r'M_1X_r\tilde{\beta}_r)]/r}{e'e/(n-k)} \sim F(r, n-k). \quad (6.70)$$

با مقایسه آماره آزمون با مقدار جدول F ، در صورتی فرضیه H_0 رد می شود که آماره آزمون در ناحیه بحرانی واقع شود. میتوان نشان داد که معادله ۶.۳۳، یعنی $F = \frac{(RSS_1 - RSS)/r}{RSS/(n-k)}$ ، حاصل معادله ۶.۷۰ است.^۲

۳. آزمون معنی دار بودن مدل رگرسیون

مدل رگرسیون زیر را ملاحظه کنید،

$$Y_i = \beta_1 + \beta_r X_{ri} + \dots + \beta_k X_{ki} + U_i.$$

می خواهیم فرضیه

$$H_0: \beta_r = \beta_r = \dots = \beta_k = 0$$

۱. به قاعده معکوس کردن ماتریسهای افراز شده در پیوست «۵ - الف» مراجعه شود.

۲. برای اثبات به مسأله ۶.۹ مراجعه کنید.

را آزمون کنیم. ابتدا R و r را در رابطه ۶-۶۴ برای این حالت تعریف می‌کنیم،

$$R = \begin{bmatrix} 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{bmatrix}, \quad r = 0,$$

که در آن $R \rightarrow (k-1) \times k$ و $r \rightarrow (k-1) \times 1$. برای ساختن آماره عمومی آزمون، یعنی معادله ۶-۶۴، باید مقادیر $(R\hat{\beta} - r)$ و $R(X'X)^{-1}R'$ را برای این حالت به دست آورد. خواهیم داشت

$$(R\hat{\beta}_r - r) = \begin{bmatrix} \hat{\beta}_r \\ \hat{\beta}_r \\ \vdots \\ \hat{\beta}_k \end{bmatrix} = \hat{\beta}_r. \quad (6.71)$$

همچنین با توجه به ماتریس R ، می‌توان بسادگی نشان داد که $R(X'X)^{-1}R'$ (که یک ماتریس $(k-1) \times (k-1)$ است) با آخرین $(k-1)$ سطر و ستون ماتریس $(X'X)^{-1}$ برابر است. یادآوری می‌کنیم که بردار $\hat{\beta}_r$ در معادله ۶-۷۱ دارای $(k-1)$ عنصر بوده، فاقد $\hat{\beta}_1$ است. بدین ترتیب آماره عمومی آزمون ۶-۶۴ برای فرضیه H_0 عبارت است از

$$F = \frac{\{\hat{\beta}_r' [R(X'X)^{-1}R']^{-1} \hat{\beta}_r\} / (k-1)}{e'e / (n-k)} \sim F(k-1, n-k). \quad (6.72)$$

اگر آماره آزمون فوق معنی‌دار باشد، فرضیه H_0 رد می‌شود.^۱

۴. نواحی اطمینان در رگرسیون چند متغیره

فاصله اطمینان برای یک پارامتر، به راحتی از طریق آماره F به دست می‌آید. در قسمت

۱. می‌توان نشان داد که معادله ۶-۷۲ دقیقاً همان معادله ۶-۳۶ یا ۶-۴۱ است. برای اثبات به مسأله ۶-۱۰ مراجعه شود.

۶۶ به بررسی ناحیه اطمینان برای دو پارامتر پرداختیم. در این قسمت این مسأله را در حالت کلی بررسی خواهیم کرد.
رابطه ۶۶۰ را یک بار دیگر می‌نویسیم،

$$\mathbf{R}(\hat{\beta} - \beta) \sim N[0, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'] .$$

یا توجه به معادله ۶۶۲ خواهیم داشت:

$$[\mathbf{R}(\hat{\beta} - \beta)]' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} [\mathbf{R}(\hat{\beta} - \beta)] \sim \chi^2(r) ,$$

که در آن r درجات آزادی توزیع χ^2 بوده، با تعداد ردیفهای ماتریس \mathbf{R} برابر است. از معادله ۶۳۸ نیز می‌دانیم که

$$\frac{\mathbf{e}'\mathbf{e}}{\sigma_u^2} \sim \chi^2(n-k) .$$

اگر دو توزیع مستقل χ^2 را بر درجات آزادی خود تقسیم نموده، سپس نسبت آنها را تشکیل دهیم نتیجه، یک توزیع F خواهد شد. اگر دقیقاً مشابه روش به کار رفته در معادله ۶۶۴ عمل کنیم، خواهیم داشت

$$F = \frac{[\mathbf{R}(\hat{\beta} - \beta)]' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} [\mathbf{R}(\hat{\beta} - \beta)] / r}{\mathbf{e}'\mathbf{e} / (n-k)} \sim F(r, n-k) . \quad (6.73)$$

به کمک معادله ۶۷۳ می‌توان برای هر گروهی از پارامترهای موجود در یک مدل رگرسیون چندمتغیره، یک ناحیه اطمینان ساخت. بنابراین از نظر محاسباتی، باید اولاً، پارامترهایی را که می‌خواهیم برای آنها ناحیه اطمینان مشترک بسازیم، تعیین کنیم و ثانیاً، ماتریس \mathbf{R} را به تناسب پارامترهایی که مورد نظر است کامل کنیم و مقدار موجود در جدول F را به دست آوریم. ناحیه اطمینان، فضایی است که آماره آزمون، یعنی مقدار محاسبه شده F ، در معادله ۶۷۳ از مقدار F در جدول F کمتر باشد؛ بنابراین اگر نامساوی زیر برقرار باشد، H_0 صحیح است،

$$\frac{[\mathbf{R}(\hat{\beta} - \beta)]' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} [\mathbf{R}(\hat{\beta} - \beta)] / r}{\mathbf{e}'\mathbf{e} / (n-k)} < F(r, n-k) . \quad (6.74)$$

مثال ۶-۱۶ مدلی که در مثال ۶-۲ مطرح کردیم را یک بار دیگر در نظر می‌گیریم،

$$Y_i = \alpha + \beta X_i + \gamma Z_i + U_i.$$

با استفاده از محاسبات مثال ۵-۲، این مدل را به صورت زیر تخمین زده‌ایم،

$$\hat{Y}_i = 4 + 2/5 X_i - 1/5 Z_i.$$

اولاً، برای β و γ ناحیه اطمینان مشترک در سطح احتمال ۹۵ درصد بسازید.

ثانیاً، با استفاده از آماره عمومی آزمون، یعنی معادله ۶-۶۴، فرضیه

$$H_1: \beta + \gamma = 0$$

را در مقابل فرضیه $H_0: \beta + \gamma \neq 0$ در سطح معنی‌دار ۵ درصد آزمون کنید.

۱. برای ساختن نواحی اطمینان باید از معادله ۶-۷۳ استفاده کنیم. ابتدا، ماتریس R

را تشکیل می‌دهیم. با توجه به اینکه می‌خواهیم برای دو پارامتر β و γ ناحیه اطمینان بسازیم، خواهیم داشت

$$R = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

باید $R(\hat{\beta} - \beta)$ و $R(X'X)^{-1}R'$ را محاسبه کنیم. داریم

$$R(\hat{\beta} - \beta) = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} (\hat{\alpha} - \alpha) \\ (\hat{\beta} - \beta) \\ (\hat{\gamma} - \gamma) \end{bmatrix} = \begin{bmatrix} (\hat{\beta} - \beta) \\ (\hat{\gamma} - \gamma) \end{bmatrix},$$

در نتیجه با استفاده از $\hat{\beta} = 2/5$ و $\hat{\gamma} = 1/5$ نتیجه می‌شود که

$$R(\hat{\beta} - \beta) = \begin{bmatrix} (2/5 - \beta) \\ (-1/5 - \gamma) \end{bmatrix}.$$

مفید است به این نکته توجه کنیم که چون نمی‌خواهیم برای α ناحیه اطمینان بسازیم؛

بنابراین می‌توان ستون اول ماتریس R و سطر اول بردار $(\hat{\beta} - \beta)$ را حذف کرد.

برای محاسبه $R(X'X)R'$ ، با توجه به اینکه ماتریس $(X'X)^{-1}$ در مثال ۶-۱ برحسب مقادیر انحراف از میانگین ارائه شده است، باید ستون اول ماتریس R و سطر اول بردار $(\hat{\beta}-\beta)$ را حذف کنیم. به این ترتیب

$$R(X'X)^{-1}R' = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} \\ = \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix}$$

باید $R(X'X)^{-1}R'$ را معکوس کرد

$$[R(X'X)^{-1}R']^{-1} = \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix}^{-1} = \begin{bmatrix} 10 & 6 \\ 6 & 4 \end{bmatrix}$$

محاسبات فوق را در معادله ۶-۷۳ قرار داده، آماره آزمون به شرح زیر محاسبه می‌شود:

$$F = \left\{ \begin{bmatrix} (2/5-\beta) & (-1/5-\gamma) \end{bmatrix} \begin{bmatrix} 10 & 6 \\ 6 & 4 \end{bmatrix} \begin{bmatrix} (2/5-\beta) \\ (-1/5-\beta) \end{bmatrix} \right\} / 2 + \frac{1/5}{5-3} \\ = \frac{26/5 - 32\beta - 18\gamma + 12\beta\gamma + 10\beta^2 + 4\gamma^2}{1/5}$$

مقدار F در سطح معنی دار ۵ درصد و با درجات آزادی ۲ و ۲ برابر است با

$$F_{\%5}(2, 2) = 19.$$

بنابراین براساس معادله ۶-۷۴ می‌توان چنین نوشت،

$$\frac{10\beta^2 + 12\beta\gamma + 4\gamma^2 - 32\beta - 18\gamma + 26/5}{1/5} - 19 = 0.$$

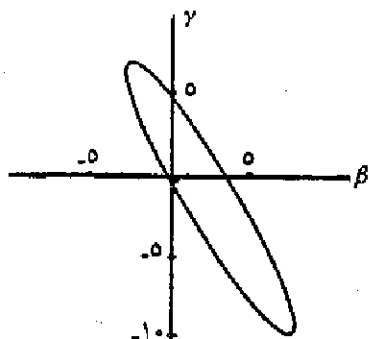
یا بعد از ساده کردن داریم

$$10\beta^2 + 12\beta\gamma + 4\gamma^2 - 32\beta - 18\gamma - 2 = 0.$$

معادله فوق یک بیضی است که در نمودار ۶-۳ ترسیم شده است. مرکز این بیضی در نقاط تخمین $(\hat{\beta} = 2/5, \hat{\gamma} = -1/5)$ قرار دارد. در مثال ۶-۱۳ و در بحث نمودار ۶-۲ دیدیم که اگر کواریانس بین پارامترها مثبت باشد، بیضی به سمت راست متمایل است، در غیر این صورت متمایل به چپ دارد. با توجه به محاسبات موجود در مثال ۶-۱ ملاحظه می‌شود که کواریانس $\hat{\beta}, \hat{\gamma}$ منفی است، زیرا

$$\text{Cov}(\hat{\beta}, \hat{\gamma}) = \hat{\sigma}^2(a_{12}) = 0.75(-1/5) = -1/125,$$

که در آن a_{12} عنصر ردیف اول و ستون دوم ماتریس $(X'X)^{-1}$ است؛ بنابراین بیضی ناحیه اطمینان برای β و γ به سمت چپ متمایل است.



نمودار ۶-۳ ناحیه اطمینان برای β و γ

به وسیله ناحیه اطمینان مشترک می‌توان فرضیه‌های مشترک در مورد پارامترها را آزمون کرد. مثلاً در مثال ۶-۴ فرضیه $H_0: \beta = \gamma = 0$ را آزمون کردیم و ملاحظه شد که H_0 رد نمی‌شود. در نمودار ۶-۳ نیز می‌توان گفت که چون $\beta = \gamma = 0$ درون بیضی قرار دارد؛ بنابراین در ناحیه اطمینان قرار گرفته است و در نتیجه H_0 رد نخواهد شد. بدیهی است در مواردی که می‌خواهیم ناحیه اطمینان را برای بیش از دو پارامتر تعیین کنیم، ترسیم هندسی مشکل یا غیرممکن است. در این موارد کافی است بعد از به دست آوردن رابطه ناحیه اطمینان، یعنی نامساوی ۶-۷۴، فرضیه مشترک پارامترها را در آن، قرار دهیم. اگر H_0 در این نامساوی صدق کند صحیح است و برعکس.

۲. برای آزمون فرضیه $H_0: \beta + \gamma = 0$ به کمک آماره عمومی آزمون، یعنی معادله ۶-۶۴، باید $(R\hat{\beta} - r)$ و $R(X'X)^{-1}R'$ را محاسبه کنیم. با توجه به اینکه $(X'X)^{-1}$ در مثال ۶-۱ برحسب انحراف از میانگین گزارش شده است، به نوشتن عنصر اول در بردار R نیازی نیست. خواهیم داشت

$$R = [1 \quad 1] \quad \text{و} \quad r = 0.$$

بنابراین $(R\hat{\beta} - r)$ برابر است با

$$(R\hat{\beta} - r) = \hat{\beta} + \hat{\gamma}.$$

با توجه به ساختار بردار R ملاحظه می شود که $R(X'X)^{-1}R'$ یک ماتریس 2×2 بوده و برابر است با

$$R(X'X)^{-1}R' = [1 \quad 1] \begin{bmatrix} 1 & -1/5 \\ -1/5 & 2/5 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0/5.$$

بدین ترتیب آماره آزمون عبارت است از

$$F = \frac{\{(R\hat{\beta} - r)' [R(X'X)^{-1}R']^{-1} (R\hat{\beta} - r)\} / r}{e'e / (n - k)},$$

$$= \frac{(\hat{\beta} + \hat{\gamma})' [0/5]^{-1} (\hat{\beta} + \hat{\gamma}) / 1}{1/5 / (5 - 3)} = \frac{(\hat{\beta} + \hat{\gamma})^2}{0/375} = 2/66.$$

مقدار جدول F در سطح معنی دار ۵ درصد و با درجات آزادی $(1, 2)$ برابر است با ۱۸/۵۱. نتیجه می گیریم که H_0 رد نمی شود.

مسائل فصل ششم

۶-۱ با روش حداقل مربعات معمولی و به کمک ۸۰ مشاهده فصلی ($n=80$)، یک مدل رگرسیون را به صورت زیر تخمین زده‌ایم،

$$\hat{Y}_t = 2/20 + 0/104 X_{1t} + 3/48 X_{2t} + 0/34 X_{3t}$$

(3/4) (0/005) (2/2) (0/15)

اعداد داخل پرانتز، مقادیر انحراف از معیار تخمین پارامترهاست. مجموع مربع مقادیر توضیح داده شده و مجموع مربعات پسماند به ترتیب برابر است با تغییرات توضیح داده شده $= 112/5$ و تغییرات توضیح داده نشده $= 19/5$.

الف) در سطح معنی‌دار ۵ درصد، کدام یک از ضرایب متغیرهای توضیحی به طور معنی‌دار با صفر متفاوتند؟

ب) مقدار R^2 را برای این مدل حساب کنید.

ج) مقدار \bar{R}^2 را نیز به دست آورید.

۶-۲ مدل رگرسیون زیر مفروض است،

$$Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + \beta_3 X_{3t} + U_t$$

با استفاده از ۲۶ مشاهده، این مدل را تخمین زده و نتایج زیر را به دست آورده‌ایم،

$$\hat{Y}_t = 2 + 3/5 X_{1t} - 0/7 X_{2t} + 2 X_{3t}$$

(1/9) (2/2) (1/5)

اعداد داخل پرانتز، مقادیر آماره آزمون t برای هر یک از پارامترهاست. همچنین مقدار R^2 برابر است با $0/982$. یک بار دیگر همین مدل را با فرض $\beta_1 = \beta_2$ تخمین زده و نتیجه زیر را ملاحظه کرده‌ایم،

$$\hat{Y}_t = 1/5 + 3(X_{1t} + X_{2t}) - 0/6 X_{3t}$$

(2/7) (2/4)

همچنین مقدار R^2 برای این مدل برابر $۰/۸۷۶$ است.

الف) فرضیه $H_0: \beta_1 = \beta_2$ را آزمون کنید.

ب) اگر متغیر توضیحی X_{2i} را از این معادله حذف کنیم، آیا به نظر شما \bar{R}^2 زیاد می شود یا کم؟

ج) آیا با حذف X_{2i} مقدار R^2 کم می شود؟

۶-۳ مدل تغییرات نرخ دستمزد به صورت زیر مفروض است،

$$Y_i = \alpha + \beta + X_i + \gamma X_i^2 + U_i,$$

که در آن Y_i نرخ تغییرات دستمزد و X_i نرخ بیکاری است. با استفاده از ۱۷ مشاهده، تخمین زیر را به دست آورده ایم،

$$\hat{Y}_i = \frac{23}{5112} - \frac{21}{6710} X_i + \frac{5}{8207} X_i^2.$$

(۹/۵۲) (۱۲/۶۸) (۴/۱۱)

همچنین می دانیم $R^2 = ۰/۴۰۶۲$.

الف) دو فرضیه $H_0: \beta = 0$ و $H_0: \gamma = 0$ را در مقابل فرضیه مخالف صفر در سطح معنی دار ۵ درصد به طور جداگانه آزمون کنید.

ب) فرضیه $H_0: \gamma = \beta = 0$ ، یعنی معنی دار بودن کل مدل را در سطح ۵ درصد آزمون کنید.

ج) به نظر شما تفاوت بین دو استنتاج آماری در بندهای «الف» و «ب» را چگونه می توان توضیح داد؟

۶-۴ مدل زیر مفروض است،

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + U_i.$$

محاسبات زیر برحسب انحراف از میانگین و با استفاده از ۲۳ مشاهده انجام شده است،

$$\sum x_{1i}^2 = 12, \quad \sum x_{1i} x_{2i} = 8, \quad \sum x_{2i}^2 = 12,$$

$$\sum x_{1i} y_i = 10, \quad \sum x_{2i} y_i = 8, \quad \sum y_i^2 = 10.$$

با استفاده از مشاهدات فوق، مدل را به صورت زیر تخمین زده ایم،

$$\hat{Y}_t = 4 + 0.17X_{1t} + 0.2X_{2t} \quad , \quad R^2 = 0.176.$$

فرضیه زیر را در سطح معنی دار ۵ درصد با آماره F و آماره t آزمون کنید.

$$H_1: \frac{\beta_1}{\beta_2} = \frac{5}{3} \quad , \quad H_0: \frac{\beta_1}{\beta_2} \neq \frac{5}{3}.$$

۶-۵ مدل رگرسیون زیر را که بر حسب انحراف از میانگین است ملاحظه نمایید،

$$y_t = \beta_1 x_{1t} + \beta_2 x_{2t} + U_t.$$

اطلاعات زیر موجود است،

$$n = 100 \quad , \quad \sum y_t^2 = \frac{493}{3} \quad , \quad \sum x_{1t}^2 = 30 \quad , \quad \sum x_{2t}^2 = 3 \quad ,$$

$$\sum x_{1t} y_t = 30 \quad , \quad \sum x_{2t} y_t = 20 \quad , \quad \sum x_{1t} x_{2t} = 0.$$

الف) پارامترهای β_1 و β_2 را تخمین بزنید. R^2 را نیز به دست آورید.

ب) فرضیه $H_1: \beta_2 = 7$ را در مقابل $H_0: \beta_2 \neq 7$ آزمون کنید.

ج) فرضیه $H_1: \beta_1 = \beta_2 = 0$ را در مقابل $H_0: \beta_1 \neq 0$ یا $\beta_2 \neq 0$ آزمون کنید.

د) فرضیه $H_1: \beta_2 = 7\beta_1$ را در مقابل $H_0: \beta_2 \neq 7\beta_1$ آزمون کنید.

می توان سطح معنی دار بودن هریک از آزمونها را ۵ درصد گرفت. تمام محاسبات را به زبان ماتریسی انجام دهید.

۶-۶ مدل رگرسیون زیر به عنوان یک تابع تولید تخمین زده شده است،

$$\ln Q_t = 1/37 + 0.132 \ln K_t + 0.452 \ln L_t \quad ,$$

$$(0.207) \quad (0.219)$$

$$R^2 = 0.98 \quad , \quad \text{Cov}(\beta_K, \beta_L) = 0.055.$$

اعداد داخل پرانتز، مقادیر انحراف از معیار تخمین پارامترها هستند.

فرضیه های زیر را در سطح معنی دار ۵ درصد آزمون کنید.

الف) کشش تولید نسبت به سرمایه K و نسبت به نیروی کار L با یکدیگر برابر است.

ب) بازده ثابت نسبت به مقیاس وجود دارد. همان گونه که ملاحظه می شود، در این مسأله تعداد مشاهدات داده نشده است. آیا این امر در استنتاجات شما تأثیری دارد؟
۶-۷ مدل رگرسیون زیر مفروض است،

$$Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + \beta_3 X_{3t} + U_t.$$

محاسبات زیر را براساس ۲۴ مشاهده و برحسب انحراف از میانگین داریم،

$$\begin{aligned} \sum y_t^2 &= 60, & \sum x_{1t} y_t &= 7, & \sum x_{2t} x_{3t} &= 10, \\ \sum x_{1t}^2 &= 10, & \sum x_{2t} y_t &= -7, & \sum x_{1t} x_{3t} &= 0, \\ \sum x_{2t}^2 &= 30, & \sum x_{3t} y_t &= -26, & \sum x_{2t} x_{3t} &= 15, \\ \sum x_{3t}^2 &= 20, & & & & \end{aligned}$$

الف) فرضیه های $\beta_1 = 1$ ، $\beta_2 = 1$ و $\beta_3 = -2$ را آزمون کنید.

ب) فرضیه $\beta_1 + \beta_2 + \beta_3 = 0$ را آزمون کنید.

ج) سه فرضیه بند «الف» را به طور همزمان آزمون کنید؛ یعنی

$$H_0: [\beta_1 \quad \beta_2 \quad \beta_3] = [1 \quad 1 \quad -2].$$

فرضیه ها را در سطح معنی دار ۵ درصد آزمون نمایید.

۶-۸ از یک مجموعه اعداد، نتایج تخمینهای زیر برای دو مدل رگرسیون به دست آمده است،

$$Y_t = 1 + 2 X_t, \quad R^2 = 0.7, \\ (1) \quad (2)$$

$$Y_t = 1/5 + 3 X_t + 4 W_t, \quad R^2 = 0.6. \\ (1/A) \quad (1) \quad (2)$$

اعداد داخل پرانتز مقادیر آماره t مربوط به پارامترهاست. تعداد مشاهدات را به دست آورید.

۶-۹ در مدل رگرسیون $y = X\beta + u$ شامل k پارامتر، می‌خواهیم فرضیه صفر بودن r پارامتر را آزمون کنیم، یعنی

$$H_0: \beta_{k-r+1} = \beta_{k-r+2} = \dots = \beta_k = 0.$$

در معادله ۶-۷۰ نشان دادیم که برای این آزمون می‌توان از آماره زیر استفاده کرد،

$$F = \frac{[\hat{\beta}_r' (X_r' M_r X_r) \hat{\beta}_r] / r}{e'e / (n-k)} \sim F(r, n-k).$$

نشان دهید که از رابطه فوق، می‌توان به آماره زیر رسید،

$$F = \frac{(RSS_r - RSS) / r}{RSS / (n-k)} \sim F(r, n-k).$$

۶-۱۰ در مدل رگرسیون $y = X\beta + u$ شامل k پارامتر، می‌خواهیم فرضیه معنی دار بودن کل مدل را آزمون کنیم، یعنی

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0.$$

در معادله ۶-۷۲ نشان دادیم که برای این آزمون می‌توان از آماره زیر استفاده کرد،

$$F = \frac{\{\hat{\beta}_r' [R(X'X)^{-1}R']^{-1} \hat{\beta}_r\} / (k-1)}{e'e / (n-k)} \sim F(k-1, n-k).$$

نشان دهید که از رابطه فوق می‌توان به آماره زیر رسید،

$$F = \frac{\hat{y}' \hat{y} / (k-1)}{e'e / (n-k)} = \frac{ESS / (k-1)}{RSS / (n-k)} \sim F(k-1, n-k).$$

حل مسائل فصل ششم

۶-۱ الف) درجات آزادی برابر است با $df = 80 - 4 = 76$. مقدار t از جدولی در سطح معنی دار ۲۵ درصد برابر است با $1/98$. هرگاه آماره آزمون از $1/98$ بیشتر باشد، H_0 رد می شود؛ یعنی ضریب مورد نظر به طور معنی داری متفاوت با صفر است. اگر مدل مفروض را به صورت زیر بنویسیم،

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + U_i$$

با فرض $H_0: \beta_i = 0$ و برای مقادیر $i = 2, 3, 4$ داریم

$$t_2 = \frac{\hat{\beta}_2}{SE(\hat{\beta}_2)} = \frac{0/104}{0/005} = 20/8,$$

بنابراین t_2 معنی دار بوده، H_0 رد می شود، یعنی β_2 معنی دار است.

$$t_3 = \frac{\hat{\beta}_3}{SE(\hat{\beta}_3)} = \frac{3/48}{2/2} = 1/581,$$

پس آماره آزمون معنی دار نیست و H_0 رد نمی شود.

$$t_4 = \frac{\hat{\beta}_4}{SE(\hat{\beta}_4)} = \frac{0/34}{0/15} = 2/26,$$

در نتیجه H_0 خواهد شد.

ب) می دانیم تغییرات توضیح داده نشده + تغییرات توضیح داده شده = کل تغییرات
بنابراین $132 = 112/5 + 19/5 =$ کل تغییرات. در نتیجه

$$R^2 = \frac{ESS}{TSS} = \frac{112/5}{132} = 0/85.$$

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k}\right) (1 - R^2) = 1 - \frac{80-1}{80-4} (1 - 0/85) = 0/84. \quad \text{ج}$$

۶-۲ الف) مدل دوم را - که براساس درستی H_0 به دست آمده است - مدل مقید می‌نامیم. آمارهٔ آزمون F را با توجه به مبحث ۶-۷ و نیز معادله‌های ۶-۳۳ یا ۶-۵۱ محاسبه می‌کنیم. می‌دانیم

$$F = \frac{(RSS_r - RSS) / r}{RSS / (n - k)} \sim F(r, n - k),$$

که در آن RSS_r و RSS به ترتیب مجموع مربعات پسماند در مدل مقید و مدل اولیه است. کل تغییرات $\sum y_i^2 = \sum y_i'^2$ برای هر دو مدل یکی است. برای مدل اول و مدل مقید داریم

$$RSS_r = (1 - R^2) \sum y_i^2 = (1 - 0/876) \sum y_i^2 = 0/124 \sum y_i^2,$$

$$RSS = (1 - R^2) \sum y_i^2 = (1 - 0/982) \sum y_i^2 = 0/018 \sum y_i^2.$$

بنابراین با جایگزینی معادله‌های فوق در آمارهٔ آزمون F خواهیم داشت

$$F = \frac{(0/124 \sum y_i^2 - 0/018 \sum y_i^2) / 1}{(0/018 \sum y_i^2) / (26 - 4)} = 129/6.$$

با درجات آزادی ۱، ۲۲ و در سطح معنی‌دار ۵ درصد، از جدول F داریم

$$F_{\%5}(1, 22) = 4/30.$$

نتیجه می‌گیریم که آمارهٔ آزمون معنی‌دار است؛ بنابراین فرضیهٔ $H_0: \beta_1 = \beta_0$ رد می‌شود. (ب) در بند ۶ قسمت ۵-۶ به طور خلاصه به این نکته اشاره کردیم که هنگامی یک متغیر توضیحی می‌تواند \bar{R}^2 را زیاد کند که مقدار عددی آمارهٔ آزمون t پارامتر آن متغیر از یک بیشتر باشد و برعکس. در این مسأله مقدار t مربوط به X_{11} برابر ۰/۷- است. بنابراین با حذف X_{11} از مدل مقدار \bar{R}^2 افزایش خواهد یافت. (ج) با حذف یک متغیر توضیحی مقدار R^2 معمولاً کاهش می‌یابد؛ بنابراین پاسخ این قسمت می‌تواند مثبت باشد.

۶-۳ الف) مقدار جدول t با $14 = 17 - 3 = n - k$ درجهٔ آزادی و در سطح معنی‌دار ۲۵ درصد برابر است با $\pm 2/145$. مقادیر آمارهٔ t برای β و γ به ترتیب عبارت است از

$(1/171)$ و $(1/42)$ ؛ بنابراین از مقدار جدول t کمتر است و معنی دار نیست. پس فرضیه‌های $\beta = 0$ و $\gamma = 0$ رد نخواهد شد.

ب) وج) آماره F را به کمک R^2 تشکیل می‌دهیم. می‌دانیم اگر H_0 صحیح باشد، آماره زیر دارای توزیع F با درجات آزادی $(k-1)$ و $(n-k)$ خواهد بود،

$$F = \frac{R^2 / (k-1)}{(1-R^2) / (n-k)} \sim F(k-1, n-k),$$

$$F = \frac{0/4062}{(1-0/4062)} \cdot \frac{(17-3)}{(3-1)} = 4/7879.$$

اما مقدار جدول F برابر است با

$$F_{\alpha}(2, 14) = 3/74,$$

بنابراین آماره آزمون از این مقدار بیشتر می‌شود و معنی دار است. بنابراین فرضیه $H_0: \beta = \gamma = 0$ رد می‌شود.

ملاحظه می‌شود که هر یک از پارامترها در آزمونهای جداگانه معنی دار نیست، در حالی که در آزمون همزمان معنی دار است. علت این امر می‌تواند همبستگی بین متغیرهای توضیحی باشد.

۶-۴ فرضیه H_0 را می‌توان به صورت زیر نوشت،

$$H_0: \beta_2 = 0/6 \beta_1, \quad H_1: \beta_2 \neq 0/6 \beta_1.$$

فرض می‌کنیم H_0 صحیح است. با جایگزینی H_0 در مدل مفروض، مدل مقید را به دست می‌آوریم،

$$\begin{aligned} Y_i &= \alpha + \beta_1 X_{1i} + 0/6 \beta_1 X_{2i} + U_i, \\ &= \alpha + \beta_1 (X_{1i} + 0/6 X_{2i}) + U_i. \end{aligned}$$

باید مجموع مربعات پسماند برای این مدل مقید را حساب کنیم. ابتدا مدل مقید را تخمین می‌زنیم،

$$\hat{\beta}_1 = \frac{\sum x_i^* y_i}{\sum x_i^{*2}},$$

که در آن علامت * بیان کننده این نکته است که محاسبات برای مدل مقید انجام شده است. باید صورت و مخرج کسر فوق را محاسبه کنیم. می دانیم

$$x_i^* = x_{1i} + 0/6 x_{2i}.$$

دو طرف رابطه فوق را مجذور کرده و برای تمام مشاهدات جمع می کنیم،

$$\begin{aligned} \sum x_i^{*2} &= \sum x_{1i}^2 + (0/6)^2 \sum x_{2i}^2 + 2(0/6) \sum x_{1i} x_{2i}, \\ &= 12 + 0/36 (12) + 1/2 (8) = 25/92. \end{aligned}$$

به همین ترتیب خواهیم داشت

$$\begin{aligned} \sum x_i^* y_i &= \sum x_{1i} y_i + 0/6 \sum x_{2i} y_i, \\ &= 10 + 0/6 (8) = 14/8. \end{aligned}$$

در نتیجه، β_1^* به دست می آید،

$$\beta_1^* = \frac{14/8}{25/92} = 0/57.$$

مجموع مربعات پسماند برای مدل مقید، یعنی RSS_T را حساب می کنیم،

$$\begin{aligned} RSS_T &= \sum y_i^2 - \hat{\beta}_1^* \sum x_i^* y_i, \\ &= 10 - 0/57 (14/8) = 1/564. \end{aligned}$$

حال باید مجموع مربعات پسماند را برای مدل اصلی، یعنی RSS ، نیز به دست آوریم.

می دانیم

$$\begin{aligned} RSS &= \sum y_i^2 - \hat{\beta}_1 \sum x_{1i} y_i - \hat{\beta}_2 \sum x_{2i} y_i, \\ &= 10 - 0/7 (10) - 0/2 (8) = 1/4. \end{aligned}$$

اگر H_0 صحیح باشد، آماره زیر توزیع F خواهد داشت

$$\begin{aligned} F &= \frac{(RSS_T - RSS) / 1}{RSS / (23 - 2)} \sim F(1, 20), \\ &= \frac{(1/564 - 1/4) / 1}{1/4 / 20} = \frac{0/164}{0/07} = 2/342. \end{aligned}$$

مقدار جدول F در سطح معنی دار ۵ درصد برابر است با

$$F_{\%5}(1, 20) = 4/35,$$

در نتیجه آماره آزمون معنی دار نبوده و فرضیه H_1 رد نمی شود. راه حل دیگر، استفاده از آزمون t است. فرضیه H_0 را به صورت زیر می نویسیم،

$$H_0: \beta_2 - 0/6 \beta_1 = 0.$$

باید $\hat{\beta}_2 - 0/6 \hat{\beta}_1$ را استاندارد کنیم. ابتدا محاسبات زیر را انجام می دهیم،

$$E(\hat{\beta}_2 - 0/6 \hat{\beta}_1) = \beta_2 - 0/6 \beta_1,$$

$$\text{Var}(\hat{\beta}_2 - 0/6 \hat{\beta}_1) = \text{Var}(\hat{\beta}_2) + (0/6)^2 \text{Var}(\hat{\beta}_1) - 2(0/6) \text{Cov}(\hat{\beta}_1, \hat{\beta}_2),$$

$$\text{Var}(\hat{\beta}_2) = \frac{\sigma^2}{(1 - r_{12}^2) \sum x_{1i}^2}, \quad r_{12}^2 = \frac{(\sum x_{1i} \sum x_{2i})^2}{\sum x_{1i}^2 \sum x_{2i}^2},$$

$$r_{12}^2 = \frac{(8)^2}{12(12)} = \frac{74}{144} = \frac{4}{9},$$

$$\text{Var}(\hat{\beta}_2) = \frac{\sigma^2}{(1 - \frac{4}{9}) 12} = \frac{3}{20} \sigma^2, \quad \text{Var}(\hat{\beta}_1) = \frac{3}{20} \sigma^2,$$

$$\begin{aligned} \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) &= \frac{-\sigma^2 r_{12}^2}{(1 - r_{12}^2) \sum x_{1i} x_{2i}} \\ &= \frac{-\sigma^2 (4)(6)}{9(0)8} = -\frac{1}{10} \sigma^2. \end{aligned}$$

بدین ترتیب خواهیم داشت

$$\text{Var}(\hat{\beta}_2 - 0/6 \hat{\beta}_1) = \frac{3}{20} \sigma^2 + (0/6)^2 \frac{3}{20} \sigma^2 + 2(0/6) \frac{1}{10} \sigma^2 = 0/324 \sigma^2,$$

$$\hat{\sigma}^2 = \frac{\sum e_i^2}{(n - k)} = \frac{1/4}{20} = 0/07.$$

با اطلاعات فوق می توان آماره t را به شرح زیر محاسبه کرد،

$$t = \frac{(\hat{\beta}_2 - 0/6 \hat{\beta}_1) - E(\hat{\beta}_2 - 0/6 \hat{\beta}_1)}{\sqrt{\text{Var}(\hat{\beta}_2 - 0/6 \hat{\beta}_1)}} = \frac{[0/2 - 0/6(0/7)] - 0}{\sqrt{0/324(0/07)}}$$

$$= \frac{-0/22}{0/1500} = -1/461.$$

مقدار موجود در جدول t با $20 = n - 3$ درجه آزادی و در سطح معنی دار $2/5$ درصد برابر است با $2/086 \pm$ و در نتیجه آماره t معنی دار نبوده و فرضیه H_0 رد نمی شود. انتظار چنین نتیجه ای را نیز داشتیم زیرا

$$t^2 = \frac{(-0/22)^2}{0/324(0/07)} = 2/134,$$

که تقریباً برابر آماره آزمون F است. اختلاف ناشی از تقریب در محاسبات است. ۶.۵ الف) ماتریس $(X'X)$ و بردار $X'y$ را تشکیل می دهیم،

$$(X'X) = \begin{bmatrix} 30 & 0 \\ 0 & 3 \end{bmatrix}, \quad X'y = \begin{bmatrix} 30 \\ 20 \end{bmatrix}.$$

در نتیجه خواهیم داشت

$$\begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = (X'X)^{-1} X'y = \begin{bmatrix} 1/30 & 0 \\ 0 & 1/3 \end{bmatrix} \begin{bmatrix} 30 \\ 20 \end{bmatrix} = \begin{bmatrix} 1 \\ 20/3 \end{bmatrix},$$

بنابراین $\hat{\beta}_1 = 1$ و $\hat{\beta}_2 = 20/3$ همچنین داریم

$$R^2 = \frac{\hat{\beta}' X'y}{y'y} = \frac{\hat{\beta}_1 \sum x_{1t} y_t + \hat{\beta}_2 \sum x_{2t} y_t}{\sum y_t^2}$$

$$= \frac{1(30) + \frac{20}{3}(20)}{\frac{493}{3}} = \frac{490}{493} = 0/99.$$

ب) ابتدا واریانس U_t را به شرح زیر حساب می‌کنیم،

$$\begin{aligned} \text{RSS} &= (1 - R^2) \text{TSS} , \\ &= \left(1 - \frac{490}{493}\right) \left(\frac{493}{3}\right) = 1 , \end{aligned}$$

$$\hat{\sigma}_u^2 = \frac{\text{RSS}}{n - k} = \frac{1}{100 - 3} = \frac{1}{97} .$$

حال آماره t را به دست می‌آوریم،

$$t = \frac{\hat{\beta}_2 - \beta_2}{\hat{\sigma} \sqrt{a_{22}}} = \frac{\frac{20}{3} - 7}{\sqrt{\frac{1}{97}} \sqrt{\frac{1}{3}}} = -0.69 .$$

یادآوری می‌شود که a_{22} به معنی دومین عنصر قطری ماتریس $(X'X)^{-1}$ است. مقدار t از جدول t برابر است با $\pm 1/98$ ؛ بنابراین آماره آزمون در ناحیه بحرانی قرار گرفته و معنی‌دار است. در نتیجه فرضیه H_0 رد می‌شود.

ج) این فرضیه به آزمون معنی‌دار بودن کل رگرسیون مربوط است. از فرمول F برحسب R^2 در معادله ۶-۴۱ استفاده می‌کنیم،

$$F = \left(\frac{n - k}{k - 1}\right) \frac{R^2}{1 - R^2} \sim F(k - 1, n - k) ,$$

$$F = \left(\frac{100 - 3}{2}\right) \frac{490}{493(1 - 490/493)} = \frac{97}{2} \cdot \frac{490(493)}{493(3)} = 7921/6 .$$

ملاحظه می‌شود که آماره آزمون معنی‌دار است و فرضیه H_0 رد می‌شود.

د) آزمون فرضیه $H_0: \beta_1 = 7\beta_2$ از راه ماتریسی می‌تواند به دو صورت انجام شود. استفاده از فرمول ۶-۵۰ با استاندارد کردن β یا نوشتن H_0 به صورت $R\beta = r$ و استفاده از آماره عمومی آزمون، یعنی فرمول ۶-۶۴. از هر دو راه این مسأله را حل می‌کنیم تا رابطه بین این دو فرمول نیز روشن شود.

راه حل اول: از معادله ۶-۵۰ می‌دانیم که

$$t = \frac{\mathbf{C}'\hat{\beta} - r}{\sigma\sqrt{\mathbf{C}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}}} \sim t(n-k).$$

H_0 را به صورت $\beta_1 - \gamma\beta_2 = 0$ نوشته، داریم

$$\begin{bmatrix} -\gamma & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = 0,$$

در نتیجه $\mathbf{C}'\hat{\beta}$ برابر است با

$$\mathbf{C}'\hat{\beta} = \begin{bmatrix} -\gamma & 1 \end{bmatrix} \begin{bmatrix} 1 \\ \frac{2}{3} \end{bmatrix} = -\frac{1}{3}.$$

صورت کسر t برابر خواهد بود با

$$\mathbf{C}'\hat{\beta} - r = -\frac{1}{3} - 0 = -\frac{1}{3}.$$

برای سهولت محاسبات، ابتدا مجذور معخرج کسر t را حساب می‌کنیم،

$$\sigma' \mathbf{C}' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C} = \frac{1}{97} \begin{bmatrix} -\gamma & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{30} & 0 \\ 0 & \frac{1}{3} \end{bmatrix} \begin{bmatrix} -\gamma \\ 1 \end{bmatrix} = \frac{1}{97} \left(\frac{59}{30} \right) = 0.020274.$$

آماره t به شرح زیر محاسبه می‌شود،

$$t = \frac{-1}{3\sqrt{0.020274}} = -2/34.09,$$

که معنی‌دار بوده و فرضیه H_0 رد می‌شود.

راه حل دوم: از معادله ۶-۶۴ می‌دانیم که

$$F = \frac{\{(\mathbf{R}\hat{\beta} - r)' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - r)\} / r}{e'e / (n-k)} \sim F(r, n-k),$$

فرمول فوق می‌تواند به طور همزمان چند ترکیب خطی را با هم آزمون کند. با توجه

به اینکه در این مسأله، فقط یک ترکیب خطی را می‌خواهیم آزمون کنیم، بنابراین ماتریس R به همان بردار C' تبدیل خواهد شد و در نتیجه، دقیقاً همان $(R\hat{\beta} - r)$ می‌شود. یعنی

$$R\hat{\beta} - r = -\frac{1}{3}.$$

به ترتیبی مشابه خواهیم داشت

$$R(X'X)^{-1}R' = C'(X'X)^{-1}C = \frac{0.9}{30}.$$

بنابراین، صورت کسر F برابر است با

$$\left[\left(-\frac{1}{3}\right) \left(\frac{0.9}{30}\right)^{-1} \left(-\frac{1}{3}\right) \right] / 1,$$

و مخرج کسر F نیز عبارت خواهد بود از

$$\frac{e'e}{n-k} = \hat{\sigma}_u^2 = \frac{1}{97}.$$

در نتیجه آماره F برابر خواهد بود با

$$F = \frac{\left(\frac{1}{3}\right)^2 / 1}{\frac{0.9}{30} \left(\frac{1}{97}\right)}.$$

ملاحظه می‌شود $\sqrt{0.9/48} = 2/34.09$ ، یعنی $\sqrt{F} = t$.

۶-۶ این مسأله اهمیت فراوان دارد؛ زیرا تعداد مشاهدات و واریانس جمله اختلال داده نشده است. آماره عمومی آزمون، یعنی معادله ۶-۶، را می‌نویسیم. می‌دانیم مخرج کسر، یعنی $e'e / (n-k)$ ، در واقع همان تخمین واریانس جمله اختلال، یعنی $\hat{\sigma}_u^2$ ، است.

$$F = \frac{\{(R\hat{\beta} - r)' [R(X'X)^{-1}R']^{-1} (R\hat{\beta} - r)\} / r}{\hat{\sigma}_u^2}.$$

استفاده از فرمول فوق مستلزم داشتن مقادیر جداگانه برای $\hat{\sigma}_u^2$ و $[R(X'X)^{-1}R']^{-1}$

است که متأسفانه در این مسأله موجود نیست؛ با وجود این، اگر آماره F را به صورت زیر بنویسیم،

$$F = \{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}\hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})\} / r, \quad (1)$$

می توان، $[\mathbf{R}\hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1}$ را به راحتی محاسبه کرد. برای اثبات، ابتدا براساس معادله ۷۷، معادله واریانس - کوواریانس $\hat{\beta}$ را می نویسیم،

$$\Sigma_{\hat{\beta}} = \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}.$$

ماتریس فوق را از سمت چپ و راست به ترتیب در \mathbf{R} و \mathbf{R}' ضرب می کنیم،

$$\mathbf{R}\hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'. \quad (2)$$

با معکوس کردن ماتریس (۲)، جمله مجهول آماره F در (۱) به دست می آید. بدیهی است مقدار $(\mathbf{R}\hat{\beta} - \mathbf{r})$ کاملاً معلوم و قابل محاسبه است. محاسبه آماره F ، مستلزم داشتن $\hat{\sigma}^2$ و $(\mathbf{X}'\mathbf{X})^{-1}$ به طور جداگانه نیست. نتیجه کلی این است که اگر ماتریس واریانس - کوواریانس $\hat{\beta}$ را داشته باشیم کافی است با استفاده از \mathbf{R} ، ماتریس (۲) را ساخته و آن را معکوس کنیم تا آماره F در رابطه (۱) قابل محاسبه شود. با استفاده از صورت مسأله می توان ماتریس واریانس - کوواریانس $\hat{\beta}$ را به دست آورد.

$$\Sigma_{\hat{\beta}} = \begin{bmatrix} \text{Var}(\hat{\beta}_K) & \text{Cov}(\hat{\beta}_K, \hat{\beta}_L) \\ \text{Cov}(\hat{\beta}_L, \hat{\beta}_K) & \text{Var}(\hat{\beta}_L) \end{bmatrix},$$

$$= \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} (0/207)^2 & 0/000 \\ 0/000 & (0/219)^2 \end{bmatrix} = \begin{bmatrix} 0/066 & 0/000 \\ 0/000 & 0/048 \end{bmatrix},$$

الف) با توجه به اینکه در مدل مفروض، متغیرها برحسب لگاریتم نوشته شده است، کشش تولید نسبت به L و K به ترتیب برابر است با β_L و β_K . بنابراین فرضیه

تساوی دو کَشش عبارت است از فرضیه $H_0: \beta_K = \beta_L$. از آماره عمومی آزمون استفاده کرده، $R\beta = r$ را می‌سازیم،

$$R = \begin{bmatrix} 1 & -1 \end{bmatrix}, \quad r = 0,$$

در نتیجه خواهیم داشت

$$R\hat{\beta} - r = \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 0/632 \\ 0/452 \end{bmatrix} - 0 = 0/18.$$

ماتریس (۲) را تشکیل می‌دهیم،

$$R\hat{\sigma}^2 (X'X)^{-1} R' = \begin{bmatrix} 1 & -1 \end{bmatrix} \begin{bmatrix} 0/066 & 0/000 \\ 0/000 & 0/048 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 0/004.$$

آماره F از رابطه (۱) به سهولت قابل محاسبه است،

$$F = \{ (0/18) [0/004]^{-1} (0/18) \} / 1 = \frac{(0/18)^2}{0/004} = 8/1.$$

در این قسمت باید آماره آزمون را با مقدار F از جدول مقایسه کنیم که خود مستلزم داشتن تعداد مشاهدات است که در این مسأله موجود نیست؛ البته بدون داشتن تعداد مشاهدات نیز می‌توان فرضیه مورد نظر را در این مسأله خاص آزمون کرد. می‌دانیم مقدار $R^2 = 0/98$ و نیز مقادیر نسبتاً کوچک واریانس تخمین پارامترها می‌توانند در واقع دلالت بر وجود حجم قابل قبولی از مشاهدات داشته باشد. با توجه به اینکه در این مدل سه پارامتر تخمین زده می‌شود انتظار این است که حجم مشاهدات از ۱۰ بیشتر باشد. حتی اگر حجم مشاهدات را ۷ بگیریم، درجات آزادی منخرج کسر F در معادله ۶-۶۴ برابر است با $n - k = 7 - 3 = 4$. مقدار F به دست آمده از جدول را برای کمترین و بیشترین مقادیر ممکن درجات آزادی منخرج به دست می‌آوریم،

$$F_{\%0}(1, 4) = 7/71, \quad F_{\%0}(1, \infty) = 3/84.$$

در هر دو حالت مقدار آماره آزمون بیشتر از مقدار جدول F می‌شود و فرضیه H_0 ، یعنی برابری کَششها رد می‌شود.

ب) فرضیهٔ بازده ثابت نسبت به مقیاس در این مسأله را می‌توان به صورت زیر نوشت،

$$\beta_K + \beta_L = 1.$$

دوباره از آمارهٔ عمومی آزمون، یعنی معادلهٔ (۱) استفاده می‌کنیم،

$$R = \begin{bmatrix} 1 & 1 \end{bmatrix}, \quad r = 1,$$

$$(R\hat{\beta} - r) = 0/084,$$

$$R \hat{\sigma}^2 (X'X)^{-1} R' = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0/066 & 0/000 \\ 0/000 & 0/048 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0/224.$$

بدین ترتیب داریم

$$F = \{(0/084) [0/224]^{-1} (0/084)\} / 1 = \frac{(0/084)^2}{0/224} = 0/315.$$

برای هر تعدادی از مشاهدات که بتوان تصور کرد، مقدار آمارهٔ آزمون معنی‌دار نیست؛ بنابراین فرضیهٔ بازده ثابت به مقیاس را نمی‌توان رد کرد.

۶-۷ ابتدا پارامترهای مدل را تخمین می‌زنیم. برای این منظور باید محاسبات زیر را

انجام داد،

$$(X'X) = \begin{bmatrix} 10 & 10 & 0 \\ 10 & 30 & 10 \\ 0 & 10 & 20 \end{bmatrix}, \quad X'y = \begin{bmatrix} 7 \\ -7 \\ -26 \end{bmatrix},$$

$$(X'X)^{-1} = \begin{bmatrix} 0/15 & -0/00 & 0 \\ -0/00 & 0/07 & -0/04 \\ 0 & -0/04 & 0/08 \end{bmatrix}.$$

در نتیجه بردار $\hat{\beta}$ عبارت است از

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix} = \begin{bmatrix} 0/15 & -0/05 & 0 \\ -0/05 & 0/07 & -0/04 \\ 0 & -0/04 & 0/08 \end{bmatrix} \begin{bmatrix} 7 \\ -7 \\ -26 \end{bmatrix} = \begin{bmatrix} 1/4 \\ 0/2 \\ -1/8 \end{bmatrix}$$

برای تخمین واریانس جمله اختلال، ابتدا ESS را حساب می‌کنیم،

$$ESS = \hat{\beta}' X' y = [1/4 \quad 0/2 \quad -1/8] \begin{bmatrix} 7 \\ -7 \\ -26 \end{bmatrix} = 55/2,$$

$$RSS = e'e = y'y - \hat{\beta}' X'y,$$

$$= 70 - 55/2 = 4/8,$$

و بدین ترتیب $\hat{\sigma}^2$ عبارت است از

$$\hat{\sigma}^2 = \frac{e'e}{n-k} = \frac{4/8}{24-4} = 0/24, \quad \hat{\sigma}^2 = 0/490.$$

(الف) از آزمون t استفاده می‌کنیم،

$$t = \frac{\hat{\beta}_i - a}{SE(\hat{\beta}_i)}.$$

برای آزمون $\beta_1 = 1$ داریم

$$t = \frac{1/4 - 1}{0/490 \sqrt{0/15}} = 2/11.$$

با توجه به اینکه مقدار t از جدول با ۲۰ درجه آزادی برابر است با ۲/۰۸۶؛ بنابراین آماره آزمون معنی‌دار است و فرضیه $\beta_1 = 1$ رد می‌شود.

برای آزمون $\beta_2 = 1$ داریم

$$t = \frac{0/2 - 1}{0/490 \sqrt{0/07}} = -6/17,$$

بنابراین آماره آزمون معنی دار است و فرضیه $\beta_2 = 1$ رد می شود.
برای آزمون $\beta_2 = 2$ داریم

$$t = \frac{-1/8 - (2)}{0/490 \sqrt{0/08}} = 1/44,$$

در نتیجه آماره آزمون معنی دار نبوده و فرضیه $\beta_2 = 2$ رد نمی شود.

ب) برای آزمون $H_0: \beta_1 + \beta_2 + \beta_3 = 0$ ، از آماره عمومی آزمون، یعنی معادله ۶-۶۴ استفاده می کنیم،

$$F = \frac{\{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})\} / r}{e'e / (n-k)}$$

می دانیم

$$\mathbf{R}\hat{\beta} - \mathbf{r} = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1/4 \\ 0/2 \\ -1/8 \end{bmatrix} = -0/2,$$

$$\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}' = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0/15 & -0/00 & 0 \\ -0/00 & 0/07 & -0/04 \\ 0 & -0/04 & 0/08 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 0/12,$$

بنابراین آماره آزمون به شرح زیر محاسبه می شود،

$$F = \frac{\{(-0/2) [0/12]^{-1} (-0/2)\} / 1}{4/8 / (24 - 4)} = \frac{0/04}{0/0288} = 1/389.$$

مقدار F از جدول برابر است با

$$F_{\%} (1, 20) = 4/35,$$

بنابراین آماره آزمون معنی دار نبوده و فرضیه H_0 ، یعنی $\beta_1 + \beta_2 + \beta_3 = 0$ رد نمی شود.

ج) برای آزمون همزمان

$$H_0: [\beta_1 \quad \beta_2 \quad \beta_3] [1 \quad 1 \quad -2],$$

از آماره عمومی آزمون استفاده می‌کنیم. داریم

$$R = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad r = \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix}.$$

بنابراین، خواهیم داشت

$$R\hat{\beta} - r = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1/4 \\ 0/2 \\ -1/8 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix} = \begin{bmatrix} 0/4 \\ -0/8 \\ 0/2 \end{bmatrix}.$$

همچنین می‌دانیم که

$$[R(X'X)^{-1}R']^{-1} = X'X,$$

زیرا R یک ماتریس واحد است. بدین ترتیب، آماره عمومی آزمون از معادله ۶-۶۴ به شرح زیر محاسبه می‌شود،

$$F = \frac{\{(R\hat{\beta} - r)' [X'X]^{-1} (R\hat{\beta} - r)\} / (k-1)}{e'e / (n-k)}.$$

ابتدا صورت کسر F را حساب می‌کنیم.

$$\left\{ (0/4 \quad 0/8 \quad 0/2) \begin{bmatrix} 0/16 & -0/00 & 0 \\ -0/00 & 0/07 & -0/04 \\ 0 & -0/04 & 0/08 \end{bmatrix} \begin{bmatrix} 0/4 \\ -0/8 \\ 0/2 \end{bmatrix} \right\} / 3 = \frac{11/2}{3},$$

مقدار آماره عمومی آزمون عبارت خواهد بود از

$$F = \frac{11/2/3}{4/8/20} = 10/06.$$

یا توجه به مقدار $F_{\%}(3, 20) = 3/10$ ، نتیجه می‌گیریم که آماره آزمون معنی‌دار است و فرضیه H_0 ، یعنی $[\beta_1 \ \beta_2 \ \beta_3] = [1 \ 1 \ 2]$ ، رد می‌شود.

۶-۸ مدل اول حالت مقید مدل دوم است؛ به عبارت دیگر اگر در مدل دوم بخواهیم فرضیه تساوی ضریب متغیر W_1 برابر صفر را آزمون کنیم و سپس این قید را در مدل دوم جایگزین کنیم، مدل اول به دست می‌آید. در مدل دوم، مقدار آماره t متعلق به W_1 برابر ۳ محاسبه شده است. با توجه به اینکه $F = t^2$ ، بنابراین مقدار F متعلق به همین آزمون برابر ۹ خواهد بود. اما می‌دانیم آماره F برابر است با

$$F = \frac{(RSS_r - RSS)}{RSS / (n - k)}$$

از طرف دیگر، $r = 1$ و $(n - k) = n - 3$ ؛ بنابراین

$$F = \frac{(RSS_r - RSS)}{RSS / (n - 3)} = (n - 3) \left(\frac{RSS_r}{RSS} - 1 \right).$$

مقادیر کل تغییرات برای هر دو مدل متساوی است. می‌دانیم

$$RSS_r = (1 - R_r^2) TSS = (1 - 0/2) TSS = 0/8 TSS,$$

$$RSS = (1 - R^2) TSS = (1 - 0/6) TSS = 0/4 TSS,$$

بنابراین

$$F = (n - 3) \left(\frac{0/8 TSS}{0/4 TSS} - 1 \right).$$

همچنین با توجه به اینکه $F = 9$ ، پس

$$9 = (n - 3) (2 - 1), \quad n = 12.$$

۶-۹ برای اثبات، باید نشان دهیم که

$$\hat{\beta}'_r (X'_r M_r X_r) \hat{\beta}_r = RSS_r - RSS = e'_r e_r - e' e.$$

معادله ۶-۶۷ را در نظر گرفته و دو طرف آن را از سمت چپ در M_2 ضرب می‌کنیم،

$$M_2 y = M_2 X_2 \hat{\beta}_2 + M_2 X_1 \hat{\beta}_1 + M_2 e. \quad (1)$$

با توجه به تعریف ماتریس M_2 ، یعنی

$$M_2 = I - X_2 (X_2' X_2)^{-1} X_2', \quad (2)$$

می‌توان نشان داد که $M_2 X_2 = 0$ ؛ زیرا

$$M_2 X_2 = I X_2 - X_2 (X_2' X_2)^{-1} (X_2' X_2) = 0 \quad (3)$$

همچنین $M_2 e = e$ ؛ زیرا

$$M_2 e = I e - X_2 (X_2' X_2)^{-1} X_2' e. \quad (4)$$

اما با توجه به تعریف ماتریس X در معادله ۶-۶۷، داریم

$$X' e = [X_2' \quad X_1'] e = [X_2' e \quad X_1' e] = [0 \quad 0],$$

زیرا می‌دانیم جمله‌های پسماند از متغیرهای توضیحی مستقل است. با جایگزینی معادله‌های (۳) و (۴) در معادله (۱)، خواهیم داشت

$$M_2 y = M_2 X_1 \hat{\beta}_1 + e.$$

معادله فوق را ترانهاد کرده، مجذور می‌کنیم، در نتیجه

$$y' M_2 y = \hat{\beta}_1' X_1' M_2 X_1 \hat{\beta}_1 + e' e. \quad (5)$$

نشان می‌دهیم که $y' M_2 y$ برابر با $e' e$ است. برای این منظور دو طرف معادله (۲) را از سمت راست در y و از سمت چپ در y' ضرب می‌کنیم،

$$\begin{aligned} y' M_2 y &= y' y - y' X_2 (X_2' X_2)^{-1} X_2' y \\ &= y' y - y' X_2 \hat{\beta}_2, \\ &= y' y - y' \hat{y}_2. \end{aligned} \quad (6)$$

اما با توجه به معادله ۵.۵۲ می دانیم که

$$= \mathbf{y}' \hat{\mathbf{y}}_e = \hat{\mathbf{y}}_e' \hat{\mathbf{y}}_e,$$

و با جایگزینی در معادله (۶)، داریم

$$\mathbf{y}' \mathbf{M}_z \mathbf{y} = \mathbf{y}' \mathbf{y} - \hat{\mathbf{y}}_e' \hat{\mathbf{y}}_e = \mathbf{e}_e' \mathbf{e}_e.$$

با جایگزینی معادله فوق در معادله (۵)، خواهیم داشت

$$\begin{aligned} \hat{\beta}_r' \mathbf{X}_r' \mathbf{M}_z \mathbf{X}_r \hat{\beta}_r &= \mathbf{e}_e' \mathbf{e}_e - \mathbf{e}' \mathbf{e}, \\ &= \text{RSS}_r - \text{RSS}. \end{aligned}$$

رابطه فوق را در آماره F جایگزین می کنیم،

$$F = \frac{(\text{RSS}_r - \text{RSS}) / r}{\text{RSS} / (n - k)} \sim F(r, n - k).$$

۶-۱۰ می دانیم در معادله ۶.۷۲، ماتریس R به صورت زیر تعریف شده است،

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{bmatrix},$$

که در آن $k \times (k-1) \rightarrow \mathbf{R}$. به کمک این ماتریس باید نشان دهیم که

$$\hat{\beta}_r' [\mathbf{R} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{R}']^{-1} \hat{\beta}_r = \hat{\mathbf{y}}' \hat{\mathbf{y}} = \text{ESS}.$$

ملاحظه می شود که $\mathbf{R}' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{R}$ ، برابر با آخرین $(k-1)$ سطر و ستون ماتریس

$(\mathbf{X}' \mathbf{X})^{-1}$ است. برای اینکه ساختار این ماتریس را بررسی کنیم ابتدا ماتریس X به

صورت زیر افراز می شود،

$$\mathbf{X} = [\mathbf{i} \quad \mathbf{X}_r],$$

که در آن i یک بردار ستونی شامل n عدد یک است و X_r یک ماتریس $(k-1) \times n$ از مشاهدات تمام متغیرهای توضیحی است. ماتریس $(X'X)$ را تشکیل می‌دهیم،

$$(X'X) = \begin{bmatrix} n & i' X_r \\ X_r' i & X_r' X_r \end{bmatrix} .$$

با استفاده از قاعده معکوس کردن ماتریسهای افزاشده^۱، می‌دانیم که اگر $(X'X)$ را معکوس کنیم، زیر ماتریس تشکیل شده از آخرین $(k-1)$ سطر و ستون آن عبارت خواهد بود از

$$(X_r' X_r - X_r' i \frac{1}{n} i' X_r)^{-1} .$$

پیش از این نیز گفتیم که ماتریس $R' (X'X)^{-1} R$ برابر با آخرین $(k-1)$ سطر و ستون ماتریس $(X'X)^{-1}$ است؛ بنابراین

$$R (X'X)^{-1} R' = (X_r' X_r - X_r' i \frac{1}{n} i' X_r)^{-1} .$$

با توجه به ساختار بردار i و با مراجعه به تعریف ماتریس تبدیل Λ که در معادله ۵-۷۱ تعریف شده است، داریم

$$(X_r' X_r - X_r' i \frac{1}{n} i' X_r)^{-1} = (X_r' \Lambda X_r)^{-1} ,$$

در نتیجه

$$R (X'X)^{-1} R' = (X_r' \Lambda X_r)^{-1} , \quad (۱)$$

که در آن عناصر ماتریس X_r بر حسب مقادیر اصلی است. معادله (۱) را در آماره عمومی آزمون، یعنی معادله ۶-۷۲ جایگزین می‌کنیم،

$$F = \frac{\{\hat{\beta}_r' (X_r' \Lambda X_r)^{-1} \hat{\beta}_r\} / (k-1)}{e'e / (n-k)} \sim F(k-1, n-k) . \quad (۲)$$

۱. به پیوست «۵-الف» مراجعه شود.

با مراجعه به معادله ۵.۸۱ می‌دانیم

$$\tilde{\beta}'_r X'_r \wedge X_r \hat{\beta}_r = \hat{y}' \hat{y} = ESS,$$

که در آن \hat{y} برحسب انحراف از میانگین است؛ بنابراین با جایگزینی در معادله (۲) داریم

$$F = \frac{ESS / (k - 1)}{RSS / (n - k)} \sim F(k - 1, n - k).$$