



دانشگاه پیام نور

مرکز تهران

دانشکده علوم پایه

پایان نامه

برای دریافت درجه کارشناسی ارشد در رشته آمار ریاضی

عنوان پایان نامه:

برآورد پارامترهای مدل‌های رگرسیون مرتبه ای

با روشهای متفاوت

استاد راهنمای اول:

دکتر علی شادرخ

استاد راهنمای همکار:

دکتر پرویز نصیری

نگارش:

علی خاشعی ورname=خواستی

تیر ماه ۱۳۸۹

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِيْمِ



دانشگاه سام نور
مرکز تهران

دانشکده علوم پایه

پایان نامه

برای دریافت درجه کارشناسی ارشد
در رشته آمار ریاضی

عنوان پایان نامه:

برآورد پارامترهای مدل‌های رگرسیون مرتبه ای با روشهای متفاوت

استاد راهنمای اول:

دکتر علی شادرخ

استاد راهنمای همکار:

دکتر پرویز نصیری

نگارش:

علی خاشعی ورname=خواستی

تیر ماه ۱۳۸۹



جمهوری اسلامی ایران
وزارت علوم تحقیقات و فناوری

جمعیت علم مپاپ و کشاورزی



تصویب نامه

پایان نامه کارشناسی ارشد در رشته: آمار ریاضی

تحت عنوان:

"برآورد پارامترهای مدل‌های رگرسیون مرتبه ای
با روشهای متفاوت"

ساعت: ۱۰-۱۱

تاریخ دفاع: ۸۹/۰۴/۲۹

درجه ارزشیابی: ۱۸۱۵ نمره پایان نامه: ۱۷۰ رُمم

اعضاء	مرتبه علمی	نام و نام خانوادگی	اساتیدهای داوران
		دکتر علی شادرخ	استاد راهنمای
		دکتر پرویز نصیری	استاد مشاور
		دکتر احسان جمالی	استاد داور
		دکتر مسعود یار محمدی	نمائنده علمی گروه

تهران، خیابان استاد
نجات الهی، خیابان
شهید فلاح پور، پلاک
۲۷
تلفن: ۰۲۵۲-۸۸۸۰۰
دورنگار: ۸۸۳۱۹۴۷۵
www.tpnu.ac.ir
science.agri@tpnu.ac.ir

باتشکر و قدردانی از اساتید گرامی

جناب آقای دکتر علی شادرخ

و

جناب آقای دکتر پرویز نصیری

که مرا در تهیه و تنظیم این تحقیق یاری نموده و اینجانب را مورد لطف
و عنایت خود قرار دادند.

نام خانوادگی دانشجو: خاشعی ورنامخواستی نام: علی
عنوان پایان نامه: برآورد پارامترهای مدل‌های رگرسیون مرتبه ای با روشهای متفاوت
استاد راهنما: دکتر علی شادرخ
استاد راهنما همکار: دکتر پرویز نصیری
قطعه تحصیلی: کارشناسی ارشد گرایش: آمار ریاضی دانشگاه: پیام نور مرکز تهران
دانشکده: علوم پایه تاریخ فارغ التحصیلی: تعداد صفحه:

چکیده:

در بسیاری از مطالعات علمی و پژوهشی مایل به بررسی تأثیرگذاری متغیرها بر هم و یا پیش‌بینی یک متغیر به وسیله گروهی از متغیرهای دیگر هستیم. در مباحث آماری این عمل اغلب با استفاده از مدل‌های رگرسیونی امکان پذیر است. اما برای استفاده از این مدل‌های رگرسیونی پیش فرضهای نیاز است. از جمله این پیش فرضها استقلال مشاهدات و داده هاست، که باید مدنظر گرفته شود. البته ممکن است که این فرض برای همه مشاهدات بر قرار نباشد. یعنی مشاهدات نسبت به هم وابسته بوده و یا ساختاری تودرتو داشته باشند. پس در این شرایط نیاز به استفاده از روشهای رگرسیونی مناسب دیگری غیر از روشهای استاندارد می‌باشد.

در حالت وجود همبستگی بین مشاهدات و یا داشتن ساختاری تودرتو، مدل‌های جدیدی به نام مدل‌های رگرسیون مرتبه ای و یا به عبارت دیگر رگرسیون سلسله مراتبی می‌تواند برای برآش داده ها و کسب نتایج قابل قبولتر، مفید باشد. در این تحقیق به معرفی مدل‌های رگرسیون مرتبه ای (سلسله مراتبی) و پارامترهای موجود در آن پرداخته ایم. سپس با استفاده از روشهای مختلف آماری از قبیل برآوردهای درستنمایی ماکسیمم، کمترین مربعات، بیز تجربی، نیم بیز و تمام بیز پارامترها را برآورد کرده و مورد بررسی و مقایسه قرار داده ایم و سعی در بهبود و اصلاح آنها نموده ایم.

در نهایت با ارائه چند نمونه عملی و کاربردی از داده های واقعی، تلاش کرده ایم که با مدل‌سازی رگرسیون مرتبه ای و برآورد پارامترهای آن، نتایج قابل قبولتری نسبت به مدل‌های رگرسیون معمولی نشان دهیم. نتایج کامل این تحقیق در متن پایان نامه مورد توجه قرار گرفته است.

واژه های کلیدی: رگرسیون سلسله مراتبی - داده های تودرتو - برآورد کمترین مربعات - برآورد

درستنمایی ماکسیمم - برآورد بیز تجربی - برآورد نیم بیز - برآورد تمام بیز.

فهرست مطالب

۱	فصل اول: مدل‌های رگرسیونی
۲	۱-۱ مقدمه
۳	۲-۱ رگرسیون خطی ساده
۴	۱-۲-۱ مدل رگرسیون خطی
۵	۲-۲-۱ تابع رگرسیون خطی
۶	۳-۲-۱ برآورد α و β و σ^2
۷	۴-۱ رگرسیون چندگانه
۸	۱-۳-۱ مدل با دو متغیر مستقل
۹	۲-۳-۱ مدل با k متغیر مستقل
۱۰	۳-۳-۱ برآورد β و σ^2
۲۵	فصل دوم: رگرسیون مرتبه ای (سلسله مراتبی)
۲۶	۱-۲ مقدمه
۲۹	۲-۲ ساختارهای تودرتوی داده ها
۳۲	۱-۲-۲ واپسگی داده ها
۳۲	۲-۲-۲ رفتار غیر چند سطحی از داده های تودرتو
۳۴	۳-۲ مدل رگرسیون دو سطحی
۳۷	۴-۲ زیر مدل‌های ساده تر
۳۷	۱-۴-۲ ANOVA یکطرفه با اثرهای تصادفی
۳۹	۲-۴-۲ روش رگرسیون میانگینها

۴۰	یکطرفه با اثرهای تصادفی ANCOVA ۳-۴-۲
۴۱	۴-۴-۲ مدل رگرسیون ضرایب تصادفی
۴۳	۵-۴-۲ روش رگرسیون عرض از مبدأها و شیبهایا
۴۳	۶-۴-۲ مدل با شیبهای مختلف غیر تصادفی
۴۵	۵-۲ مدل خطی چندگانه سلسله مراتبی
۴۵	۱-۵-۲ چندین X و چندین W
۴۷	۲-۵-۲ تعمیم ساختارهای خطای سطح ۱ و سطح ۲
۴۸	۶-۲ هدف از آنالیز چند سطحی

۴۹	فصل سوم: برآوردهای رگرسیون سلسله مراتبی
۵۰	۱-۳ مقدمه
۵۲	۲-۳ برآوردهای ثابت
۵۲	۱-۲-۳ برآوردهای نقطه‌ای
۵۵	۲-۲-۳ برآوردهای فاصله‌ای
۵۶	۳-۲-۳ داده‌ها با رتبه ناقص
۵۸	۳-۳ برآوردهای ضرایب تصادفی سطح ۱
۵۸	۱-۳-۳ برآوردهای درستنمائی ماکسیمم
۶۰	۲-۳-۳ برآوردهای بیزی
۶۳	۳-۳-۳ برآوردهای بیزی تجربی (EB)
۶۵	۴-۳-۳ برآوردهای نیم بیزی (SB)
۶۶	۵-۳-۳ برآوردهای تماماً بیزی (FB)
۶۶	۶-۳-۳ برآوردهای فاصله‌ای
۶۷	۷-۳-۳ قابلیت اعتماد ضرایب برآوردهای سطح ۱
۶۸	۸-۳-۳ مشاهده باقیمانده

۷۰	فصل چهارم: کاربرد مدل‌های رگرسیون سلسله مراتبی
۷۱	۱-۴ مقدمه

۷۴	۲-۴ تحقیق در مورد رابطه موفقیت و وضعیت اجتماعی اقتصادی در مدارس
۷۴	۱-۲-۴ بررسی رابطه موفقیت و وضعیت اجتماعی اقتصادی در یک مدرسه
۷۶	۲-۲-۴ بررسی رابطه موفقیت و وضعیت اجتماعی اقتصادی در دو مدرسه
۷۷	۳-۲-۴ بررسی رابطه موفقیت و وضعیت اجتماعی اقتصادی در ۱ مدرسه
	۳-۴ ارتباط بین درس ریاضیات و معدل درسی دانش آموزان مدارس راهنمائی شهرستان
۸۱	لنجان
۸۵	۱-۳-۴ پکترفه ANOVA
۸۹	۲-۳-۴ روش رگرسیون میانگینها
۹۲	۳-۳-۴ مدل ضرایب تصادفی
۹۹	۴-۳-۴ مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای $Type$ و Nst برای سطح مدرسه
۱۱۳	۵-۳-۴ مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای $Type$ و $Space$ برای سطح مدرسه
۱۲۲	۶-۳-۴ مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای $Type$ ، Nst ، $Space$ و $Space$ برای سطح مدرسه
۱۳۰	۴-۴ نتیجه گیریها

منابع

۱۳۶	پیوست
۱۳۶	پیوست الف: جدول مربوط به بخش ۳-۴

فهرست جداول

- ۱-۲ نمونه هائی از واحدها در سطح ۱ و سطح ۲.
- ۲-۲ خلاصه ای از عبارات بیان شده برای وصف واحدها در هر یک از دو سطح.
- ۱-۴ نتایج از رگرسیون نمره ریاضیات و معدل ۲۰۰ دانش آموز.
- ۲-۴ آماره های توصیفی برای داده های سطح دانش آموز.
- ۳-۴ آماره های توصیفی برای سطح مدرسه.
- ۴-۴ برآورد ضرایب سطح ۱ در مدل ANOVA یکطرفه.
- ۵-۴ نتایج از مدل ANOVA یکطرفه.
- ۶-۴ نتایج از روش رگرسیون میانگینها.
- ۷-۴ برآورد ضرایب سطح دانش آموز در مدل ضرایب تصادفی.
- ۸-۴ نتایج از مدل ضرایب تصادفی.
- ۹-۴ قابلیت اعتماد ضرایب سطح ۱ برای مدل ضرایب تصادفی.
- ۱۰-۴ نتایج از مدل با پیش بین *MAT* برای سطح دانش آموز و پیش بینهای *Type* و *N St* برای سطح مدرسه.
- ۱۱-۴ قابلیت اعتماد ضرایب سطح ۱ برای مدل مدل با پیش بین *MAT* برای سطح دانش آموز و پیش بینهای *Type* و *N St* برای سطح مدرسه.
- ۱۲-۴ برآوردهای بیزی تجربی ضرایب سطح دانش آموز برای مدل با پیش بین *MAT* برای سطح دانش آموز و پیش بینهای *Type* و *N St* برای سطح مدرسه.
- ۱۳-۴ قابلیت اعتماد ضرایب سطح ۱ برای مدل با پیش بین *MAT* برای سطح دانش آموز و پیش بینهای *Type* و *N St* برای سطح مدرسه با توجه به $\tau_{..} = 0.5$ و $\tau_{11} = 0.005$.

۱۴-۴ برآوردهای نیم بیز ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش آموز.

۱۰۹ و پیش بینهای $N St Type$ و برای سطح مدرسه.

۱۵ ۴ برآوردهای تمام بیز ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش

۱۱۱ آموز و پیش بینهای $N St Type$ و برای سطح مدرسه.

۱۶-۴ نتایج از مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای $Space Type$ و $St Type$ برای

۱۱۴ سطح مدرسه.

۱۷-۴ قابلیت اعتماد ضرایب سطح ۱ برای مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای

۱۱۶ $Space Type$ و برای سطح مدرسه.

۱۸-۴ برآوردهای بیز تجربی ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش

۱۱۷ آموز و پیش بینهای $Space Type$ و $St Type$ برای سطح مدرسه.

۱۹-۴ برآوردهای نیم بیز ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش آموز

۱۱۹ و پیش بینهای $Space Type$ و $St Type$ برای سطح مدرسه.

۲۰-۴ برآوردهای نیم بیز ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش آموز

۱۲۱ و پیش بینهای $Space Type$ و $St Type$ برای سطح مدرسه.

۲۱-۴ نتایج از مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای $St N. Type$ و $Space$

۱۲۲ برای سطح مدرسه.

۲۲-۴ قابلیت اعتماد ضرایب سطح ۱ برای مدل با پیش بین MAT برای سطح دانش آموز و پیش بینهای

۱۲۵ $Space Type$ و $St N. Type$ برای سطح مدرسه.

۲۳-۴ برآوردهای بیز تجربی ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش

۱۲۶ آموز و پیش بینهای $Space Type$ و $St N. Type$ برای سطح مدرسه.

۲۴-۴ برآوردهای نیم بیز ضرایب سطح دانش آموز برای مدل با پیش بین MAT برای سطح دانش

۱۲۸ آموز و پیش بینهای $Space Type$ و $St N. Type$ برای سطح مدرسه.

الف-۱ مشخصات مربوط به نمونه انتخاب شده از مدارس (نوع مدرسه، دولتی (۰) و غیردولتی (۱)،

جنسیت، پسر (۰) و دختر (۱)، تعداد دانش آموزان مشغول به تحصیل در آن و مقدار فضای

آموزشی که مدرسه در اختیار دارد بر حسب مترمربع) و دانش آموزان (نمره درس ریاضیات و

۱۳۶ معدل درسی) مقطع راهنمائی شهرستان لنجان در سال تحصیلی ۸۷-۸۸

فهرست اشکال و نمودارها

- ۱-۲ ساختار داده ها برای یک مدل سلسله مراتبی دو سطحی.
۲-۲ سازگاری سطح خرد در مقابل سطح کلان.
- ۱-۴ نمودار پراکنش ارتباط بین موفقیت و SES در یک مدرسه فرضی.
۲-۴ نمودار پراکنش ارتباط بین موفقیت و SES (متمرکز شده) در یک مدرسه فرضی.
۳-۴ نمودار پراکنش ارتباط بین موفقیت و SES در دو مدرسه فرضی.
۴-۴ نمودار میانگینها و شبیهای برای ۲۰۰ مادرسه فرضی.
۵-۴ نمودار پراکنش بین درس ریاضیات و معدل ۲۰۰ دانش آموز.
۶-۴ نمودار ارتباط بین میانگین درس ریاضیات و میانگین معدل برای مدارس.
۷-۴ نمودار ارتباط بین ریاضیات و معدل درسی برای مدرسه ۶.
۸-۴ نمودار ارتباط بین ریاضیات و معدل درسی برای مدرسه ۱.
۹-۴ ارتباط بین معدل و ریاضیات برای مدارس دولتی (۰) و غیردولتی (۱).
۱۰-۴ مانده های کمترین مریعات معمولی از عرض از مبدأ و شبیه برای مدارس در مدل با پیش بینهای $St N$ و $Type$ برای سطح مدرسه.
۱۱-۴ مانده های بیز تجربی از عرض از مبدأ و شبیه برای مدارس در مدل با پیش بینهای $Type$ و $St N$ برای سطح مدرسه.
۱۲-۴ مانده های نیم بیز از عرض از مبدأ و شبیه برای مدارس در مدل با پیش بینهای $Type$ و $St N$ برای سطح مدرسه.
۱۳-۴ مانده های تمام بیز از عرض از مبدأ و شبیه برای مدارس در مدل با پیش بینهای $Type$ و $St N$ برای سطح مدرسه.

۱۴-۴ مانده های کمترین مرتبات معمولی از عرض از مبدأ و شیب برای مدارس در مدل با پیش

۱۱۸ بینهای *Type* و *Space* برای سطح مدرسه.

۱۵-۴ مانده های بیز تجربی از عرض از مبدأ و شیب برای مدارس در مدل با پیش بینهای *Type* و

۱۱۸ *Space* برای سطح مدرسه.

۱۶-۴ مانده های کمترین مرتبات معمولی از عرض از مبدأ و شیب برای مدارس در مدل با پیش

۱۲۷ بینهای *St N.* *Type* و *Space* برای سطح مدرسه.

۱۷-۴ مانده های بیز تجربی از عرض از مبدأ و شیب برای مدارس در مدل با پیش بینهای *Type*

۱۲۷ *Space* و *St N.* برای سطح مدرسه.

۱

مدلهای رگرسیونی

۱-۱ مقدمه

تحلیل رگرسیونی ما را قادر می سازد که بین یک متغیر مورد نظر که آن را متغیر وابسته یا پاسخ می نامیم و یک یا چند متغیر مستقل یا پیش بین رابطه ای را معین کنیم و از آن استفاده نمائیم. مثلاً، در یک مطالعه رگرسیونی، متغیر پاسخ، قد افراد بر حسب سانتیمتر و متغیر مستقل، وزن افراد بر حسب کیلوگرم در یک زمان معین است.

تحلیل رگرسیونی را غالباً برای پیش بینی متغیر پاسخ و یا اساساً برای بررسی وابستگی بین متغیرهای مستقل و متغیر پاسخ به کار می برند. مثلاً، مایلیم قد افراد را با استفاده از آگاهی‌ی که از وزن آنها داریم پیش بینی کنیم و یا بررسی کنیم که آیا با افزایش وزن افراد، قد آنها نیز افزایش می یابد یا نه.

عموماً متغیر پاسخ با Y و متغیر مستقل با X نشان داده می شود. همانطور که بیان شد تحلیل رگرسیونی به ما این امکان را می دهد که بین دو متغیر X و Y رابطه ای برقرار کنیم. وقتی مقدار X مشخص باشد و مقدار Y به صورتی یکتا معین شود، بین X و Y رابطه ای دقیق برقرار شده است که این رابطه، یک رابطه تابعی بین دو متغیر X و Y است. اگر براساس مقادیر X و Y نمودار پراکنشی رسم شود، بدلیل اینکه مقدار Y به صورتی یکتا از روی مقدار X تعیین شده، تمام مشاهدات بر خط رابطه بین X و Y می افتد. حال اگر مقدار X مشخص باشد و مقدار Y به صورتی یکتا معین نشود، رابطه بین X و Y دقیق نخواهد بود که این رابطه یک رابطه آماری بین دو متغیر X و Y است. برای توصیف این رابطه اگر نمودار پراکنش مشاهدات رسم شود و از محل تمرکز و تجمع نقاط خطی رسم شود، دیده می شود که این رابطه، رابطه ای کامل نیست و

مشاهدات حول این خط پراکنده اند. این دو نوع رابطه تابعی و آماری بین دو متغیر، نسبت به مشاهدات X و Y به صورت خطی یا غیر خطی نشان داده می شوند. در این فصل، از برخی مفاهیم پایه ای تحلیل رگرسیونی بحث خواهیم کرد.

۲-۱ رگرسیون خطی ساده

اگر رابطه بین دو متغیر X و Y آماری باشد، مدل رگرسیونی برای مشاهده Y_i به صورت زیر است:

$$Y_i = \alpha + \beta X_i + \varepsilon_i \quad (1-1)$$

مؤلفه $\alpha + \beta X_i$ که رابطه آماری را منعکس می کند، مؤلفه رگرسیونی می خوانند و مؤلفه ε_i که پراکنش تصادفی را منعکس می کند، خطأ می نامند. Y_i هم متغیری تصادفی است، به این دلیل که هر مشاهده Y_i ، شامل یک مؤلفه تصادفی ε_i است.

۲-۱ مدل رگرسیون خطی^{۱)}

حال مدل رگرسیونی را که در آن رابطه آماری خطی است مورد بررسی قرار می دهیم:

$$Y_i = \alpha + \beta x_i + \varepsilon_i \quad i = 1, 2, \dots, n \quad (2-1)$$

که در آن:

- Y_i ، پاسخ در i امین مورد است.
- x_i ، مقدار متغیر مستقل در i امین مورد، مقداری ثابت و معلوم فرض شده است.
- α و β ، پارامترنند.

¹⁾) Linear regression model

• ε_i ها، متغیرهای مستقل دارای توزیع نرمال $N(0, \sigma^2)$ هستند.

مدل (۲-۱) را مدل رگرسیون خطی ساده می‌نامند. این مدل شامل سه پارامتر α ، β و σ^2 می‌باشد.

۲-۲-۱ تابع رگرسیون خطی^{۱)}

می‌خواهیم مقدار امید $E(Y_i)$ را بدست آوریم. با توجه به شرایطی که در بالا ذکر کردیم، داریم:

$$E(Y_i|x) = E(\alpha + \beta x_i + \varepsilon_i) = E(\alpha + \beta x_i) + E(\varepsilon_i) = \alpha + \beta x_i + E(\varepsilon_i)$$

از آنجاییکه $E(\varepsilon_i) = 0$ است، در پایان با توجه به معلوم بودن x ، به نتیجه زیر می‌رسیم.

$$E(Y_i|x) = \alpha + \beta x_i \quad (3-1)$$

رابطه بین x و $E(Y_i)$ را تابع رگرسیونی می‌نامند. پارامتر α ، عرض از مبدأ خط رگرسیونی و پارامتر β ، هم شیب خط رگرسیونی است.

حال می‌خواهیم واریانس Y_i را بیابیم.

$$\begin{aligned} Var(Y_i|x) &= Var(\alpha + \beta x_i + \varepsilon_i) \\ \text{از آنجاییکه } \alpha + \beta x_i &\text{ مقداری ثابت است، داریم:} \end{aligned}$$

$$Var(Y_i|x) = Var(\varepsilon_i) = \sigma^2 \quad (4-1)$$

پس صرفنظر از مقادیر x_i ها، Y_i ها تغییرپذیری یکسانی دارند.

چون Y_i تابعی خطی از متغیر تصادفی نرمال ε_i است، می‌توان گفت که هر Y_i هم به صورت نرمال توزیع شده است و سرانجام، چون ε_i ها برای مشاهدات مختلف مستقل فرض شده بودند، Y_i ها نیز چنین خواهند بود. بنابراین Y_i ها متغیرهای تصادفی نرمال مستقلی هستند.

¹⁾) Linear regression function

$$Y_i \sim N(\alpha + \beta x_i, \sigma^2) \quad (5-1)$$

۳-۲-۱ براورد α و β و σ^2

معمولًاً پارامترهای رگرسیونی α و β و σ^2 مجهول اند و باید از داده‌های نمونه براورد شوند. بنابراین برای براورد این پارامترها از براوردهایی مانند براوردهای کمترین مربعات^۱ و براوردهای درستنمایی ماکسیمم^۲ استفاده خواهیم کرد.

الف) براوردهای کمترین مربعات

اگر نمودار پراکنشی از مشاهدات را در نظر بگیریم و بخواهیم خط راستی را به عنوان بهترین برازش از این داده‌ها رسم کنیم، روش کمترین مربعات، مجموع توانهای دوم انحرافهای مشاهدات Y از خط راست را بعنوان معیار برازش در نظر می‌گیرد. در اینجا این مجموع را با Q نشان می‌دهیم. روش کمترین مربعات، خطی را بعنوان بهترین خط راست برازش داده شده در نظر می‌گیرد که برای آن مقدار Q یعنی مجموع توانهای دوم انحرافها کمترین مقدار ممکن را نشان دهد.

براوردهای کمترین مربعات α و β را به ترتیب با a و b نشان می‌دهیم، بنابراین خط رگرسیونی برازش داده شده با براوردهای کمترین مربعات برابر خواهد بود با:

$$\hat{Y}_i = a + bX_i$$

پس برای مقدار Q داریم:

$$Q = \sum_{i=1}^n [Y_i - (a + bX_i)]^2 \quad (6-1)$$

مطابق با اصل کمترین مربعات، مقادیر a و b را باید طوری تعیین کنیم که Q کمینه شود.

^۱) Least squares estimator

^۲) Maximum likelihood estimator

برای بدست آوردن برآوردهای a و b ، یک راه تحلیلی وجود دارد. برای ارائه ساده‌تر این راه حل مجموعهای زیر را تعریف می‌کنیم:

$$S_X^r = \sum(X_i - \bar{X})^r = \sum X_i^r - n\bar{X}^r \quad (7-1)$$

$$S_Y^r = \sum(Y_i - \bar{Y})^r = \sum Y_i^r - n\bar{Y}^r \quad (8-1)$$

$$S_{XY}^r = \sum(X_i - \bar{X})(Y_i - \bar{Y}) = \sum X_i Y_i - n\bar{X}\bar{Y} \quad (9-1)$$

که در آن،

• \bar{Y} و \bar{X} ، میانگینهای مقادیر i ها و Y_i ها،

• S_Y^r و S_X^r ، مجموع توانهای دوم انحرافها از میانگینها،

• S_{XY}^r ، مجموع حاصلضربهای برداری انحرافها،

ابتدا می‌نویسیم

$$Y_i - (a + bX_i) = (Y_i - \bar{Y}) - b(X_i - \bar{X}) + (\bar{Y} - a - b\bar{X})$$

با مربع کردن دو طرف داریم:

$$(Y_i - a - bX_i)^r = (Y_i - \bar{Y})^r + b^r(X_i - \bar{X})^r + (\bar{Y} - a - b\bar{X})^r$$

$$-rb(X_i - \bar{X})(Y_i - \bar{Y}) - rb(X_i - \bar{X})(\bar{Y} - a - b\bar{X}) + r(Y_i - \bar{Y})(\bar{Y} - a - b\bar{X})$$

هر دو طرف رابطه را بر روی مقادیر $i = 1, 2, \dots, n$ جمع بندی می‌کنیم. چون

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

$$\sum(Y_i - \bar{Y}) = 0$$

با توجه به معادلات (7-1)، (8-1) و (9-1) بنابراین داریم:

$$Q = S_Y^r + b^r S_X^r + n(\bar{Y} - a - b\bar{X})^r - rb S_{XY}^r$$

می‌توانیم جملات را اینگونه نیز مرتب کنم،

$$Q = n(\bar{Y} - a - b\bar{X})^r + (b^r S_X^r - rb S_{XY}^r + \frac{S_{XY}^r}{S_X^r}) + S_Y^r - \frac{S_{XY}^r}{S_X^r}$$

$$Q = n(\bar{Y} - a - b\bar{X})^2 + (b^2 S_X^2 - 2b S_{XY} + \frac{S_{XY}^2}{S_X^2}) + S_Y^2 - \frac{S_{XY}^2}{S_X^2}$$

اگر a و b موجود در رابطه را برابر با $a = \bar{Y} - b\bar{X}$ و $b = \frac{S_{XY}}{S_X^2}$ قرار دهیم، دو جمله اول را به کوچکترین مقدار، یعنی صفر کاهش خواهیم داد.

بنابراین برآوردهای کمترین مربعات برای α و β که از اینجا به بعد آنها را با $\hat{\alpha}$ و $\hat{\beta}$ نشان می‌دهیم، عبارت اند از:

$$\hat{\beta} = \frac{S_{XY}}{S_X^2} \quad (10-1)$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X} \quad (11-1)$$

که با توجه به تعریف S_{XY} و S_X^2 ، برای (10-1) داریم

$$\hat{\beta} = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2}$$

سرانجام خط رگرسیون کمترین مربعات خواهد بود از:

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i \quad (12-1)$$

معادله (12-1) برآورده از تابع رگرسیونی $E(Y_i|X) = \alpha + \beta X_i$ است که آن را تابع رگرسیونی برآورد شده می‌نامند.

بطور کلی در روش کمترین مربعات، برآوردهای $\hat{\alpha}$ و $\hat{\beta}$ را به قسمی پیدا می‌کنیم که مجموع توانهای دوم انحرافها $Y_i - \hat{Y}_i$ را مینیمم سازد. برای این مینیمم سازی روش دیگری نیز وجود دارد که در اینجا به ذکر آن می‌پردازیم.

برای یافتن مقادیر $\hat{\alpha}$ و $\hat{\beta}$ که Q را مینیمم سازد، نسبت به $\hat{\alpha}$ و $\hat{\beta}$ مشتق می‌گیریم و حاصل را برابر صفر قرار می‌دهیم: