



حمایت از حقوق پدیدآورندگان

پایان نامه حاضر، حاصل پژوهشهای نگارنده در دوره کارشناسی ارشد رشته آمار گرایش آمار ریاضی است که در بهمن ۱۳۹۲ در دانشکده علوم پایه دانشگاه یاسوج به راهنمایی دکتر حیدرعلی مردانی فرد و مشاوره دکتر انوشیروان غفاری پور از آن دفاع شده است و کلیه حقوق مادی و معنوی آن متعلق به دانشگاه یاسوج است.



دانشکده علوم پایه
گروه ریاضی

پایان نامه کارشناسی ارشد رشته آمار گرایش آمار ریاضی

برآورد کمترین مربعات جریمه‌دار در مدل‌های جزئاً خطی با ابعاد بالا

استاد راهنما

دکتر حیدرعلی مردانی فرد

پژوهشگر

زهرا فریدونی

بهمن ۱۳۹۲



برآورد کمترین مربعات جریمه‌دار در مدل‌های جزئاً خطی با ابعاد بالا

به وسیله

زهرا فریدونی

پایان نامه

ارائه شده به تحصیلات تکمیلی دانشگاه به عنوان بخشی از فعالیت‌های تحصیلی لازم برای اخذ

درجه کارشناسی ارشد

در رشته:

آمار

در تاریخ توسط هیأت داوران زیر بررسی و با درجه به تصویب نهایی رسید.

- | | | | |
|------------------------------------|--------------------------|------------------------|-------|
| ۱- استاد راهنما: | دکتر حیدرعلی مردانی فرد | با مرتبه علمی استادیار | امضاء |
| ۲- استاد مشاور: | دکتر انوشیروان غفاری پور | با مرتبه علمی استادیار | امضاء |
| ۳- استاد داور داخل گروه: | دکتر آرش اردلان | با مرتبه علمی استادیار | امضاء |
| ۴- استاد داور خارج گروه: | دکتر علیرضا نعمت‌الهی | با مرتبه علمی استادیار | امضاء |
| ۵- نماینده تحصیلات تکمیلی دانشگاه: | دکتر زهرا رفیعی | با مرتبه علمی استادیار | امضاء |

تقدیم به:

همسرم، اسطوره زندگی ام، پناه خستگی ام

و

فرزندم، امید بودنم

قدردانی

سپاس خدای را که سخنوران در ستودن او بمانند و شمارندگان شمردن نعمت‌های او ندانند و کوشندگان حق او را گزاردن نتوانند و سلام و درود بر محمد و خاندان پاک او، طاهران معصوم، هم‌آنان که وجودمان وامدار وجودشان است. وظیفه خود می‌دانم از استاد راهنمای خود، جناب آقای دکتر حیدرعلی مردانی‌فرد، که در کمال سعه صدر، حسن خلق و فروتنی از هیچ کمکی در این عرصه بر من دریغ نمودند صمیمانه تشکر و قدردانی کنم. از استاد صبور، جناب آقای دکتر انوشیروان غفاری پور که زحمت مشاوره این پایان‌نامه را بر عهده گرفتند کمال تشکر و قدردانی را دارم. از استاد گرامی دکتر آرش اردلان سپاسگزارم، چرا که زحمت داوری این پایان‌نامه را در حالی بر عهده گرفتند که بدون راهنمایی‌های ارزنده ایشان این مجموعه به انجام نمی‌رسید.

بر خود لازم می‌دانم از استاد عالی‌قدر دانشگاه شیراز جناب آقای دکتر نعمت‌اللهی قدر دانی نمایم. به پاس قدردانی با قلبی آکنده از عشق و معرفت از همسرم و پسرم که صبورانه و صادقانه مرا همراهی نمودند تا بتوانم در کمال آرامش و آسایش به تهیه و تنظیم این اثر بپردازم تشکر و قدردانی می‌نمایم. باشد که این خردترین، بخشی از زحمات آنان را سپاس گوید.

زهرا فریدونی

بهمن ۱۳۹۲

چکیده

در مطالعات رگرسیونی، زمانی که بین متغیرهای مستقل همبستگی بالایی وجود داشته باشد استفاده از روش‌های معمول از جمله روش کمترین مربعات معمولی باعث ناپایداری واریانس برآوردها می‌شود. یک راه حل معمول، استفاده از روش کمترین مربعات جریمه‌دار است که در آن برای مقادیر بزرگ برآوردها، جریمه بالایی در نظر گرفته می‌شود و به نوعی تغییرات برآوردها تحت کنترل در می‌آید.

مورد دیگر استفاده از رگرسیون جریمه‌دار در مدل‌های با ابعاد بالا یعنی مدل‌هایی با تعداد زیادی متغیر مستقل است. در این مدل‌ها تلاش می‌شود از ضرائب "نزدیک به صفر" حتی‌الامکان صرف‌نظر گردد تا فقط متغیرهایی در مدل باقی بمانند که تاثیر کاملاً معنی‌داری در متغیر پاسخ دارند.

در این پژوهش تلاش شده است ضمن مرور مختصری بر روش کمترین مربعات جریمه‌دار، رگرسیون جریمه‌دار و نحوه عمل این روش در هموارسازی نمودارهای رگرسیونی، استفاده از این روش در برازش مدل‌هایی با ابعاد بالا مورد مطالعه و بررسی قرار گیرد. با استفاده از شبیه‌سازی، درستی بعضی از روابط و برتری این روش در مقایسه با سایر روش‌ها تحقیق شده است که با ارائه چند سری داده واقعی، این روش در تحلیل آنها مورد استفاده قرار گرفته است.

فهرست مطالب

iii	فهرست علائم اختصاری
iv	فهرست تصاویر
۱	فصل ۱: مقدمه و تاریخچه
۴	فصل ۲: مفاهیم و تعاریف
۱۹	فصل ۳: رگرسیون جریمه‌دار
۱۹	۱-۳ برآورد کمترین مربعات جریمه‌دار
۲۱	۲-۳ مدل‌های برآورد کمترین مربعات جریمه‌دار
۲۲	۱-۲-۳ رگرسیون جریمه‌دار ریج
۳۱	۲-۲-۳ رگرسیون جریمه‌دار لاسو
۳۶	۳-۲-۳ رگرسیون جریمه‌دار اسکاد
۳۸	۴-۲-۳ انتخاب پارامتر همواری
۴۲	فصل ۴: چندجمله‌ایها و رگرسیون تکه تکه
۴۲	۱-۴ چندجمله‌ایها
۴۳	۲-۴ چندجمله‌ای تکه‌ای
۴۴	۳-۴ رگرسیون چندجمله‌ای
۵۰	۴-۴ رگرسیون شکسته
۵۱	۵-۴ مدل‌های رگرسیون تکه تکه
۵۳	۶-۴ کاربرد رگرسیون جریمه‌دار ریج در هموارسازی اسپلاین

۵۷	فصل ۵: برآوردهای جریمه دار در مدل‌های جزئاً خطی و ابعاد بالا
۵۷	۱-۵ مدل‌های نیمه پارامتری
۶۲	۱-۱-۵ برآورد پارامترها در مدل‌های جزئاً خطی
۶۳	۲-۱-۵ هموارسازی در مدل‌های جزئاً خطی
۷۰	۲-۵ برآورد در مدل‌های جزئاً خطی با ابعاد بالا
۷۴	فصل ۶: شبیه‌سازی
۷۴	۱-۶ برآورد پارامتر برای داده‌های شبیه‌سازی شده
۷۴	۱-۱-۶ برآورد پارامتر توسط روش انتخاب بهترین زیرمجموعه
۷۵	۲-۱-۶ برآورد پارامترها توسط رگرسیون ریبج و لاسو
۷۸	۳-۱-۶ برآورد پارامترها به روش ریبج در ابعاد بالا
۷۹	پیوست آ: تعاریف و مفاهیم ریاضی
۹۱	پیوست ب: برنامه‌های نوشته شده برای برآورد پارامترها با استفاده از نرم‌افزار R
۹۱	ب-۱ برنامه شبیه‌سازی شده
۹۲	ب-۲ برنامه نمودارهای فصل‌های قبل
۱۱۰	واژه‌نامه فارسی به انگلیسی
۱۱۱	واژه‌نامه انگلیسی به فارسی
۱۱۲	مراجع

فهرست علائم اختصاری

<i>MSE</i>	میانگین مربع خطا
<i>PE</i>	خطای پیش‌بینی
<i>MAE</i>	میانگین قدر مطلق خطا
<i>SSE</i>	مجموع توان دوم خطا
<i>OLS</i>	کمترین مربعات معمولی
<i>SS_R</i>	مجموع توان دوم رگرسیونی
<i>SS_T</i>	مجموع توان دوم کل
<i>MSR</i>	میانگین توان دوم رگرسیونی
<i>MSE</i>	میانگین توان دوم خطا
<i>CV</i>	اعتبارسنجی متقابل
<i>PRSS</i>	کمترین مربعات جریمه‌دار
<i>PLM</i>	مدل جزئاً خطی

فهرست تصاویر

۹	۱-۲ نمودار پراکنش برای داده‌های cars
۹	۲-۲ تابع موجدار برای داده‌های cars
۱۰	۳-۲ پایه رگرسیون خطی
۱۱	۴-۲ رگرسیون خطی
۱۱	۵-۲ پایه مدل رگرسیون درجه دو
۱۲	۶-۲ مدل رگرسیون درجه دو
۱۳	۷-۲ مدل با یک شکستگی
۱۳	۸-۲ پایه مدل با یک شکستگی
۱۴	۹-۲ پایه مدل با چندین شکستگی
۱۷	۱۰-۲ R^2 برای داده‌های پروستات
۱۷	۱۱-۲ خطای پیش‌بینی کمترین مربعات و بهترین زیرمجموعه
۱۸	۱۲-۲ بهترین زیرمجموعه
۲۰	۱-۳ رگرسیون خطی و تابع موجدار
۲۶	۲-۳ فضای کانتور ریج
۲۸	۳-۳ خطای پیش‌بینی ریج و کمترین مربعات
۳۰	۴-۳ برآورد ریج در برابر کمترین مربعات
۳۲	۵-۳ فضای رگرسیون لاسو
۳۳	۶-۳ مقایسه لاسو و کمترین مربعات
۳۵	۷-۳ رگرسیون لاسو
۳۵	۸-۳ خطای پیش‌بینی لاسو و کمترین مربعات
۳۷	۹-۳ برآورد اسکاد

۴۰ GCV۱۰-۳
۴۱ CV ۱۱ نمودار ۳-۳
۴۳ ۱-۴ چند جمله‌ای تکه‌ای
۴۵ ۲-۴ چند جمله‌ای درجه دو برای داده‌های اتانول
۴۷ ۳-۴ رگرسیون خطی برای داده‌های لیدار
۴۷ ۴-۴ رگرسیون درجه دو برای داده‌های لیدار
۴۸ ۵-۴ رگرسیون درجه چهار برای داده‌های لیدار
۴۸ ۶-۴ رگرسیون درجه ده برای داده‌های لیدار
۴۹ ۷-۴ اسپلاین با یک گره برای داده‌های لیدار
۵۲ ۸-۴ اسپلاین درجه یک با ۳ گره برای cars
۵۳ ۹-۴ اسپلاین درجه دو
۵۴ ۱۰-۴ اسپلاین خطی با ۴ گره
۵۴ ۱۱-۴ اسپلاین خطی با ۲۰ گره
۵۶ ۱۲-۴ اسپلاین با $\lambda = 10$
۵۸ ۱-۵ نمودار پراکنش شوری تالاب
۵۹ ۲-۵ باقی‌مانده در برابر شوری آب تالاب
۶۰ ۳-۵ مقایسه نمودارهای خطی با حضور داده پرنفوذ
۶۰ ۴-۵ مقایسه منحنی‌ها با حضور داده پرنفوذ
۶۵ ۵-۵ نمودار پراکنش داده‌های ragweed
۶۷ ۶-۵ برازش با استفاده از اسپلاین در مدل نیمه‌پارامتری
۷۰ ۷-۵ نمودار اعتبارسنجی متقابل با استفاده از برازش لاسو
۷۳ ۸-۵ نمودار اعتبارسنجی متقابل با استفاده از برازش اسکاد
۸۶ ۱-آ تصویر بردار روی بردار دیگر

فصل ۱

مقدمه و تاریخچه

هدف استنباط آماری، نتیجه‌گیری در مورد یک یا چند جمعیت بر اساس اطلاعات موجود در نمونه است. یکی از مسائل مهم در استنباط آماری، بحث برآوردیابی می‌باشد که شامل سه شاخه

۱ - برآوردهای کلاسیک

۲ - برآوردهای بیزی

۳ - برآوردهای دنباله‌ای است

که هرکدام از برآوردهای کلاسیک و برآوردهای بیزی شامل برآوردهای نقطه‌ای و برآورد فاصله‌ای است. در برآورد نقطه‌ای برآوردگر باید تا حد امکان به پارامتر نامعلوم نزدیک باشد. از ساده‌ترین و در عین حال کاراترین روش‌های برآورد، روش کمترین مربعات معمولی است که اولین بار توسط کارل فردریش گاوس^۱ در سال ۱۷۹۴ میلادی مطرح شد. این روش برای برآورد پارامترهای رگرسیون خطی مرتبه‌های بالاتر از یک نیز قابل استفاده می‌باشد و از جمله کسانی که به استفاده از این روش در رگرسیون چندجمله‌ای پرداخت جرگونه^۲ می‌باشد. یکی از ایرادهای وارد به روش کمترین مربعات معمولی، توانائی تفسیر آن است. هنگامیکه با وضعیتی روبه‌رو هستیم که به دلیل تعداد زیاد پیش‌بینی کننده‌ها، شناسائی گروه کوچک و

^۱ gauss

^۲ Gergonne

موثری از متغیرها مفید می‌باشد، رگرسیون خطی چنین ویژگی را ندارد و همچنین رگرسیون خطی هنگامیکه تعداد متغیرها بیشتر از مشاهدات باشد، تعریف نشده است و چون در بسیاری از شاخه‌های علوم زیست‌شناسی و پزشکی با مواردی روبه‌رو هستیم که در آن تعداد متغیرها بیش از نمونه است اهمیت برآورد در چنین نمونه‌هائی آشکار می‌شود.

بنابراین رگرسیون جریمه‌دار، به عنوان یک توسیع از پیش‌بینی کننده‌های خطی معرفی می‌شود. رگرسیون جریمه‌دار ریچ توسط هورل^۳ در سال ۱۹۶۲ معرفی شد که در آن اظهار داشت، با وجود همبستگی بین متغیرهای توضیحی، بکارگیری روش حداقل مربعات باعث ایجاد خطا در برآورد می‌شود و رگرسیون ریچ را به عنوان جایگزینی برای این روش، توسعه داد که اجازه می‌دهد که برآوردها با واریانس کمتر از روش حداقل مربعات با داشتن اریبی، محاسبه شوند. اندازه‌گیری اصلی این روش در آمار توسط هورل و کنالد^۴ در سال ۱۹۷۰ مطرح شد.

در سال ۱۹۹۸ ولیلا^۵ و هسیج^۶ نشان دادند که اگر متغیرهای پیش‌بینی کننده از همبستگی بالائی برخوردار باشند، برآورد ریچ عملکرد خوبی نخواهد داشت بنابراین سعی شد از نوع دیگر رگرسیون جریمه‌دار که یک نوع روش انتخاب متغیرهاست و توسط تیشیرانی^۷ در سال ۱۹۹۶ معرفی و توسیع داده شد استفاده نمایند.

در سال ۲۰۰۱ فان و لی انتخاب متغیر از طریق درست‌نمایی جریمه‌شده غیرمقعر و ویژگی‌های آن را ارائه کردند [۱]. که ایکسی^۸ و هونگ^۹ برای انتخاب متغیر و برآورد ضرائب رگرسیون در یک مدل جزئاً خطی از آن استفاده نمودند. فان و پنگ در سال ۲۰۰۴ درست‌نمایی جریمه‌دار غیرمقعر با مشتق‌پذیری تعدادی از پارامترها را معرفی نمودند [۲]. همچنین هانگ و هورویتز در سال ۲۰۰۸ ویژگی‌های مجانبی از برآوردهای بریچ را در مدل‌های رگرسیونی بحث کردند [۳]. رگرسیون جریمه‌دار اسکد در مدل‌های جزئاً خطی توسط ایکس و هانتز بحث

Horel^۳Kenald^۴velilla^۵Hsich^۶Tibshirani^۷Xie^۸Huang^۹

شد [۴].

فصل ۲

مفاهیم و تعاریف

در این فصل مفاهیم و تعاریف مورد نیاز ارائه خواهند شد.

فرض می‌کنیم $\mathbf{X}_i = (x_{i1}, \dots, x_{ip})$ متغیرهای مستقل و Y_i که $1 \leq i \leq n$ هستند متغیرهای وابسته باشند و $\epsilon_i \stackrel{iid}{\sim} (0, \sigma^2)$. بررسی اثر متغیرهای مستقل روی متغیرهای وابسته، هدف تحلیل رگرسیون است و رگرسیون را از لحاظ نوع تابعی که رابطه‌ی بین متغیرهای مستقل و وابسته را تعیین می‌کند به دو گروه رگرسیون پارامتری و رگرسیون ناپارامتری دسته‌بندی می‌کنیم.

در رگرسیون پارامتری رابطه بین متغیرهای مستقل و وابسته از طریق تابع پارامتری f به صورت $Y_i = f(x_i) + \epsilon_i$ بیان می‌شود. $f(x_i)$ می‌تواند یک تابع تک متغیره یا چندمتغیره از متغیرهای مستقل و بردار پارامترهای مجهول $\beta' = (\beta_1, \dots, \beta_p)$ باشد.

اما در رگرسیون ناپارامتری رابطه بین متغیرهای مستقل و وابسته از طریق تابع ناپارامتری $m(\cdot)$ به صورت $Y_i = m(x_i) + \epsilon_i$ نشان داده می‌شود که $m(\cdot)$ تابعی نامعلوم است. برای برآورد توابع پارامتری و ناپارامتری روش‌های بسیاری وجود دارد به طوریکه برای برآورد $m(x)$ می‌توان به رگرسیون چندجمله‌ای موضعی، هموارسازی اسپلاین و روش هسته‌ای هموارسازی اشاره کرد [۶].

در این پایان‌نامه با تمرکز بر رگرسیون پارامتری، انواعی از روش‌های برآورد تابع پارامتری بیان می‌شود. به همین منظور ابتدا مروری کوتاه بر برآورد پارامترها در رگرسیون پارامتری

خواهیم داشت .

رگرسیون پارامتری

فرض کنید برای n فرد از افراد یک جامعه آماری مقادیر $(x_i, y_i), \dots, i = 1, \dots, n$ مشاهده شده باشد. مدل $Y_i = f(x_i) + \epsilon_i$ را در نظر می‌گیریم. هدف برآورد تابع $f(x_i)$ می‌باشد که با $\hat{f}(x_i)$ نمایش می‌دهیم. اگر $f(x_i)$ را بصورت $\sum_{j=0}^q \beta_j x_i^j + \epsilon_i$ باشد آنگاه $Y_i = \sum_{j=0}^q \beta_j x_i^j + \epsilon_i$ خواهد بود و برآورد $f(x_i)$ به برآورد پارامترهای مجهول منجر خواهد شد.

روش کمترین مربعات^۱

روشی برای برازش یک مدل به داده‌ها است و از طریق مینیمم کردن

$$Q = \sum_{i=1}^n \left(\sum_{j=0}^q \beta_j x_i^j - y_i \right)^2$$

بدرست می‌آید. با بکارگیری نماد ماتریسی، مدل رگرسیونی بصورت $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ بازنویسی می‌شود.

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^q \\ 1 & x_2 & x_2^2 & \cdots & x_2^q \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^q \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_q \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

بامشتق‌گیری از Q نسبت به β_j ها و به کارگیری نماد ماتریسی برآورد پارامترهای مجهول بصورت

$$\hat{\boldsymbol{\beta}}^{ols} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

و برآورد مدل بصورت $\hat{Y} = \mathbf{X} \hat{\boldsymbol{\beta}}$ خواهد شد.

تعریف ۱-۲. باقی‌مانده (e_i) : اختلاف بین i امین مقدار مشاهده شده (y_i) و i امین مقدار برازش شده (\hat{y}_i) می‌باشد و بصورت ماتریسی $e = \mathbf{Y} - \hat{\mathbf{Y}}$ نوشته می‌شود.

□

^۱ordinary least-square

تعریف ۲-۲. اریبی: فرض کنید $T(\mathbf{X})$ یک برآوردگر برای $\gamma(\theta)$ باشد، میزان اریبی $T(\mathbf{X})$ با b نشان داده می‌شود.

$$b(\theta) = E[T(\mathbf{X}) - \gamma(\theta)]$$

□

تعریف ۲-۳. نارایی: برآوردگر $T(\mathbf{X})$ را برای $\gamma(\theta)$ نارایب گویند اگر و تنها اگر

$$E_{\theta}[T(\mathbf{X})] = \gamma(\theta) \quad b(\theta) = 0 \quad \forall \theta \in \Theta$$

□

تعریف ۲-۴. میانگین مربع خطا^۲: فرض کنید $T(\mathbf{X})$ یک برآوردگر برای $\gamma(\theta)$ باشد در این صورت میانگین مربع خطای برآورد $T(\mathbf{X})$ به صورت زیر تعریف می‌شود

$$MSE_{\theta} = E_{\theta}(T(\mathbf{X}) - \gamma(\theta))^2 = Var(T(\mathbf{X})) + \underbrace{[E(T(\mathbf{X}) - \gamma(\theta))]^2}_b$$

□

میانگین مربع خطا شامل دو مقدار، واریانس $T(\mathbf{X})$ و یک مقدار نامنفی می‌باشد.

تعریف ۲-۵. مقدار پیش‌بینی شده: فرایند تعیین مقدار تابع در یک نقطه جدید مانند x_0 به کمک تابع $\hat{Y}(x_0)$ را مقدار پیش‌بینی شده گویند.

□

تعریف ۲-۶. اختلاف مقدار واقعی و مقدار پیش‌بینی شده در یک نقطه جدید مانند x_0 که با عبارت $\hat{Y} - \hat{Y}(x_0)$ نشان داده می‌شود را خطای پیش‌بینی گویند و با نماد PE ^۳ نشان داده می‌شود.

بررسی مناسب بودن مدل

برای بررسی مناسب بودن مدل رگرسیونی، به مجموعه بزرگی از تکنیک‌ها نیازمندیم که مجموعه این تکنیک‌ها درستی تشخیص^۴ نام دارد. دو مولفه اصلی این تکنیک‌ها،

Mean Square Error^۲

prediction error^۳

diagnostic checking^۴

مقادیر برازش شده و باقی مانده‌ها می‌باشند. بیشترین اطلاعات راجع به مناسب بودن مدل، در باقی مانده‌ها است. زیرا به کشف هر نوع ناسازگاری بین داده‌ها و مدل برازش شده کمک می‌کند. اگر مدل مناسب باشد باقی مانده‌ها الگوی خاصی نخواهند داشت. در صورت وجود یک الگوی خاص در بین باقی مانده‌ها می‌توانیم روی داده‌ها تبدیل مناسبی قرار دهیم. برای اطلاعات بیشتر می‌توان به [۷] مراجعه کرد.

معیارهای سنجش مناسب بودن یک مدل را می‌توان از دو دیدگاه مورد بررسی قرار داد:

۱ - میزان دوری یا نزدیکی مقدار y برازش شده و مقدار واقعی.

۲ - پیش‌بینی مقدار y برای مقادیر جدید توسط \hat{f} .

سنجش دوری یا نزدیکی مقدار y برازش شده و مقدار واقعی از طریق خطای برآورد (e_i) صورت می‌گیرد. اگر $\rho(t)$ هر تابع غیر نزولی در $|t|$ باشد آنگاه $\sum \rho(e_i)$ را می‌توان به عنوان یک معیار برای مقایسه دو مدل بکار برد. از رایج‌ترین توابع غیرنزولی می‌توان تابع درجه دوم و تابع قدرمطلق را نام برد که به ترتیب معیارهای MSE و MAE را تولید می‌کنند.

$$MSE = \frac{1}{n} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$MAE = \frac{1}{n} \sum_{i=1}^N |y_i - \hat{y}_i|$$

برازش‌هایی که دارای خطای برآورد مینیمم باشند برازش مناسب محسوب می‌شوند. خطای پیش‌بینی در یک نقطه خاص را زمانی می‌توان محاسبه کرد که مقدار y آن نقطه را مشاهده کرده باشیم.

ویژگی‌های برآورد کمترین مربعات

۱- $\hat{\beta}^{ols}$ ناریب با حداقل واریانس می‌باشد.

۲- واریانس برآوردگر بصورت زیر است:

$$Var(\hat{\beta}) = \hat{\sigma}^2 (X^T X)^{-1}$$

۳- میانگین توان دوم خطا، (MSE) نیز به روش زیر محاسبه می‌شود:

$$\begin{aligned} MSE(\hat{\beta}^{ols}) &= \hat{\sigma}^2 \text{trace}(X^T X)^{-1} \\ &= \hat{\sigma}^2 \sum_{i=1}^p \frac{1}{d_j} \end{aligned}$$

که در آن مقادیر ویژه ماتریس $(X^T X)$ می باشد.
نمودار پراکنش:

در مطالعه رابطه بین دو متغیر، اولین و ساده ترین قدم رسم داده ها به صورت نقاطی روی یک صفحه نمودار است. شکل حاصله که نمودار پراکنش^۵ نامیده می شود، یک نمودار ریاضی با استفاده از مختصات دکارتی است و با استفاده از این نمودار نحوه و میزان پراکندگی و هم پراشی (پراکندگی توام) متغیرها آشکار می شود. یکی از جنبه های قدرتمند نمودار پراکنش، توانائی نشان دادن رابطه غیر خطی است. همچنین با استفاده از این نمودار داده های پرت شناسائی می شوند.

مثال ۲-۷. نمودار پراکنش داده های Cars : این داده ها در بسته نرم افزاری datasets در نرم افزار R موجود است و مربوط به سرعت و فاصله طی شده از لحظه ترمز کردن تا توقف کامل برای ۵۰ ماشین می باشد که در سال ۱۹۲۰ ثبت شده اند.

scatter plot^۵