

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

۲۲۴/۲



دانشگاه شهید بهشتی

۱۳۷۹ / ۱۱ / ۲۰

دانشگاه شهید بهشتی

دانشکده علوم ریاضی

گروه آمار

پایان نامه کارشناسی ارشد آمار

عنوان:

نیرومندی در رگرسیون پیری

استاد راهنما:

دکتر سیامک نوربلوچی

۱۹۸۲۴

نگارش:

علی بابایی

شهریور ۱۳۷۸

۳۲۴۱۳

بیتالی

تاریخ

ردیف

پوست

صور تجلسه دفاع از پایان نامه

جلسه هیئت داوران ارزیابی پایان نامه آقای / خانم / علی بابایی

به شناسنامه شماره ۲۰ صادره از اراک متولد ۱۳۵۲

دانشجوی دوره کارشناسی ارشد ناپیوسته رشته آمار گرایش محض

با عنوان نیرومندی بیزی در رگرسیون بیزی

به راهنمایی آقای دکتر نوریلوچی طبق دعوت قبلی در تاریخ ۷۸/۶/۳۱

تشکیل گردید و براساس رأی هیأت داوران و با عنایت به ماده ۲۰ آئین نامه

کارشناسی ارشد مورخ ۷۳/۱۰/۲۵ پایان نامه مزبور بانمره ۱۸ تمام (هیچجه تمام)

و درجه عالی مورد تصویب قرار گرفت.

۱- آقای دکتر سیامک نوریلوچی

۲- " " عبدالرحیم شهبایی

۳- " " علی عمیدی

۴- " " محمد نکایی

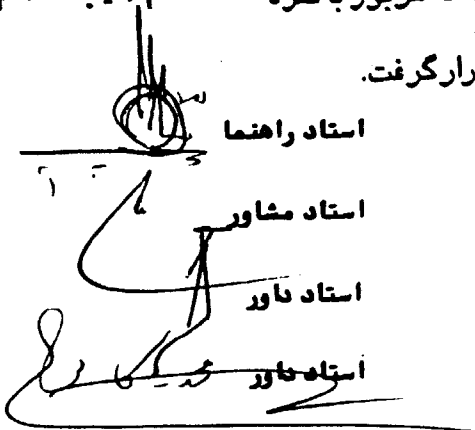
۵-

استاد راهنما

استاد مشاور

استاد داور

استاد داور



تقدیریم:

به طاهریم که در لفظه لفظه دوران تحصیلم چشمهای
نگرانش بدرقه راهم بود، او که همواره امید را در دلم
زنده نگه میداشت و فردایی روشن را نوید میداد.

و پیریم که همواره چون کوهی استوار تکیه گاهم بودند.

و برادریم که همه تحصیلاتم را مدیون محبت های ایشان
هستم.

و خواهرانم که همواره غمخوارم بودند.

سپاس و قدردانی

دکتر شریعتی در وصیت به پسرش می‌نویسد: پسر من تنها چیزی که برای تو در این مسیری که عمر نام دارد آرزو می‌کنم، برخورد با یکی دو دل بزرگ، یکی دو روح عظیم است، یک بالاترین عدد ممکن است، دو را برای وزن کلام آورده‌ام و من بقدر خوشبختی که به آن یک، برخوردم، آن هم در بهترین ایام زندگی، جوانی. آری آن یک، کسی نیست جز استادی بزرگوار، انسانی متعهد نیک سرشت و دوست داشتنی، جناب آقای دکتر نوریلوپی که در تمام ایام که در خدمت ایشان بودم نه تنها در زمینه علمی از مضر ایشان کسب فیض نمودم بلکه به مثابه شاگردی که تازه پا در کلاس درس می‌گذارد قدم در کلاس درس اخلاق ایشان گذاشتم و بهره‌های فراوان از حسن خلق و حس انسان دوستی ایشان بردم. امیدوارم این عزیز همواره پاینده باشند تا فرزندان این مرز و بوم از وجود ایشان بهره‌مند شوند.

جا دارد از کلیه اساتید گروه آمار خصوصاً آقای دکتر عبدالرحیم شهلایی بفاخر قبول زحمت مشاورت و از آقای دکتر ذکایی که مرحمت نموده و داوری پایان نامه اینجانب را پذیرفته‌اند تشکر و قدردانی نمایم.

همچنین آقای دکتر علی عمیدی و آقای دکتر محمد رضا مشکانی که براسستی تبسم عینی معلم بودن هستند چه بفاخر زحماتی که در مطالعه و داوری پایان نامه و چه در دوران تصحیح اینجانب متحمل شده‌اند صمیمانه تشکر و قدردانی می‌نمایم.

بر خود لازم میدانم که از کمکهای بی‌شائبه دوست عزیزم آقای محمدرضا اناری صمیمانه تشکر کنم. همچنین از مساعدت دوستان عزیزم در دوران تصحیح خصوصاً آقایان قدرت روشنایی، حبیب اسماعیلی، مهدی امانی، سعدا... مرادی و مهدی ابراهیم زاده قدردانی می‌نمایم.

پیشگفتار

با پیدایش مکتب بیزین‌ها در علم آمار، همه شاخه‌های آماری تحت تأثیر این مکاتب قرار گرفت در این رهگذر مساله رگرسیون نیز از این تأثیر مصون نماند و مساله رگرسیون بیزی مطرح شد.

عموماً در رگرسیون بیزی از دو نوع توزیع پیشین برای پارامترهای مدل استفاده شده است: پیشین جفریز و پیشین نرمال - گاما یا نرمال - گامای وارون. در مساله رگرسیون بیزی به خاطر تعداد زیاد پارامترها انتخاب پیشین مشکل است، این امر انگیزه‌ای برای ورود به بحث نیرومندی بیزی در رگرسیون شده است. در این رساله چند صورت از توزیعهای پیشین که نتایج رگرسیون بیزی نسبت به آنها نیرمند است را مورد مطالعه قرار داده‌ایم. ساختار این رساله به شرح زیر است:

در فصل اول با مروری کوتاه بر رگرسیون کلاسیک وارد رگرسیون بیزی می‌شویم ابتدا مساله رگرسیون بیزی یک متغیره با پیشین ناآگاهنده (جفریز) را مورد بررسی قرار داده و سپس مساله را در حالت رگرسیون چند گانه با پیشین مذکور مطرح می‌کنیم. در ادامه فصل با در نظر گرفتن پیشین نرمال - گاما برای پارامترهای مدل تحلیل‌های مربوطه را ارائه کرده‌ایم.

در فصل دوم برای کامل بودن بحث، نیرومندی بیزی را بطور مستقل مطرح کرده‌ایم بدین صورت که ابتدا مفهوم نیرومندی بیزی را بیان کرد، و سپس نتایج نیرومندی بیزی را در رده‌های خاص مثل رده ϵ - آلوده‌ها با انواع آلاینده‌ها بررسی کرده‌ایم.

در فصل سوم مساله نیرومندی در رگرسیون بیزی در رده بخصوص از پیشین‌ها را مطرح کرده‌ایم. ابتدا رگرسیون بیزی را با پیشین آمیخته نرمال - گامای

واردن و نا آگاهنده عنوان کرده و نیرومندی نتایج مدل نسبت به انتخاب پیشین پایه را نشان داده‌ایم. در قسمت بعد مساله را با انتخاب رده ϵ - پیشین توزیع پایه و آلاینده یک نوع خاص از g - پیشین هاست پی گیری کرده و نتایج نیرومندی را ارائه داده‌ایم.

امید است که این رساله بتواند مبحث نیرومندی بیزی در رگرسیون بیزی را به

جامعه آماری کشور معرفی کرده باشد.

فهرست

صفحه	عنوان
	فصل اول: رگرسیون خطی بیزی پارامتری
۱	مقدمه:
۲	۱-۱: رگرسیون کلاسیک ساده و چندگانه
۲	۱-۱-۱: رگرسیون کلاسیک ساده
۵	۲-۱-۱: رگرسیون کلاسیک چندگانه
۸	۲-۱: رگرسیون بیزی یک متغیره با پیشین ناآگاهنده
۱۵	۳-۱: رگرسیون بیزی چندگانه با پیشین ناآگاهنده
۲۰	۱-۳-۱: به هنگام سازی پسین
۲۲	۲-۳-۱: تابع چگالی پیش گو
۲۵	۳-۳-۱: تحلیل مدل وقتی $X X$ ویزه است
۲۶	۴-۱: رگرسیون بیزی چندگانه با پیشین مزدوج نرمال - گاما
۲۷	۱-۴-۱: تحلیل پسین
۳۰	۲-۴-۱: استنباط برای مدل رگرسیونی
۳۱	۱-۲-۴-۱: برآورد نقطه ای
۳۲	۲-۲-۴-۱: برآورد ناحیه ای و فاصله ای
۳۴	۳-۲-۴-۱: آزمون فرضیه ها
۳۶	۳-۴-۱: تحلیل پیش گو

صفحه

عنوان

۳۹	۱-۴-۴: g- پیشین ها.....
۴۳	۱-۴-۵: برآورد پارامترهای توزیع پیشین نرمال - گاما به روش بیز تجربی.....
۴۷	فصل دوم: آشنایی با نیرومندی بیزی.....
۴۷	۲-۱: مقدمه.....
۴۷	۲-۲: دیدگاه نیرومندی بیزی.....
۵۶	۲-۳: مروری بر معیارهای نیرومندی بیزی.....
۶۰	۲-۴: پسین نیرومندی نسبت به رده توزیعیهای ϵ - آلوده.....
۶۰	۲-۴-۱: رده ϵ - آلودهها: تمام توزیعیها.....
۶۲	۲-۴-۲: رده ϵ - آلودهها: توزیعیهای تک مدی.....
۶۳	۲-۴-۳: رده ϵ - آلودهها: توزیعیهای تک مدی و متقارن.....
۶۴	۲-۵: پیشین‌های نیرومند.....
۶۶	فصل سوم: نیرومندی در رگرسیون بیزی پارامتری.....
۶۶	۳-۱: مقدمه.....
	۳-۲: بررسی نیرومندی بیزی استنباطهای مدل رگرسیون خطی تحت پیشین آمیخته
۶۷	نرمال - گامای وارون و ناآگاهنده.....
۷۰	۳-۲-۱: تحلیل پسین - پیشین.....

صفحه	عنوان
۷۷	۲-۲-۳: بررسی نیرومندی پسین در ردهٔ پیشین‌های نرمال - گامای وارون
	۳-۳: بررسی نیرومندی بیزی در ردهٔ ۴- آلوره، وقتی توزیع پایه و آلاینده، g - پیشین باشد.
۸۱	۱-۳-۳: انتخاب رده‌ای مناسب برای توزیع پیشین پارامترها
۸۲	۲-۳-۳: تحلیل پسین
۹۱	۳-۳-۳: پیش بینی نقطه‌ای
۹۶	۴-۳: تحلیل حساسیت در مدل رگرسیون بیزی
۱۰۲	فصل چهارم: کاربرد

فصل اول

رگرسیون خطی بیزی پارامتری

مقدمه:

رگرسیون کلاسیک مبتنی بر روش کمترین توانهای دوم و نرمال بودن توزیع خطاها گسترش یافته است. نظریه بیزی در برخورد با مسأله رگرسیون نیز حداقل سابقه ای ۳۰ ساله دارد. مسأله نیرومندی بیزی در رگرسیون بیزی موضوع اصلی این رساله است. در فصل حاضر با رگرسیون کلاسیک و همچنین نظریه بیزی آن آشنا می‌شویم. مراجع اصلی ما در این فصل زلنر^(۱) (۱۹۷۱)، بروملینگ^(۲) (۱۹۸۵) و پرس^(۳) (۱۹۸۹) است.

این فصل دارای ساختار زیر است. بخش اول را با توجه به دانش عامی که درباره رگرسیون کلاسیک موجود است صرفاً جهت کامل بودن بحث و معرفی نمادها معرفی کردیم و در بخش دوم با توجه به نمادهای بخش قبل، رگرسیون بیزی پارامتری را مطرح نمودیم. به صورت دقیقتر این بخش‌ها عبارتند از:

I: رگرسیون کلاسیک

(۱) رگرسیون یک متغیره کلاسیک

(۲) رگرسیون چندگانه کلاسیک

II: رگرسیون بیزی

الف: وقتی σ^2 معلوم است.

1- Zellner

2- Bromelling

3- Press

ب: وقتی σ^2 مجهول است.

(۱) رگرسیون ساده بیزی (با پیشین نا آگاهنده)

(۲) رگرسیون چندگانه بیزی پارامتری (ناتجربی) }
 با پیشین نا آگاهنده
 با پیشین مزدوج نرمال - گاما

در بخش اخیر مباحث مربوطه به روشهای کلاسیک استنباطی یعنی مسأله‌ی برآورد کردن و آزمون فرضیه‌ها و همچنین مسأله پیشگوئی مقادیر آینده را مطرح می‌کنیم و نتایج بدست آمده در این حوزه‌ها را مرور می‌نمائیم. قضایا و روشهای مطرح در این فصل ابزاری ضروری برای فصلهای آینده این رساله است.

۱-۱- رگرسیون کلاسیک ساده و چندگانه

۱-۱-۱- رگرسیون کلاسیک ساده

در مدل رگرسیون خطی ساده یک متغیر وابسته و یک متغیر مستقل داریم و مدل رگرسیونی به صورت زیر مطرح است.

$$Y_i = \beta_1 + \beta_2 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1-1)$$

Y_i مقدار متغیر وابسته از i -امین مشاهده و x_i مقدار i -امین متغیر مستقل است. فرض می‌شود که متغیر X تحت کنترل آزمایشگر است. ε_i ها را خطای تصادفی یا مانده مدل رگرسیونی می‌نامند.

فرضهای مدل رگرسیونی ساده

۱- متغیر تصادفی ε_i دارای توزیعی با میانگین صفر و واریانس ثابت σ^2 است.

۲- ε_i ها دویبدو ناهمبسته‌اند: $\forall i \neq j \text{ Cov}(\varepsilon_i, \varepsilon_j) = 0$

β_1 را عرض از مبدأ و β_2 را شیب خط رگرسیونی می‌نامند. اگر $x=0$ در حوزه مدل باشد

β_1 مقدار Y به ازای $x=0$ است و β_2 نتغییر Y به ازای یک واحد تغییر در x است.

از مقدمات فوق توزیع y به شرط x عبارتست از:

$$Y_i | x_i \sim N(\beta_1 + \beta_2 X_i, \sigma^2) \quad , \quad i = 1, \dots, n$$

که در آن داریم: $\text{Cov}(Y_i | x_i, Y_j | x_j) = 0 \quad \forall i \neq j$

برآورد پارامترهای مدل به روش کمترین توانهای دوم خطا^(۱) (LSE)

در این روش ایده اصلی این است که خط رگرسیونی را طوری برآورد کنیم که

فاصله عمودی مشاهدات از خط بدست آمده کمترین باشد. معیاری که بدین منظور

تعریف می‌کنیم به صورت زیر است:

$$Q(\beta_1, \beta_2) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_1 - \beta_2 x_i)^2$$

Q تابعی از β_1 و β_2 است. مقادیر $\hat{\beta}_1$ و $\hat{\beta}_2$ رایج در رگرسیون یعنی:

$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x} \quad (2-1)$$

$$\hat{\beta}_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

با مشتق گیری از $Q(\beta_1, \beta_2)$ نسبت به β_1 و β_2 بدست آمده است.

برآورد پارامترها به روش ماکسیم درستنمایی⁽¹⁾ (ML):

اگر فرض نرمال بوده ε_i ها را بپذیریم، لگاریتم تابع درستنمایی عبارتست از:

$$\log L(\beta_1, \beta_2, \sigma^2) = -\frac{n}{2} (\log 2\pi + \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_1 - \beta_2 X_i)^2)$$

که چنانچه آن را نسبت به β_1 و β_2 و σ^2 ماکسیم کنیم. خواهیم داشت:

$$\hat{\sigma}^2 = \text{MSE} = \frac{1}{n} \sum (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_i)^2 = \frac{\text{SSE}}{n}$$

(۳-۱)

$$\hat{\beta}_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}, \quad \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}$$

که همان برآوردهای کمترین توانهای دوم قبلی است.

به علاوه می توان نشان داد که اگر σ^2 معلوم باشد داریم:

$$\hat{\beta}_2 \sim N\left(\beta_2, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$$

(۵-۱)

$$\hat{\beta}_1 \sim N\left(\beta_1, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2}\right)\right)$$

از نتایج فوق می توان برای انجام آزمون فرضیه و به دست آوردن فواصل اطمینان

برای β_1 و β_2 استفاده کرد.

۱-۱-۲- رگرسیون کلاسیک چندگانه

در این مدل تغییرات متغیر وابسته Y به چندین متغیر مستقل مربوط می‌شود. فرض کنید Y متغیر وابسته و X_1, \dots, X_{k-1} متغیرهای مستقل باشند. مدل رگرسیونی به صورت زیر است.

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_{k-1} X_{i,k-1} + \varepsilon_i, \quad i = 1, \dots, n \quad (6-1)$$

که برای سهولت آن را به دو صورت ماتریسی زیر می‌نویسیم.

$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon} \quad (7-1)$$

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \underline{X}_n = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1,k-1} \\ 1 & X_{21} & X_{22} & \dots & X_{2,k-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{n,k-1} \end{bmatrix}, \quad \underline{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \quad \underline{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{k-1} \end{bmatrix}$$

فرض می‌کنیم:

$$\underline{\varepsilon} \sim N(0, \sigma^2 I_n)$$

$$E(\underline{y}) = \underline{X}\underline{\beta}, \quad \text{Var}(\underline{y}) = \sigma^2 I_n, \quad \underline{\varepsilon} = \underline{y} - \underline{X}\underline{\beta} \quad \text{در نتیجه داریم:}$$

$$Q(\underline{\beta}) = \underline{\varepsilon}'\underline{\varepsilon} = (\underline{y} - \underline{X}\underline{\beta})'(\underline{y} - \underline{X}\underline{\beta}) \quad \text{گیریم:}$$

که از حداقل کردن آن داریم:

$$\hat{\underline{\beta}} = (\underline{X}'\underline{X})^{-1} \underline{X}'\underline{y} \quad (8-1)$$

$$\text{SSE} = \sum_{i=1}^n e_i^2 = (\underline{y} - \underline{X}\hat{\underline{\beta}})'(\underline{y} - \underline{X}\hat{\underline{\beta}}) = \underline{y}'\underline{y} - \hat{\underline{\beta}}'\underline{X}'\underline{y} \quad (9-1)$$

با فرض نرمال بودن توزیع $\underline{\varepsilon}$ خواهیم داشت:

$$\hat{\underline{\beta}} \sim N(\underline{\beta}, \sigma^2 (\underline{X}'\underline{X})^{-1}) \quad (10-1)$$