



98947

دانشگاه پیام نور

مرکز تهران

دانشکده علوم پایه گروه: آمار

پایان نامه:

برای دریافت درجه کارشناسی ارشد

گرایش: آمار ریاضی

عنوان:

عنوان:

تحلیل رگرسیون معکوس قطعه ای استوار

استاد راهنما:

دکتر مسعود یارمحمدی

استاد مشاور:

دکتر مجتبی گنجعلی

نگارش:

هادی عبدالحسینی

خرداد ۱۳۸۲

۱۳۸۷ / ۲ / ۱۴

۱۳۸۷ / ۴ / ۱۳



۹۵۹۲۷

«قدر دانی»

با سپاس ایزد دانا و توانا در اینجا بر خود لازم می بینم که از زحمات و همکاری و همفکری اساتید محترم و گرانمایه که طی مراحل مختلف این پایان نامه مرا یاری نموده اند کمال تشکر و قدر دانی خود را ابراز نمایم.

استاد محترم آقای دکتر مسعود یارمحمدی که راهنمایی این پایان نامه را بر عهده داشته و از هیچ کمکی در این مورد دریغ نفرمودند از زحمات و همکاری ایشان نهایت تشکر را می نمایم.

از استاد محترم آقای دکتر مجتبی گنجعلی که در سمت استاد مشاور قبول زحمت فرموده و در انجام این پایان نامه اینجانب را یاری فرمودند نهایت تشکر و سپاسگزاری را دارم.

از اساتید محترم آقای دکتر پرویز نصیری و آقای دکتر علی اکبر رحیم زاده ثانی که داوری این پایان نامه را به عهده گرفته اند کمال تشکر را دارم.

در پایان هم از تمامی دوستانی که مرا در تهیه منابع و مراجع این پایان نامه یاری رسانده اند بویژه آقای رضا ندیمی کمال تشکر و امتنان را دارم.

تقديم:

به آنان که پیمانۀ «لا» زدند.....

فهرست مندرجات

فصل اول: کلیات و مسائل مربوط

- ۱-۱ مقدمه..... ۱
- ۲-۱ تاریخچه ، ضرورت و کاربردهای روش..... ۲
- ۳-۱ داده های دور افتاده در مسائل رگرسیونی..... ۳
- ۴-۱ هم وردایی..... ۵
- ۵-۱ نقطه فروریزش..... ۵
- ۶-۱ نقطه فروریزش نمونه متناهی..... ۷
- ۷-۱ روابط کمی بین **ABP** و **RBP**..... ۸
- ۸-۱ توزیع بیضوی..... ۱۱
- ۹-۱ فضای پوچ..... ۱۲
- ۱۰-۱ تحلیل همبستگی کانونی..... ۱۲
- ۱۱-۱ تحلیل مولفه اصلی..... ۱۴

فصل دوم: رگرسیون معکوس قطعه ای

- ۱-۲ مقدمه..... ۱۹
- ۲-۲ چرا کاهش بعد؟..... ۲۰
- ۳-۲ فضاهای کاهش داده شده معمولی و روش های کاهش بعد..... ۲۱
- ۱-۳-۲ فضای خطی وزن دار شده..... ۲۱
- ۲-۳-۲ فضای خطی بعد از تبدیل وابسته..... ۲۲
- ۳-۳-۲ تابع جمع پذیر بعد از تبدیل مستقل..... ۲۲
- ۴-۳-۲ فضای جمع پذیر بعد از تبدیل مستقل و وابسته..... ۲۳

۲۳.....	۵-۳-۲ روشهای کاهش بعد متقارن
۲۴.....	۴-۲ فضای کاهش بعد مؤثر (<i>edr</i>)
۲۵.....	۵-۲ روشهای برآورد <i>edr</i>
۲۵.....	۱-۵-۲ رگرسیون خطی
۲۵.....	۲-۵-۲ روش SIR - I
۲۹.....	۳-۵-۲ روش SIR - II
۳۰.....	۴-۵-۲ اجرای روش SIR - I
۳۲.....	۵-۵-۲ اجرای روش SIR - II
۳۶.....	۶-۵-۲ روش واریانس قطعه ای (SAVE)
۳۷.....	۷-۵-۲ روش راستای هسه ای اصلی (PHD)
۳۸.....	۸-۵-۲ برآورد دنباله ای
۳۸.....	۶-۲ فرضهای اساسی
۴۰.....	۷-۲ خواص نمونه ای SIR
۴۲.....	۸-۲ ارزیابی تعداد مؤلفه های مدل
۴۳.....	۹-۲ مثال عملی

فصل سوم: روشهای استوار با نقطه فروریزش بالا

۴۶.....	۱-۳ مقدمه
۴۶.....	۲-۳ M برآوردگرها
۴۸.....	۳-۳ S- برآوردگرها
۵۰.....	۴-۳ خواص S- برآوردگرهای استوار چند متغیره مکانی و کوواریانسی
۵۱.....	۵-۳ فروریزش و رده دور افتاده ها
۵۶.....	۶-۳ فرآیندهای دوبار اصلاح شده
۵۶.....	۳-۶-۱ تابع t- دووزنی (دووزنی انتقال داده شده)
۵۸.....	۳-۶-۲ تابع مسطح مضاعف

۷-۳ برآورد تعقیب تصویر..... ۶۱

فصل چهارم: رگرسیون معکوس قطعه ای استوار

۱-۴ مقدمه ۶۴

۲-۴ حساسیت SIR به داده های دورافتاده..... ۶۵

۳-۴ SIR ی تعمیم یافته..... ۶۶

۴-۴ هم بردای آفین GSIR..... ۶۷

۵-۴ نقطه فروریزش نمونه متناهی GSIR..... ۶۸

۶-۴ یک روش خاص از GSIR..... ۶۸

پیوست الف ۷۴

چکیده

روش رگرسیون معکوس قطعه ای، روشی است که برای کاهش بعد، در مدل‌های رگرسیونی استفاده می‌شود. عمل کاهش بعد براساس اطلاعاتی که در منحنی رگرسیون معکوس قرار می‌گیرد انجام می‌شود. در مراحل مختلف روش رگرسیون معکوس قطعه ای از برآوردهای کلاسیک استفاده می‌شود. بنابراین نتایج روش در مقابل داده های دور افتاده ناستوار است. یک روش برای استوار کردن این روش جایگزینی برآوردهای کلاسیک با برآوردهای استوار است. در این رساله روش رگرسیون معکوس قطعه ای به عنوان یک روش کاهش بعد معرفی شده و سپس با ارائه یک مثال کاربردی نحوه استفاده از این روش مورد بررسی قرار می‌گیرد. در پایان به استوار سازی این روش با استفاده از برآوردهای استوار می‌پردازیم.

واژه های کلیدی: کاهش بعد، رگرسیون معکوس قطعه ای، دور افتاده ها، نقطه

فروریزش، استواری

کلیات و مسائل مربوط

۱-۱ مقدمه

با توجه به اینکه یک مسئله ویژه در تحلیل داده ها، فزونی بیش از حد اطلاعات^۱ است و این فزونی بیش از حد اطلاعات از تعداد زیاد متغیرهای مورد استفاده برای مشاهدات در مجموعه داده ها می آید، و با توجه به اصل امساک که ما را از استفاده بیش از حد متغیرها منع می کند، روشهای گوناگونی در رابطه با کاهش تعداد متغیرها و به نوعی در کاهش بعد آمده است که به ما کمک می کند با استفاده از کمترین تعداد متغیر به حداکثر اطلاعات دسترسی پیدا کنیم. این مسئله فزونی بیش از حد اطلاعات اغلب در مدل‌های وابسته - مستقل نمایان می شود. یعنی زمانی که فرد می خواهد وابستگی متغیر پاسخ Y را با یک یا بیشتر متغیر پیشگوی X_1, X_2, \dots, X_m بدست آورد.

^۱ Overabundance of Information

برای $m > 2$ این مسئله اغلب بدون استفاده از برخی صورتهای کاهش بعد غیرممکن می شود. در تحلیل داده ها، این مسئله اغلب بوسیله ساختن تصاویر پایین بعدی از داده ها حل می شود. این تصاویر ابعاد پایین تر، درمدهایی استفاده می شوند که متغیر پاسخ، تابعی از این تصاویر پایین بعدی از فضای متغیرهای کمکی است. این امیدواری وجود دارد که رابطه بین پاسخ و متغیرهای کمکی در فضای پایین بعدی حفظ شود. معمولاً تابعی که آماردانه انتخاب می کنند و یا ترجیح می دهند انتخاب کنند، تابعهای خطی است که آن هم به علت سادگی و تفسیر پذیری آنهاست. با این حال برخی مواقع یک ساختار صحیح از رابطه بین پاسخ و متغیرهای کمکی نمی تواند به طور مناسبی بوسیله یک تابع خطی حفظ شود. روشهای متفاوت بسیاری تهیه شده اند که این امیدواری را می دهند که اجازه انعطاف بیشتری را برای به مدل درآوردن این رابطه، با حفظ سادگی و تفسیر پذیری ایجاد کنند، که از جمله این روشها می توان به رگرسیون پیگیری تصویر^۲ اشاره کرد. از روشهای اخیر هم می توان به روش رگرسیون معکوس قطعه ای (SIR)^۳ اشاره کرد.

۲-۱ تاریخچه ، ضرورت و کاربردهای روش

روش رگرسیون معکوس قطعه ای و کاربردهای آن نخستین بار توسط لی^۴ در سال ۱۹۹۱ مطرح شده است. روش رگرسیون معکوس قطعه ای روشی سودمند برای کاهش بعد در مسئله رگرسیون ناپارامتری است. این روش هنگامی استفاده می شود که متغیر پاسخ Y بستگی به k ترکیب خطی نامشخص از متغیرهای توضیحی $X = (X_1, \dots, X_m)$ داشته و صورت دقیق و کامل وابستگی معلوم نباشد. لازم به ذکر است که روش رگرسیون معکوس قطعه ای برای کاهش بعد متغیرهای توضیحی، بدون انجام یک فرایند برآزش مدل پارامتری یا ناپارامتری به کار می رود.

^۱ Projection Pursuit Regression

^۲ Sliced Inverse Regression

^۴ Li

در مورد کاربردهای این روش می توان به مباحث کاربردی آمار مانند پزشکی ، ژنتیک و زیست شناسی که در آنها عموماً با داده های با ابعاد بالا که روشهای معمولی رگرسیونی معمولاً قادر به انجام آن نبوده و تحلیل داده ها در ابعاد بالا آسان نیست، اشاره کرد. با استفاده از روش رگرسیون معکوس قطعه ای می توانیم ضمن کاهش بعد داده ها، به طور چشم گیری در تحلیل این داده ها به نتایج دقیق تری دست یابیم.

۱-۳ داده های دورافتاده در مسائل رگرسیونی

یک تعریف کلی برای مشاهدات دورافتاده به اینصورت بیان می گردد که مشاهده دورافتاده مشاهده ای است که مانده آن از نظر قدرمطلق خیلی بزرگتر از مانده های سایر مشاهدات است. البته گاهی مشاهدات دورافتاده به گونه ای دیگر خود را بروز می دهند در برخورد با داده ها برای شناسایی نقاط دورافتاده، معیارهایی وجود دارند که فرد با استفاده از آنها می تواند به دسته بندی داده ها اقدام کند. بدین معنی که برخی نقاطی که از اکثریت داده ها متفاوت هستند و به نوعی می توان آنها را ((خاص)) نامید شناسایی کرد. یکی از معروف ترین این معیارها مربوط به نقاط با نفوذ می باشد. این نقاط مهم هستند چرا که در برازش مدلها و در نتیجه در مقادیر پیش بینی شده تأثیر می گذارند. نکته ای که در اینجا بایستی به آن توجه داشت این است که لزوماً همه نقاط با نفوذ بالا داده دورافتاده نیستند و بدین مفهوم نیست که آنها را باید از داده ها جدا کرده و یا حذف نمود. چه بسا که بعضی از این نقاط حاوی اطلاعات مفیدی در مورد جامعه مورد بررسی باشند. در رگرسیون خطی دور افتادگی در چندین روش تعریف می شود.

برای مدل رگرسیونی در حالت چند متغیره که به صورت زیر بیان می شود:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m + \varepsilon \quad (1.3.1)$$

که X_1, X_2, \dots, X_m متغیرهای توضیحی و $\beta_0, \beta_1, \dots, \beta_m$ ضرایب رگرسیونی در حالت

چند متغیره هستند. نفوذ i امین نقطه در فضای چند متغیره به صورت زیر تعریف می شود:

$$h_{ii} = x_i'(XX)^{-1}x_i \quad (۲.۳.۱)$$

که X ماتریس متغیر کمکی $n \times k$ است. نفوذ یک نقطه به طور ویژه بیان می کند که چقدر یک نقطه خاص x_i از مرکز توده متغیرهای مستقل، یعنی از $\bar{x} = \frac{1}{n}I'X$ دور است (که I برداری n بعدی با مؤلفه های واحد است). در مورد توده های تک مدی از نقاط متغیر های کمکی این می تواند یک اندازه مناسب دوری یک نقطه از مرکز اکثریت نقاط باشد.

مسئله ای که در اینجا مورد توجه است این است که داده هایی که در نهایت داده دورافتاده هستند می توانند اثرات بدی روی برآوردگرها و خواص آنها داشته باشند. بنابراین معرفی برآوردگرهایی که بتوانند در مقابل داده های دورافتاده مقاومت نشان داده و کمتر تحت تأثیر این داده ها واقع شوند می تواند بسیار حائز اهمیت باشد. این برآوردگرها در قالب برآوردگرهای استوار^۵ معرفی می شوند.

بنابراین پس از ارائه یک برآوردگر، مسئله اساسی رفتار آن برآوردگر در رابطه با داده های دورافتاده و بررسی اثرات مخرب آنهاست. در مرحله بعد به منظور مقاوم سازی برآوردگر در مقابل داده های دورافتاده استفاده از روشهای استوار سازی لازم می باشد.

حال با توجه به اینکه جهت بدست آوردن برآوردگر SIR از معیارهای میانگین و تابع اتو کوواریانس استفاده می شود که به داده های دور افتاده مقاوم نمی باشند، لذا با استفاده از روشهای استوار سازی، به اصلاح این برآوردگرها پرداخته و سعی می کنیم آنها را در برابر داده های دور افتاده مقاوم تر نماییم.

در ادامه این فصل با ارائه برخی نمادها و تعریف هایی که در فصول آینده به آنها نیاز داریم، که از جمله آنها تعاریف و مباحث مربوط به استواری برآوردگرها را می توان نام برد، می پردازیم.

^۵ Robust estimators

۴-۱ هم وردایی^۶

مفهوم هم وردایی یک برآوردگر بصورت زیر تعریف می شود:

برآوردگر T هم ورداست اگر هم وردای رگرسیونی^۷، هم وردای مقیاسی^۸ و هم وردای وابسته^۹ باشد.

برآوردگر T را هم وردای رگرسیونی گوئیم هر گاه :

$$T(\{(x_i, y_i + x_i v); i=1, \dots, n\}) = T(\{(x_i, y_i); i=1, \dots, n\}) + v \quad (۳. ۴. ۱)$$

برآوردگر T را هم وردای مقیاسی گوئیم هر گاه :

$$T(\{(c x_i, y_i + x_i v); i=1, \dots, n\}) = c T(\{(x_i, y_i); i=1, \dots, n\}) \quad (۴. ۴. ۱)$$

برآوردگر T را هم وردای وابسته گوئیم هر گاه :

$$T(\{(x_i A, y_i + x_i v); i=1, \dots, n\}) = A^{-1} T(\{(x_i, y_i); i=1, \dots, n\}) \quad (۵. ۴. ۱)$$

که در آن A هر ماتریس ناویژه مربع است.

باید توجه داشت برای این که یک برآوردگر هم وردا باشد باید در هر سه تعریف هم وردایی مطرح شده در فوق صدق کند و صرف هم وردای رگرسیونی یا هم وردای مقیاسی یا هم وردای وابسته بودن و یا دو مورد از موارد فوق را دار بودن نمی تواند به هم وردا بودن برآوردگر منجر شود.

۵-۱ نقطه فروریزش^{۱۰} :

فرض کنید یک نمونه شامل n مشاهده :

$$z = \{(x_{11}, \dots, x_{1m}, y_1), \dots, (x_{n1}, \dots, x_{nm}, y_n)\}$$

^۶ Equivariant
^۷ Regression equivariant
^۸ Scale equivariant
^۹ Affine equivariant
^{۱۰} Break down point

و T برآوردگر رگرسیونی برای پارامتر θ باشد در اینصورت با به کار بردن T و نمونه z ،

$$\text{حاصل یک بردار از ضرایب رگرسیونی } \hat{\theta} = \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_m \end{pmatrix} \text{ می باشد به طوری که } T(z) = \hat{\theta} \text{ . حال اگر}$$

تمام نمونه های مخدوش شده z' که بوسیله جایگذاری m مقدار دلخواه به جای مقادیر اصلی بدست آمده را در نظر بگیریم، بیشترین اریبی که می تواند توسط ورود یک ناخالصی ایجاد شود عبارت است از :

$$\text{bias}(m, T, z) = \sup_{z'} \|T(z') - T(z)\| \quad (6.5.1)$$

اگر $\text{bias}(m, T, z)$ نامتناهی باشد، بدین معنی است که m نقطه دورافتاده روی برآوردگر T اثر زیاد داشته و در نتیجه برآوردهای حاصل نامناسب خواهند بود.

برای یک نمونه متناهی z ، نقطه فروریزش برآوردگر T توسط رابطه زیر تعریف می شود :

$$\varepsilon_v^* = \left\{ \min \frac{m}{n} \quad \text{s.t.} \quad \text{bias}(m; T, z) = \infty \right\} \quad (7.5.1)$$

یعنی $\min \frac{m}{n}$ به طوری که $\text{bias}(m; T, z)$ نامتناهی باشد.

به عبارت دیگر این کوچکترین کسری از ناخالصی داده هاست که سبب می شود برآورد T مقداری دورتر از مقدار مورد انتظار اختیار کند. در روش کمترین توان های دوم مشاهده می شود که تنها یک نقطه دورافتاده کافی است که این برآوردگر مغشوش شود. بنابراین نقطه فروریزش برآوردگر روش توان های دوم $\frac{1}{n}$ است. که اگر حجم نمونه افزایش یابد عبارت فوق به سمت صفر میل می کند. بنابراین می توان گفت که روش توان های دوم دارای نقطه فروریزش صفر درصد می باشد.

تعریف دیگری از نقطه فروریزش به صورت زیر بیان می گردد (کسلّا و برگر^{۱۱}، ۱۹۹۰، صفحه ۲۴۲):

اگر $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ یک نمونه مرتب شده با حجم n باشد و T_n یک آماره براساس این نمونه باشد، آنگاه T_n دارای نقطه فروریزش b است که $0 \leq b \leq 1$ ، هر گاه داشته باشیم:

$$\forall \varepsilon > 0 \quad \lim_{x_{(1-b)n} \rightarrow \infty} T_n < \infty \quad \lim_{x_{(1-b+\varepsilon)n} \rightarrow \infty} T_n = \infty \quad (۸.۵.۱)$$

با توجه به روابط بالا می توان نشان داد که نقطه فروریزش میانگین صفر درصد و نقطه فروریزش میانه پنجاه درصد است.

نقاط فروریزش نقش مهمی در برآوردهای استوار دارند. با توجه به تعریف نقطه فروریزش هر چه یک برآورد گر دارای نقطه فروریزش بالایی باشد آن برآوردگر در برابر نقاط دورافتاده کارایی بالایی دارد. بنا بر این جهت کاهش اثر نقاط دورافتاده بر روی برآورد پارامترها، یافتن برآوردهای با نقاط فروریزش بالا ضروری می باشد.

۱-۶ نقطه فروریزش نمونه متناهی^{۱۲}

این نماد نخستین بار در سال ۱۹۶۷ و بوسیله هاج^{۱۳} منتشر شد، و سپس بوسیله همپل^{۱۴} (۱۹۶۸، ۱۹۷۱) تعمیم پیدا کرد. نسخه نمونه متناهی از نقطه فروریزش شامل نقطه فروریزش جمعی^{۱۵} (ABP) و نقطه فروریزش جایگزین^{۱۶} (RBP) است. که بوسیله دوناهو و هوبر^{۱۷} (۱۹۸۳) معرفی شده و برای ارزیابی استواری برآوردها استفاده شده که در زیر چگونگی بدست آوردن آنها نشان داده می شود.

^{۱۱} Casella & Berger

^{۱۲} Finite sample breakdown point

^{۱۳} Hodge

^{۱۴} Hampe

^{۱۵} Addition breakdown point

^{۱۶} Replacement breakdown point

^{۱۷} Donoho and Huber

فرض کنید $X^n = \{X_1, \dots, X_n\}$ یک نمونه به اندازه n در R^d باشد. نقطه فروریزش جمعی نمونه متناهی از یک برآوردگر T در X^n به صورت زیر تعریف می شود:

$$ABP(T, X^n) = \min \left\{ \frac{m}{m+n} : \sup_{Y^m} \|T(X^n \cup Y^m) - T(X^n)\| = \infty \right\} \quad (9.6.1)$$

که Y^m دلالت بر مجموعه داده هایی از اندازه n با مقادیر دلخواه و $X^n \cup Y^m$ دلالت بر نمونه آلوده شده بوسیله مقادیر Y^m اضافه شده به X^n می باشد. نقطه فروریزش جایگزین نمونه متناهی از یک برآوردگر T در X^n به صورت زیر تعریف می شود:

$$RBP(T, X^n) = \min \left\{ \frac{m}{n} : \sup_{X_m^m} \|T(X_m^m \cup Y^m) - T(X^n)\| = \infty \right\} \quad (10.6.1)$$

که X_m^m دلالت بر نمونه آلوده شده از X^n را دارد که بوسیله جایگزینی m نقطه از X^n با مقادیر دلخواه، ایجاد شده است.

به بیان دیگر ABP و RBP ی یک برآوردگر به ترتیب کمترین کسر اضافه شده و کسر جایگزینی است که می تواند اریبی برآوردگر را به سمت بینهایت ببرد.

۷-۱ روابط کمی بین ABP و RBP

با توجه به ABP و RBP ی معرفی شده در بخش ۱-۶، این دو تعریف از نقطه فروریزش اغلب به طور جداگانه برای برآوردگرها مورد بحث واقع می شوند و به صورت دو مفهوم کاملاً مستقل رفتار می کنند. برخی نویسندگان علاقه مند به استفاده از ABP در بحث های استواری برآوردگرها هستند در حالی که برخی دیگر RBP را برای بررسی خواص استواری ترجیح می دهند و اعتقاد دارند که ساده تر، واقع گرایانه تر و به طور کلی کاربردی تر است. برای این که بتوان به قضاوت کامل تری در مورد استواری برآوردگرها دست یافت با بیان روابط کمی بین این دو در کلاس بزرگی از برآوردگرها می توان با بدست آوردن اندازه یکی از دو معیار دیگری را نیز بدست آورد.

یک نکته جالب در مورد نقاط فروریزش یک برآوردگر T در X^n این است که در بسیاری از موارد، آنها فقط به اندازه نمونه n وابسته هستند و از ترکیب و پیکربندی X^n مستقل هستند (هر چند از تعاریف این چنین به نظر می آید که آنها بایستی به نمونه X^n وابسته باشند). به بیان دیگر برای نقاط فروریزش که در حالت کلی در ارزیابی استواری یک برآوردگر استفاده می شوند این خاصیت «آزاد از نمونه» به طور قطع مطلوب است.

در زیر توجه خود را به نقاط فروریزشی منحصر می کنیم که از شکل X^n مستقل هستند و RBP هایی که به فرم $\frac{\lfloor an+b \rfloor}{n}$ یا $\frac{\lfloor an+b \rfloor}{n}$ هستند. در واقع این مورد برای بیشتر برآوردگرهای مکانی و پراکنش صادق است. در اینجا $\lfloor x \rfloor$ ، نشان دهنده بزرگترین عدد صحیح کوچکتر یا مساوی x و $\lceil x \rceil$ کوچکترین عدد صحیح بزرگتر یا مساوی x است.

قضیه ۱-۱:

فرض کنید T یک برآوردگر در R^d باشد و X^n یک مجموعه داده با اندازه n باشد.

الف) نقطه فروریزش جمعی T می تواند به صورت زیر بیان شود:

$$ABP(T, X^n) = (\lfloor \varepsilon n + \beta \rfloor + m) / (n + \lfloor \varepsilon n + \beta \rfloor + m)$$

اگر نقطه فروریزش جایگزین T بتواند بصورت زیر بیان شود:

$$RBP(T, X^n) = \lfloor (\varepsilon n + (\beta + m + \varepsilon)) / (\varepsilon + 1) \rfloor / n$$

ب) نقطه فروریزش جایگزین T می تواند به صورت زیر بیان شود:

$$RBP(T, X^n) = \lfloor (\varepsilon n + (\beta + m + r)) / (\varepsilon + 1) \rfloor / n$$

اگر نقطه فروریزش جمعی T بتواند بصورت زیر بیان شود:

$$ABP(T, X^n) = (\lfloor \varepsilon n + \beta \rfloor + m) / (n + \lfloor \varepsilon n + \beta \rfloor + m)$$

که ε و β و عدد صحیح m برخی ثابتهایی هستند که از n و X^n مستقل هستند و $0 \leq \beta \leq \varepsilon \leq 1$ ،

$$\max\{- (\varepsilon n + \beta + m) / (\varepsilon + 1), \max_N \{D^-(N)\}\} \leq r / (\varepsilon + 1) < 1 + \min_N \{D^-(N)\}$$

با

$$N = \lfloor (\varepsilon + 1)n + \beta + m \rfloor$$

و

$$D^-(n) = \lfloor (\varepsilon n + \beta + m) / (\varepsilon + 1) \rfloor - (\varepsilon n + \beta + m) / (\varepsilon + 1)$$

اثبات: زو^{۱۸} ۱۹۹۹ رابینید.

قضیه ۲-۱:

فرض کنید T یک برآوردگر در R^d باشد و X^n یک مجموعه با اندازه n باشد.

الف) نقطه فروریزش جمعی T می تواند به صورت زیر بیان شود:

$$ABP(T, X^n) = (\lceil \varepsilon n + \beta \rceil + m) / (n + \lceil \varepsilon n + \beta \rceil + m)$$

اگر نقطه فروریزش جایگزین T بتواند به صورت زیر نوشته شود:

$$RBP(T, X^n) = \lceil (\varepsilon n + \beta + m + \varepsilon) / (\varepsilon + 1) \rceil / n$$

ب) نقطه فروریزش جایگزین T می تواند به صورت زیر نوشته شود

$$RBP(T, X^n) = \lceil (\varepsilon n + (\beta + m + r)) / (\varepsilon + 1) \rceil / n$$

اگر نقطه فروریزش جمعی T بتواند به صورت زیر نوشته شود

$$ABP(T, X^n) = (\lceil \varepsilon n + \beta \rceil + m) / (n + \lceil \varepsilon n + \beta \rceil + m)$$

که ε و β و عدد صحیح m برخی ثابتهایی هستند که از n و X^n مستقل هستند

$$0 \leq \beta \leq \varepsilon \leq 1, \text{ و}$$

$$\{\max\{-(\varepsilon n + \beta + m) / (\varepsilon + 1), \max_N \{D^+(N)\} - 1\} \leq r / (\varepsilon + 1) < \min_N \{D^+(N)\}$$

با $N = \lfloor (\varepsilon + 1)n + \beta + m \rfloor$ و

$$D^+(n) = \lceil (\varepsilon n + \beta + m) / (\varepsilon + 1) \rceil - (\varepsilon n + \beta + m) / (\varepsilon + 1)$$

اثبات: زو ۱۹۹۹ رابینید.

مثال ۱-۱:

RBP ی میانه نمونه ای را در R^1 با استفاده مستقیم از ABP ی آن بدست می آوریم.

ABP ی میانه نمونه ای در دوناهو و هوپر (۱۹۸۳) داده شده است که آن مقداری برابر با $\frac{1}{4}$

است. برای استفاده نمودن از قضیه ۱-۱ نیاز است که ε و β و m را مشخص نمائیم. آنها به

صورت زیر در نظر گرفته می شوند $\beta = 0$ و $\varepsilon = 1$.

^{۱۸} Zou

حال با استفاده از قضیه ۱-۱، RBP ی برآوردگر میانه نمونه ای برابر $\lfloor (n+r)/2 \rfloor / n$ می باشد به طوری که :

$$\max \left\{ \frac{1}{2}, \max_N \{ \lfloor N/2 \rfloor - N/2 \} \right\} \leq r/2 < 1 + \min_N \{ \lfloor N/2 \rfloor - N/2 \}$$

و $N = \lfloor 2n \rfloor = 2n$. بنابراین نتیجه می گیریم که $1 \leq r < 2$. لذا، RBP ی برآوردگر می تواند به صورت $\lfloor (n+1)/2 \rfloor / n$ نوشته شود، که می توان نشان داد این RBP در واقع همانند RBP یی است که از تعریف اصلی RBP مشتق می شود.

مثال ۱-۲:

ABP ی L_1 -میانه که در نوشته ها میانه فضایی هم نامیده می شود را در R^d ، ($d \geq 1$) به طور مستقیم از RBP بدست می آوریم. L_1 -میانه به صورت بردار T تعریف می شود به طوری که :

$$T = \arg \min_{x \in R^d} \sum_{i=1}^n \|X_i - x\| \quad (11.7.1)$$

توجه شود که منظور از $\arg \min$ ی مجموع فوق، $x \in R^d$ یی است که به ازای آن مجموع کلیه فواصل اقلیدسی X_i از x به حداقل برسد.

برای هر نمونه X_1, \dots, X_n لویا و روسو^{۱۹} (۱۹۹۱) نشان داده اند که :

$$RBP(T, X^n) = \lfloor (n+1)/2 \rfloor / n$$

برای به کار بردن قضیه ۱، نیاز داریم که ε ، β و m را مشخص کنیم. به نظر می رسد برای این برآوردگر $m=0$ ، $\beta=0$ و $\varepsilon=1$ به طور منحصر به فردی مشخص می شوند. حال با استفاده از قضیه ۱ داریم که ABP ی L_1 -میانه $\frac{1}{4}$ است.

۸-۱ توزیع بیضوی

$F_{\mu, \Sigma}$ دارای توزیع بیضوی است هر گاه چگالی آن به فرم زیر باشد:

$$f(x) = (\det(\Sigma))^{-\frac{n}{2}} f \left[(x-\mu)^T \Sigma^{-1} (x-\mu) \right] \quad (12.8.1)$$

^{۱۹} Lopuha? and Rousseuw

که $f: [0, \infty) \rightarrow [0, \infty)$ یک تابع معین است و $(\mu, \Sigma) \in \Theta$ که $\Theta = R^p \times PDS(p)$ و PDS^{*} ماتریس متقارن معین مثبت است.

۹-۱ فضای پوچ

فضای پوچ یک ماتریس $m \times n$ مانند A ، به صورت مجموعه زیر تعریف می شود:

$$\text{Null}(A) = \{X \in R^n : AX = 0\}$$

که 0 دلالت بر بردار صفر با m مؤلفه دارد. معادله ماتریسی $AX = 0$ معادل، معادلات همگن زیر می باشد:

$$\begin{aligned} AX = 0 \Leftrightarrow \begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = 0 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = 0 \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = 0 \end{cases} \end{aligned}$$

با این نمایش فضای پوچ حل مجموعه معادلات همگن فوق می باشد.

۱۰-۱ تحلیل همبستگی کانونی

تحلیل همبستگی کانونی به شناخت و کمی کردن رابطه بین دو مجموعه از متغیرها پرداخته، اولین بار توسط هتلینگ (۱۹۳۵) مطرح شد. در اصل تحلیل همبستگی کانونی در مورد همبستگی بین یک ترکیب خطی از متغیرهای یک مجموعه و یک ترکیب خطی از متغیرهای مجموعه دیگر متمرکز می شود. ابتدا هدف این است که دو ترکیب خطی با بیشترین همبستگی تعیین شود. سپس دو ترکیب خطی را تعیین می کنیم که در میان تمام زوجهای ناهمبسته با زوج انتخاب شده اول دارای بیشترین همبستگی باشد. زوجهای ترکیبات خطی را متغیرهای کانونی و همبستگی آنها را همبستگی های کانونی می نامند. همبستگی های کانونی شدت ارتباط بین دو مجموعه از متغیرها را اندازه می گیرد. آنچه تحلیل همبستگی کانونی انجام می دهد این است که روابط بین دو مجموعه را که براساس تعداد زیاد متغیرها پایه ریزی شده

^{۲۰} Positive Definite Symmetric