



دانشکده علوم ریاضی

گروه آمار

پایان نامه جهت اخذ درجه کارشناسی ارشد در رشته آمار ریاضی

عنوان پایان نامه

m —برآوردگرهای استوار مکانی و مقیاسی

استاد راهنما

دکتر مجید سرمد

استاد مشاور

دکتر ناصر رضا ارقامی

به اهتمام

محمد ولی احمدی

شهریور ۸۸

ما بدان مقصد عالی نتوانیم رسید
هم مگر پیش نهد لطف شما گامی چند

پروردگار بزرگ را سپاس می گزارم، که توفیق نگارش این رساله را به بنده‌ی حقیر خویش، عنایت فرموده و همواره مرا مورد لطف و رحمت بی‌پایان خود قرار داده است.

رساله‌ی حاضر شرح شاگردی این‌جانب در محضر بزرگانی بوده است که خود را پیوسته مرهون عنایت‌ها و تشویق‌های ایشان می‌دانم و هم اکنون بر خود لازم می‌دانم که از آنان قدردانی کنم.

نخست از بزرگترین سرمایه‌های زندگی‌م، **پدر و مادر عزیزم**، که همواره مشوق و پشتیبانم بودند، متشکرم.

از استاد فرهیخته و بزرگوام، **جناب آقای دکتر مجید سرمد**، که هدایت و راهنمایی این‌جانب را قبول کردند و همواره با دقت نظر در طول نگارش این رساله، چراغ راهم بودند، تشکر می‌کنم.

از **جناب آقای دکتر ناصر رضا ارقامی**، استاد گرانقدر مشاورم، که در نگارش این رساله، همکاری صمیمانه داشتند، کمال تشکر خود را ابراز می‌کنم.

همچنین از اساتید محترم داور، **جناب آقای دکتر خنجری و سرکار خانم دکتر حبیبی** که زحمت مطالعه‌ی این رساله را بر خود هموار کردند و نکات لازم را یادآور شدند، تشکر می‌کنم.

در پایان، از همه‌ی دوستان و عزیزانی که به نحوی مرا در نوشتن این پایان‌نامه یاری نمودند، سپاس‌گزاری می‌کنم.

محمد ولی احمدی

شهریور ماه ۱۳۸۸



باسمه تعالی
مشخصات پایان نامه تحصیلی دانشجویان
دانشگاه فردوسی مشهد

عنوان پایان نامه: m -برآوردگرهای استوار مکانی و مقیاسی

نویسنده: محمد ولی احمدی استاد راهنما: مجید سرمد

دانشکده: علوم ریاضی گروه: آمار

رشته تحصیلی: آمار ریاضی تاریخ دفاع: ۱۳۸۸/۶/۲۱

مقطع تحصیلی: کارشناسی کارشناسی ارشد دکتری تعداد صفحات: ۱۳۹

چکیده پایان نامه:

داده‌های جمع‌آوری شده در قالب موارد شامل یک یا چند نقطه‌ی پرت هستند. نقاطی که از توده‌ی مشاهدات نمونه جدا قرار گرفته‌اند و یا از الگوی کلی که برای داده‌ها فرض کرده‌ایم، پیروی نمی‌کنند. علاوه بر این ثابت شده است برآوردگرهای کلاسیکی که تاکنون از ویژگی‌های بهینه برای برآورد پارامتر جامعه برخوردار بوده‌اند، در اثر انحراف از پذیره‌های اساسی، کارایی و اعتبار خود را از دست می‌دهند. در حالی که با استفاده از روش‌های آماری استوار، در صورت انحراف از پذیره‌های اساسی نیز همچنان می‌توان برآورد دقیقی را از پارامتر جامعه به‌دست آورد.

در این پایان‌نامه علاوه بر این که سعی کرده‌ایم روشی را برای شناسایی نقاط پرت در داده‌های تصادفی به‌دست آمده از جامعه‌ای نرمال بیان کنیم، به معرفی و بیان ویژگی‌های دسته‌ای از مشهورترین برآوردگرهای استوار مکانی و مقیاسی با عنوان m -برآوردگرها پرداخته‌ایم. علاوه بر این معیارهایی را برای سنجش استواری برآوردگرهای مختلف معرفی کرده و در نهایت m -برآوردگری را می‌یابیم که علاوه بر استواری از کارایی بالایی نیز برخوردار باشد. m -برآوردگر حاصل به m -برآوردگر همپل معروف است.

امضای استاد راهنما

کلید واژه‌ها: درجه‌ی بدی، مقادیر Z ، مقادیر Z اصلاح شده، توزیع اسلش،

تاریخ

m -برآوردگر، نقطه‌ی شکست، تابع تأثیر، خطای ناخالص حساسیت



In the name of God
Dissertation details
Ferdowsi University of Mashhad

Title : *Location and Scale Robust M-estimators*

Student's name : *Mohammad Vali Ahmadi* **Major Professor :** *Majid Sarmad*

School : Mathematical Sciences **Department:** Statistics

Major : Mathematical Statistics **Date of defense :** 12/9/2009

Number of pages: 139 **For degree :** BS MS Ph.D

Abstract :

The collected data often contains one or more outlier which are either isolated from the bulk of data or do not follow the general model that we have supposed for data.

Furthermore, it has been proved that the classic estimators which are used to estimate the distribution parameter and have the optimal properties, will lose their efficiency and validity due to deviation from fundamental assumptions. An exact estimation of distribution parameter can be obtained using robust statistical methods, even there exists a deviation from fundamental assumptions .

In this thesis, we want to find a method to identify outlier in random data which is given from a normal distribution. In addition, a class of robust estimator of location and scale as m-estimators is studied. Also, two quantities for measuring robustness of different estimators and finally an m-estimator with high efficiency and robustness, i.e. Hampel's m-estimator will be studied.

Signature of major professor:

Keywords : *badness, Z-scores, modified Z-scores, slash distribution, m-estimator, breakdown point, influence function, gross error sensitivity*

فهرست مندرجات

۱۰	پیش گفتار	۱-۰
۱۲	شناسایی نقاط پرت بر اساس مقادیر Z اصلاح شده	۱
۱۳	مقدمه	۱-۱
۱۷	معرفی توزیع اسلش	۲-۱
۱۷	تابع چگالی احتمال توزیع اسلش	۱-۲-۱
۱۸	تابع توزیع اسلش	۲-۲-۱
۱۹	چگونگی یافتن مقادیر برش بر اساس حجم نمونه به دست آمده	۳-۱
۲۰	نحوه شبیه سازی	۴-۱
۲۳	مراحل شبیه سازی	۵-۱
۲۴	نتایج شبیه سازی	۱-۵-۱
۲۹	مقایسه استفاده از مقادیر برش وابسته به حجم نمونه با مقدار برش ثابت	۶-۱
۳۴	معرفی دو معیار برای سنجش استواری برآوردگرها	۲
۳۵	مقدمه	۱-۲
۴۰	تابع تأثیر	۲-۲
۴۵	نقطه شکست	۳-۲
۴۵	نقطه شکست در نمونه های متناهی	۱-۳-۲
۴۷	نقطه شکست در حالت مجانبی	۲-۳-۲
۵۱	m -برآوردگرهای استوار مکانی و مقیاسی	۳
۵۲	مقدمه	۱-۳
۵۲	m -برآوردگرهای مکانی	۲-۳
۵۶	مقایسه برآوردگرهای ML چهار توزیع مختلف	۱-۲-۳
۷۱	m -برآوردگرهای مقیاسی	۳-۳
۷۶	نحوه یافتن m -برآوردگر مکانی با خاصیت هم پایایی مقیاسی	۴-۳
۷۷	استفاده از m -برآوردگرهای استودنت شده	۱-۴-۳
۷۷	برآورد همزمان پارامترهای مکان و مقیاس	۲-۴-۳

۷۹ محاسبات عددی برای یافتن m -برآوردگرهای استوار	۵-۳
۷۹ محاسبه m -برآوردگرهای مکانی	۳-۵-۱
۸۱ محاسبه m -برآوردگرهای مقیاسی	۳-۵-۲
۸۲ نقطه شکست m -برآوردگرها	۳-۶-۶
۸۲ نقطه شکست m -برآوردگرهای مکانی	۳-۶-۱
۸۵ نقطه شکست m -برآوردگرهای مقیاسی	۳-۶-۲
۸۷ تابع تأثیر m -برآوردگرها	۳-۷
۹۴	خواص مجانبی m-برآوردگرها	۴
۹۵ مقدمه	۴-۱
۹۵ وجود و یکتایی m -برآوردگرها	۴-۲
۹۸ سازگاری m -برآوردگرها	۴-۳
۱۰۳ نرمال بودن مجانبی m -برآوردگرها	۴-۴
.		
۱۰۷	رویکرد همپل در تعادل بین استواری و کارایی	۵
۱۰۸ مقدمه	۵-۱
۱۰۹ B -استوارترین m -برآوردگر	۵-۲
۱۱۳ رویکرد همپل در تعادل بین استواری و کارایی	۵-۳
۱۲۶	آینده تحقیق	
۱۲۷	کتابنامه	
۱۳۰	ضمیمه	
۱۳۱ ضمیمه الف: اثبات چند رابطه	
۱۳۳ ضمیمه ب: برنامه توابع شبیه‌سازی شده توسط نرم‌افزار آماری R مربوط به فصل یک	
۱۳۷ ضمیمه ج: برنامه‌های نوشته شده در mathematica مربوط به مثالهای فصل پنج	

فهرست جداول

۱۶	مقادیر Z مشاهدات مثال ۱-۱-۱.....	۱-۱
۱۶	مقادیر Z اصلاح شده‌ی مشاهدات مثال ۱-۱-۱.....	۲-۱
۲۵	مقادیر $\hat{\lambda}_m^S$ (مقدار برشی که در ازای آن درجه‌ی بدی مقیاس بندی شده از توزیع اسلش کمترین می‌شود.) و B_m^S (مقدار درجه‌ی بدی مقیاس بندی شده از توزیع اسلش در ازای مقدار برش $\hat{\lambda}_m^S$) وابسته به حجم نمونه.....	۳-۱
۲۶	مقادیر d (ثابتی که از تقسیم MAD_n بر آن، برآوردگری نارایب برای σ در توزیع $N(\mu, \sigma^2)$ به دست می‌آید.) وابسته به حجم نمونه.....	۴-۱
۲۸	مقادیر λ^* نهایی (مقدار برش نهایی) وابسته به حجم نمونه به همراه تعداد شبیه‌سازی‌های صورت گرفته برای رسیدن به نتیجه‌ای پایدار.....	۵-۱
۳۱	مقادیر $D^S(\lambda_n^*)$ و $D^G(\lambda_n^*)$ وابسته به حجم نمونه.....	۶-۱
۳۲	مقادیر Z اصلاح شده‌ی مشاهدات مثال ۱-۶-۱ با استفاده از مقدار ثابت d در تعریف M_i ..	۷-۱
۳۳	مقادیر Z اصلاح شده‌ی مشاهدات مثال ۱-۶-۱ با استفاده از مقدار ثابت d وابسته به حجم نمونه در تعریف M_i	۸-۱
۳۷	مقادیر $ARE(\varepsilon)$ (کارایی جانبی نسبی دو برآوردگر SD_n و D_n) در ازای ε های مختلف .	۱-۲
۳۹	n برابر واریانس‌های میانگین و میانه نمونه در مثال ۳-۱-۲.....	۲-۲
۵۷	چگالی توزیع‌های نرمال، لجستیک، نمایی دوگانه و کوشی و توابع ρ و ψ متناظر با برآوردگرهای ML این توابع توزیع.....	۱-۳
۶۴	مقادیر واریانس جانبی m -برآوردگر هیوبر (v) زمانی که مشاهدات بر اساس توزیع آمیخته‌ی $F_\mu = (1 - \varepsilon)G_\mu + \varepsilon H_\mu$ به دست آمده‌اند که در آن $G_\mu = N(\mu, 1)$ و $H_\mu = N(\mu, 100)$	۲-۳
۹۱	تأثیر اعمال تغییرات در بردار مشاهدات بر برآوردگرهای مختلف مکانی.....	۳-۳

فهرست اشکال

- ۱-۱ نمودار تابع چگالی احتمال توزیع‌های نرمال استاندارد و اسلش ۱۷
- ۲-۱ نمودار مقادیر شبیه‌سازی شده $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ (درجات بدی مقیاس‌بندی شده بر اساس توزیع‌های گوسین و اسلش) در ازای مقادیر برش $\lambda = \{0, 0.75, \dots, 4\}$ برای $n = 5 \dots$ ۲۱
- ۳-۱ بر اساس این شکل محل تلاقی مقادیر شبیه‌سازی شده $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ همان مقدار برش بهینه (λ^*) است. ۲۶
- ۴-۱ نمودار ۳ مرتبه تکرار شبیه‌سازی مقادیر $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ برای $n = 5$ به جهت اطمینان از پایداری نتایج ۲۷
- ۵-۱ مقادیر λ^* بر اساس حجم نمونه‌های زوج و فرد به صورت جداگانه ۲۸
- ۱-۲ نمودار تابع تأثیر برآوردهای میانگین و میانه نمونه در توزیع نرمال استاندارد ۴۴
- ۱-۳ نمودار توابع چگالی توزیع‌های نرمال، لجستیک، نمایی دوگانه و کوشی زمانی که پارامتر $\mu = 0$ است. ۵۹
- ۲-۳ نمودار توابع ψ متناظر با برآوردهای ML توابع توزیع جدول (۱-۳) ۵۹
- ۳-۳ توابع ρ و ψ متناظر با $-m$ برآوردهای میانگین، میانه و هیوبر ۶۱
- ۴-۳ تابع وزن هیوبر بر اساس رابطه‌ی (۳-۲-۱۳) ۶۵
- ۵-۳ نمودار توابع ρ متناظر با $-m$ برآوردهای همپل، سینوسی، دو مربعی و بل ۶۸
- ۶-۳ نمودار توابع ψ متناظر با $-m$ برآوردهای همپل، سینوسی، دو مربعی و بل ۶۹
- ۷-۳ تابع وزن دو مربعی بر اساس رابطه‌ی (۳-۲-۱۴) ۷۰
- ۸-۳ نمودار توابع g و h بنا بر رابطه‌ی (۳-۵-۴) بر اساس ρ -تابع و ψ -تابع بیسکوار با $k_i = 4 / 68$ ۸۱

۰-۱) پیش گفتار

در بسیاری از داده‌های جمع آوری شده به مشاهدات پرتی برمی‌خوریم که تحلیل داده‌ها را با مشکل مواجه می‌کنند. از این رو آماردانان همواره به دنبال راهی بوده‌اند تا تأثیر این نقاط را بر نتیجه‌گیریها کاهش دهند و در عین حال از اطلاعاتی که نقاط پرت در خود نهفته دارند، به بهترین نحو استفاده کنند. در همین راستا در این پایان‌نامه سعی کرده‌ایم تا روشهایی را برای برخورد با نقاط پرت ارائه دهیم.

در فصل اول، ابتدا مقادیر Z اصلاح شده مشاهدات نمونه را معرفی می‌کنیم. سپس بر اساس روشی کم‌بیشینه که به تفصیل آن را شرح می‌دهیم، سعی می‌کنیم مقادیر برشی را متناسب با حجم مشاهدات نمونه به دست آوریم. بر این اساس می‌توانیم از مقایسه‌ی مقادیر Z اصلاح شده‌ی مشاهدات و مقادیر برش به دست آمده، مقادیر پرت را در میان توده مشاهدات تشخیص دهیم.

در فصل دوم ابتدا با ذکر چند مثال به مقایسه برآوردهای استوار و غیر استوار می‌پردازیم. سپس دو معیار نقطه‌ی شکست و تابع تأثیر را برای سنجش میزان استواری برآوردها معرفی می‌کنیم. در توضیح برآوردهای استوار باید گفت که این برآوردها در صورتی که مشاهدات نمونه حاوی نقاط پرت باشند، نیز همچنان برآورد صحیحی را از پارامتر جامعه حاصل می‌کنند.

در فصل سوم m -برآوردهای مکانی و مقیاسی را - که دسته‌ای از برآوردهای استوار هستند - معرفی می‌کنیم و چگونگی محاسبه آنها را شرح می‌دهیم. در نهایت نقطه‌ی شکست و تابع تأثیر این برآوردها را به دست می‌آوریم.

در فصل چهارم خواص مجانبی m -برآوردهای استوار را بررسی می‌کنیم و توزیع مجانبی آنها را به دست می‌آوریم.

و در فصل پنجم به دنبال m -برآوردگری هستیم که علاوه بر استواری از کارایی بالایی نیز برخوردار باشد. در همین راستا نیز m -برآوردگر بهینه همپل را معرفی کرده و آن را به دست می آوریم.

فصل ۱

شناسایی نقاط پرت بر اساس مقادیر Z اصلاح

شده

در این فصل، ابتدا مقادیر Z اصلاح شده مشاهدات نمونه را معرفی می‌کنیم. سپس بر اساس روشی کم‌بیشینه که به تفصیل آن را شرح می‌دهیم، سعی می‌کنیم مقادیر برشی را که در تعیین نقاط پرت استفاده می‌شوند، متناسب با حجم مشاهدات نمونه به دست آوریم. بر این اساس می‌توانیم از مقایسه مقادیر Z اصلاح شده مشاهدات و مقادیر برش به دست آمده، مقادیر پرت را در میان توده مشاهدات تشخیص دهیم.

(۱-۱) مقدمه

اگرچه زمانی که صحبت از نقطه پرت می‌شود، مفهوم واحدی از آن برداشت می‌شود، اما بررسی‌ها نشان می‌دهد، تعاریف متنوع و متفاوتی درباره‌ی این نقطه از سوی آماردانان ارائه شده است.

میزان اعتماد به هر مشاهده به رابطه آن با سایر مشاهداتی بستگی دارد که تحت شرایط یکسان به دست آمده‌اند. به عقیده‌ی اکثر محققین مشاهده‌ای که دور از توده مشاهدات قرار می‌گیرد، نقطه پرت محسوب می‌شود. از این میان می‌توان به بارنت و لویس (۱۹۸۴، صفحه ۴) و هاوکینز (۱۹۸۰، صفحه ۱) اشاره کرد. بنا بر بارنت و لویس (۱۹۸۴) "نقطه‌ی پرت مشاهده‌ای است که با سایر مشاهدات نمونه ناسازگار و متفاوت است." علاوه بر این در تعریفی مشابه هاوکینز (۱۹۸۰) نقطه پرت را مشاهده‌ای معرفی می‌کند که "از الگوی کلی داده‌ها پیروی نمی‌کند، به گونه‌ای که به نظر می‌رسد بر اساس فرایندی متفاوت با سایر داده‌ها به دست آمده است."

اما بکمن و کوک (۱۹۸۳، صفحه ۱۲۱) عقیده دارند "مشاهده پرت، یا مشاهده‌ای ناهماهنگ و متفاوت است یا مشاهده‌ای آلوده". مشاهده ناهماهنگ "مشاهده‌ای است که حضور آن در میان مشاهدات غیر طبیعی به نظر می‌رسد". اما مشاهده آلوده "مشاهده‌ای است که از توزیعی غیر از توزیع فرض شده برای جامعه به دست آمده است."

در این نوشتار تعریف بارنت و لویس (۱۹۸۴) را از مشاهده پرت می‌پذیریم. وقوع این دسته از مشاهدات در اکثر نمونه‌های تصادفی تقریباً امری طبیعی است. بنابراین اولین دغدغه و نگرانی محققین پس از جمع‌آوری

مشاهدات نمونه، چگونگی برخورد با نقاط پرت است. به دلیل این که شناسایی علل و عوامل وقوع نقاط پرت، کیفیت فرآیند نمونه‌گیری را ارتقا می‌دهد و یا در برخی موارد توزیعی را که برای داده‌ها فرض کرده‌ایم بهبود می‌بخشد، بررسی این مقادیر از اهمیت بسیاری برخوردار است. اما آنچه در ابتدا مهم است، شناسایی این مقادیر در میان توده مشاهدات است.

بررسی ظاهری مشاهدات، اولین راه برای تشخیص مقادیر پرت است. اما این روش در مجموع از اعتبار کافی برخوردار نیست. بر اساس تحقیقی که کالت و لويس (۱۹۷۶) بر داده‌های تک‌متغیره انجام دادند، دریافتند شناسایی نقاط پرت در این داده‌ها به عوامل زیر بستگی دارد:

(۱) نحوه بررسی مقادیر نمونه برای شناسایی مقادیر پرت: اگر برای شناسایی نقاط پرت، داده‌های تصادفی را در سه حالت مرتب شده، مرتب نشده و نموداری بررسی کنیم، شانس تشخیص صحیح این مقادیر در داده‌های مرتب شده بسیار بیشتر است.

(۲) مقیاس مقادیر نمونه: اگر مقیاس اندازه‌گیری مشاهدات را افزایش دهیم، احتمال اینکه مشاهدات بزرگ‌تر را، پرت در نظر بگیریم افزایش خواهد یافت.

بنابراین دلایل می‌توان گفت بررسی ظاهری مشاهدات به تنهایی روش مطمئنی برای شناسایی نقاط پرت نیست. علاوه بر این در نمونه‌های پیچیده‌تر با حجم‌های بزرگ‌تر (رگرسیون چند متغیره، طرح‌های چند عاملی و ...) این کار تقریباً غیر ممکن است. پس در این موارد بهتر است به جای بررسی ظاهری داده‌ها، از معیارهای منطقی‌تری در این زمینه بهره بگیریم. واضح است که با استفاده از این معیارها نمی‌توان مشاهدات پرت را به طور قطعی پیدا کرد، بلکه تنها مقادیر مشکوکی را می‌توان یافت که ارزش بیشتری دارد تا درباره پرت بودنشان تحقیق شود.

چون اکثر پدیده‌های طبیعی از توزیع نرمال پیروی می‌کنند، در این مقاله راه‌حلی را برای شناسایی نقاط پرت در داده‌های نرمال ارائه داده‌ایم. مشهورترین روش برای شناسایی نقاط پرت در داده‌های نرمال، استفاده از مقادیر Z مشاهدات است. بر اساس نمونه تصادفی X_1, \dots, X_n این مقادیر برابرند با:

$$Z_i = \frac{X_i - \bar{X}}{S},$$

که \bar{X} و S به ترتیب میانگین و انحراف معیار نمونه‌اند. چون در توزیع گوسین (نرمال استاندارد)، ۹۹/۷ درصد مشاهدات بین -۳ و ۳ قرار می‌گیرند، مشاهداتی که $|Z_i| > 3$ ، نقطه پرت فرض می‌شوند. اما این

روش به خصوص در نمونه‌های کوچک دقیق نیست. شفلر (۱۹۸۸) نشان داد که مقدار مطلق Z_i در نمونه‌ای به حجم n ، حداکثر برابر $\sqrt{n}/(n-1)$ است. بنابراین در نمونه‌ای به حجم 10° ،

$$\forall i ; |Z_i| < 2/85 < 3$$

در حالی که غیرمنطقی است هر نمونه 10° تایی ممکن را از توزیع نرمال فرض کنیم.

برای رفع این مشکل ایگلیویکس و هوگلین (۱۹۹۳) از برآوردگرهای استوار بهره گرفتند. به طور کلی با استفاده از برآوردگرهای استوار در صورتی که درصدی از مشاهدات نمونه پرت باشند نیز می‌توان برآورد دقیقی را از پارامتر جامعه به دست آورد. ایگلیویکس و هوگلین برای برآورد پارامتر مکان در جامعه از میانه مشاهدات (\tilde{X}) و برای برآورد پراکندگی در جامعه از برآوردگری که همپل (۱۹۷۴) به صورت زیر معرفی کرده است و به میانه انحراف از میانه شهرت دارد، استفاده کردند.

$$MAD_n = Med\{|X_i - \tilde{X}|\}.$$

و در نهایت مقادیر Z اصلاح شده مشاهدات را برابر زیر تعریف کردند:

$$M_i = \frac{X_i - \tilde{X}}{\frac{MAD_n}{d}}$$

که ثابت d در این عبارت سبب می‌شود، MAD_n/d برآوردگری ناریب برای σ باشد. چون بر اساس نمونه‌های با حجم زیاد از توزیع $N(\mu, \sigma^2)$ ، $E(MAD_n) = 0.6745\sigma$ است، ایگلیویکس و هوگلین ثابت d را برابر 0.6745 در نظر گرفتند. به علاوه بر اساس بررسی‌های شبیه‌سازی بر روی داده‌های نرمال، تصمیم گرفتند مشاهداتی را که در $|M_i| > 3/5$ صدق می‌کنند، نقطه پرت در نظر بگیرند.

در مثال زیر دو روشی را که تاکنون برای شناسایی نقاط پرت معرفی کرده‌ایم، مقایسه می‌کنیم.

*** مثال ۱-۱-۱)** داده‌های فرضی زیر را در نظر بگیرید:

$$2/1, 2/6, 2/4, 2/5, 2/3, 2/1, 2/3, 2/6, 8/2, 8/3$$

میانگین و انحراف معیار داده‌ها برابرند با:

$$\bar{X} = 3/54, \quad S = 2/49$$

وجود مقادیر $2/8$ و $3/8$ در بین داده‌ها سبب کاهش دقت دو بر آوردگر میانگین و انحراف معیار شده است. مقادیر Z داده‌ها در جدول (۱-۱) درج شده است. مشاهده می‌کنیم که:

$$\forall i; |z_i| < 3$$

جدول ۱-۱: جدول مقادیر Z مشاهدات مثال ۱-۱-۱

i	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
x_i	۲/۱	۲/۶	۲/۴	۲/۵	۲/۳	۲/۱	۲/۳	۲/۶	۸/۲	۸/۳
z_i	-۰/۵۸	-۰/۳۸	-۰/۴۶	-۰/۴۲	-۰/۵۰	-۰/۵۸	-۰/۵۰	-۰/۳۸	۱/۸۷	۱/۹۱

در حالی که دو مقدار $2/8$ و $3/8$ مقادیری کاملاً پرت در میان داده‌ها هستند، اما مقادیر z_i مشاهدات، این امر را تأیید نمی‌کنند. حال اگر دو مقدار $\tilde{X} = 2/45$ و $MAD_n = 0/15$ را برای برآورد میانگین و انحراف معیار جامعه در نظر بگیریم، مقادیر M_i مشاهدات برابر جدول (۲-۱) خواهد بود.

جدول ۲-۱: جدول مقادیر Z اصلاح شده‌ی مشاهدات مثال ۱-۱-۱

i	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
x_i	۲/۱	۲/۶	۲/۴	۲/۵	۲/۳	۲/۱	۲/۳	۲/۶	۸/۲	۸/۳
m_i	-۱/۵۷	۰/۶۷	-۰/۲۲	۰/۲۲	-۰/۶۷	-۱/۵۷	-۰/۶۷	۰/۶۷	۲۵/۸۶	۲۶/۳۱

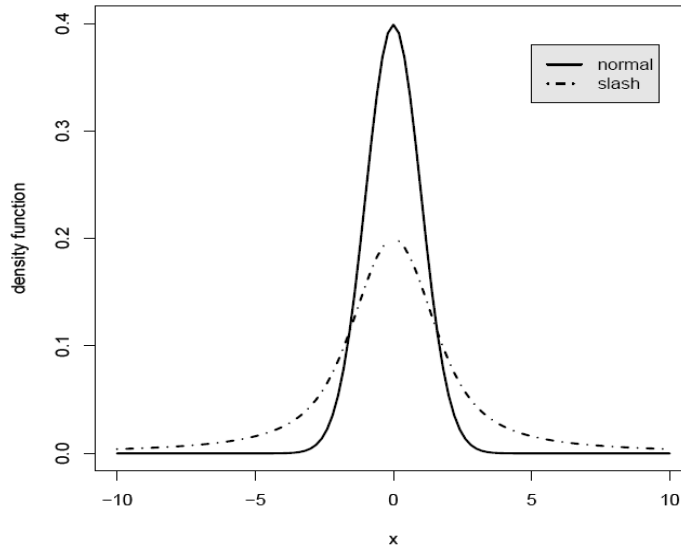
چون m_i مشاهدات $2/8$ و $3/8$ بیش از $3/5$ است، بنابراین بر اساس این معیار می‌توان دو مقدار پرت را در میان مشاهدات پیدا کرد. این مثال به وضوح برتری مقادیر Z اصلاح‌شده را بر مقادیر Z معمولی در شناسایی نقاط پرت نشان می‌دهد.

در ادامه قصد داریم، ثابت d را که ایگلوپکس و هوگلین برای نمونه‌های با حجم‌های متفاوت برابر مقدار ثابت $6745/0$ فرض کردند، برای حجم نمونه‌های تصادفی $n = 5, 6, \dots, 29, 30$ به صورت مجزا شبیه‌سازی کنیم. همچنین قصد داریم، بر خلاف ایگلوپکس و هوگلین که مقدار برش ثابتی را برای شناسایی نقاط پرت استفاده کردند، از مقادیر برش وابسته به حجم نمونه در این زمینه بهره بگیریم و بر اساس روشهای شبیه‌سازی این مقادیر را بیابیم.

قبل از پرداختن به این موضوع، توزیع اسلش را که در تولید نقاط پرت استفاده می‌شود، معرفی می‌کنیم.

(۲-۱) معرفی توزیع اسلش

توزیع اسلش همانند توزیع گوسین نسبت به مبدأ متقارن است، اما دمهای سنگین آن سبب می‌شود تا با کمک آن بتوان مقادیر بسیار بزرگی را در نمونه‌ی تصادفی تولید کرد، به گونه‌ای که این مقادیر در مقایسه با مقادیر حاصل از توزیع گوسین به راحتی پرت به نظر رسند. در شکل (۱-۱) توابع چگالی توزیع‌های گوسین و اسلش رسم شده است.



شکل ۱-۱: نمودار تابع چگالی احتمال توزیع‌های نرمال استاندارد و اسلش

تعریف (۱-۲-۱) اگر X و Y متغیرهای تصادفی مستقل و به ترتیب از توزیع‌های $N(0,1)$ و $U(0,1)$ برخوردار باشند، متغیر تصادفی

$$S = \frac{X}{Y}$$

دارای توزیع اسلش خواهد بود.

(۱-۲-۱) تابع چگالی احتمال توزیع اسلش

ابتدا توزیع توأم متغیرهای تصادفی $S = \frac{X}{Y}$ و $W = Y$ را به دست می‌آوریم. ژاکوبین این تبدیل برابر است با:

$$J = \begin{vmatrix} \frac{\partial x}{\partial s} & \frac{\partial x}{\partial w} \\ \frac{\partial y}{\partial s} & \frac{\partial y}{\partial w} \end{vmatrix} = \begin{vmatrix} w & s \\ 0 & 1 \end{vmatrix} = w .$$

بنابراین توزیع توأم متغیرهای تصادفی S و W برابر است با:

$$g_{S,W}(s,w) = f_{X,Y}(sw,w)|J| = \frac{w}{\sqrt{\nu\pi}} e^{-\frac{s^2 w^2}{\nu}} \quad ; \quad 0 \leq w \leq 1, s \in \mathbb{R}$$

و اما تابع چگالی احتمال متغیر تصادفی S برابر است با:

$$\begin{aligned} f_S(s) &= \int_0^1 g_{S,W}(s,w)dw = \frac{1}{\sqrt{\nu\pi}} \int_0^1 w e^{-\frac{s^2 w^2}{\nu}} dw \\ &= \frac{1}{s^2} \left(\frac{1}{\sqrt{\nu\pi}} - \frac{1}{\sqrt{\nu\pi}} e^{-\frac{s^2}{\nu}} \right) \\ &= \frac{1}{s^2} (\phi(0) - \phi(s)), \end{aligned}$$

که در آن ϕ تابع چگالی احتمال نرمال استاندارد است.

۱-۲-۲) تابع توزیع اسلش

ابتدا تابع توزیع شرطی S به شرط $Y = y$ را می‌یابیم:

$$\begin{aligned} P(S \leq s | Y = y) &= P\left(\frac{X}{Y} \leq s | Y = y\right) \\ &= P\left(\frac{X}{y} \leq s | Y = y\right) \\ &= P(X \leq sy | Y = y) \\ &= P(X \leq sy) = \Phi(sy), \end{aligned}$$

که در آن Φ تابع توزیع نرمال استاندارد است. بنابراین تابع توزیع اسلش برابر است با:

$$\begin{aligned} F_S(s) &= \int_0^1 P(S \leq s | Y = y) f_Y(y) dy = \int_0^1 \Phi(sy) dy \\ &= \frac{1}{s} \int_0^s \Phi(t) dt = \frac{1}{s} \left(t\Phi(t) \Big|_0^s - \int_0^s t\phi(t) dt \right) \\ &= \frac{1}{s} \left(s\Phi(s) - \frac{1}{\sqrt{\nu\pi}} e^{-\frac{t^2}{\nu}} \Big|_0^s \right) = \Phi(s) - \frac{1}{s} (\phi(0) - \phi(s)). \end{aligned}$$

(۳-۱) چگونگی یافتن مقادیر برش بر اساس حجم نمونه به دست آمده

در ابتدا تأکید می‌کنیم که همانند روش ایگلیویکس و هوگلین قصد داریم از مقادیر Z اصلاح شده‌ی مشاهدات، برای شناسایی نقاط پرت استفاده کنیم. مسأله در این جا یافتن مقادیر برش بهینه مبتنی بر حجم نمونه‌ی به دست آمده، است.

برای یافتن مقادیر برش بهینه، استراتژی کم‌بیشینه (مینیماکس) را پیشنهاد می‌کنیم. به عبارت دیگر ابتدا تابع زیانی را تعیین کرده و دو حالت را در نمونه‌ی تصادفی در نظر می‌گیریم. حالت اول زمانی که هیچ نقطه‌ی پرتی در مشاهدات وجود ندارد و حالت دوم زمانی که نمونه دارای نقاط پرت بسیاری است. مقدار برش بهینه حد وسطی است برای این دو حالت، به نحوی که ماکسیمم زیان ما در ازاء مقادیر برش مختلف کمترین گردد. ابتدا معیار زیر را به عنوان تابع زیان مسأله تصمیم تعریف می‌کنیم.

تعریف (۱-۳-۱) درجه‌ی بدی: برای نمونه‌ی تصادفی X_1, \dots, X_n از توزیعی با میانگین μ ، درجه‌ی بدی مقدار برش λ برابر است با:

$$B(\lambda) = (\bar{X}_\lambda - \mu)^2,$$

که \bar{X}_λ میانگین مشاهدات باقیمانده در نمونه است، پس از حذف مقادیری که $|M_i| > \lambda$.

در تعریف درجه‌ی بدی فرض بر این است که هدف از جمع‌آوری داده‌ها، برآورد میانگین جامعه است. به نظر می‌رسد میانگین نمونه پس از حذف مشاهدات پرت، برآوردگری منطقی برای این منظور باشد. بنابراین برای هر توزیعی همانند F ، تعریف می‌کنیم:

$$b(\lambda) = E_F(B(\lambda)),$$

که امید ریاضی بر اساس نمونه‌های تصادفی به حجم n از این توزیع به دست آمده است. فرض کنید λ_m مقداری است که در ازاء آن $b(\lambda)$ می‌نیمم می‌شود، یعنی $b(\lambda_m) = \min_\lambda b(\lambda)$. به علاوه فرض می‌کنیم $b_m = b(\lambda_m)$. بنابراین درجه‌ی بدی مقیاس‌بندی شده را به صورت زیر تعریف می‌کنیم:

$$b_{SC}(\lambda) = \frac{b(\lambda)}{b_m}.$$

بنابر دلائل زیر از درجه‌ی بدی مقیاس‌بندی‌شده برای یافتن مقدار برش بهینه استفاده می‌کنیم:

(۱) بهترین مقدار برش برای مشاهدات نمونه زمانی است که $b(\lambda)$ کمترین مقدار خود یعنی b_m را اختیار کند. از طرفی $b(\lambda)$ معیاری است که میزان ناکارایی میانگین نمونه را - پس از حذف مقادیر پرت بر اساس مقدار برش λ - برای برآورد میانگین جامعه نشان می‌دهد. بنابراین $b_{SC}(\lambda)$ ناکارایی نسبی برآوردگر میانگین را در مقایسه با بهترین حالت ممکن آن نشان می‌دهد.

(۲) به جهت این که بتوانیم به ازای λ های معین، درجات بدی توزیع‌های مختلف را با یکدیگر مقایسه کنیم، استفاده از $b_{SC}(\lambda)$ را ترجیح می‌دهیم.

حال برای یافتن مقدار برش بهینه، از دو توزیع کاملاً متفاوت در تولید نقاط پرت استفاده می‌کنیم: توزیع گوسین (نرمال استاندارد) و توزیع اسلش که نقاط پرت بسیاری را - در مقایسه با توزیع گوسین - تولید می‌کند.

بنابراین به روش کم‌بیشینه، مسأله‌ی تصمیمی را حل می‌کنیم که تابع زیان آن مطابق با درجه‌ی بدی مقیاس‌بندی‌شده تعریف می‌شود و بر این اساس مقدار برش بهینه - که آن را با λ^* نشان می‌دهیم - مقداری است که در رابطه‌ی زیر صدق می‌کند:

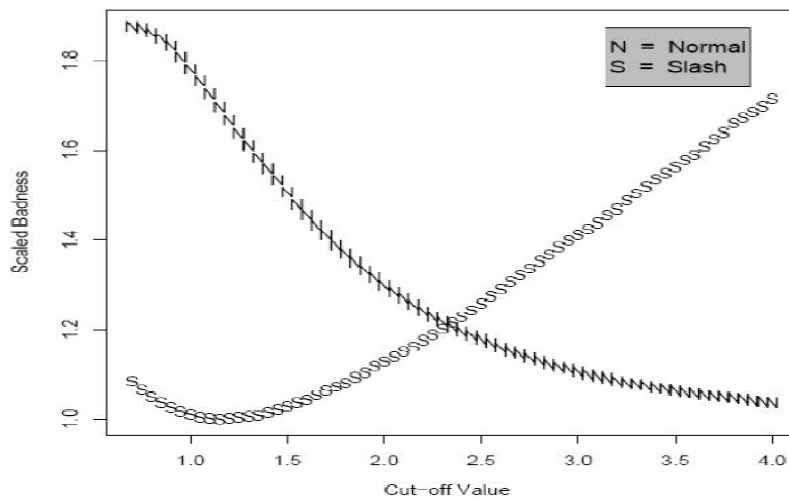
$$\max\left(b_{SC}^G(\lambda^*), b_{SC}^S(\lambda^*)\right) = \min_{\lambda} \max\left(b_{SC}^G(\lambda), b_{SC}^S(\lambda)\right),$$

که $b_{SC}^G(\cdot)$ و $b_{SC}^S(\cdot)$ به ترتیب درجات بدی مقیاس‌بندی‌شده‌ی توزیع‌های گوسین و اسلش را نشان می‌دهد. به عبارت دیگر به ازاء مقادیر برش مختلف، مقداری را انتخاب می‌کنیم که ماکزیمم درجه بدی مقیاس‌بندی‌شده‌ی آن در بین دو توزیع گوسین و اسلش کمترین باشد.

۴-۱) نحوه شبیه‌سازی

زمانی که نمونه تصادفی دقیقاً بر اساس توزیع گوسین به‌دست آمده است، برای این که میانگین نمونه - پس از حذف مقادیر پرت - برآورد دقیقی را از پارامتر μ حاصل کند، منطقی است که مقدار برش λ را بزرگترین مقدار ممکن در نظر بگیریم تا تمامی مشاهدات نمونه در برآورد پارامتر μ شرکت داشته باشند. در این صورت هر چه λ بزرگ‌تر باشد، درجه‌ی بدی $b^G(\lambda)$ کاهش خواهد یافت. (همین طور $b_{SC}^G(\lambda)$) بنابراین $b_{SC}^G(\lambda)$ تابعی نزولی بر حسب λ است. اما اگر مشاهدات نمونه دارای نقاط پرت باشند (که در اینجا برای تولید نمونه‌های تصادفی دارای نقاط پرت بسیار از توزیع اسلش استفاده کرده‌ایم) مقادیر بزرگ λ باعث می‌شود، مشاهدات پرت همچنان در نمونه باقی مانند و دقت برآورد پارامتر μ را کاهش دهند. در این

صورت با افزایش λ ، $b^S(\lambda)$ افزایش خواهد یافت. (همین طور $b_{SC}^S(\lambda)$ بنابراین $b_{SC}^S(\lambda)$ تابعی صعودی بر حسب λ است. در شکل (۲-۱) مبتنی بر شبیه‌سازی برای حجم نمونه‌ی برابر ۵، مقادیر برآوردشده $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ را در ازای λ های مختلف رسم کرده‌ایم. (نحوه رسم این شکل را در بخش‌های بعد توضیح می‌دهیم.) در این شکل صعودی بودن b_{SC}^S و نزولی بودن b_{SC}^G (به استثناء همسایگی عدد یک) کاملاً مشهود است.



شکل ۲-۱: نمودار مقادیر شبیه‌سازی شده $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ در ازای $n = 5$ برای $\lambda = \{0/7, 0/75, \dots, 3/95, 4\}$

بر اساس شکل (۲-۱) محل تلاقی توابع $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ بر روی محور افقی، مقدار برش بهینه λ^* را مشخص می‌کند. حال کافی است نحوه به‌دست آوردن λ^* را شرح دهیم. رابطه‌ی زیر را تعریف می‌کنیم:

$$\delta(\lambda) = b_{SC}^G(\lambda) - b_{SC}^S(\lambda).$$

مقدار λ^* از حل معادله‌ی $\delta(\lambda^*) = 0$ حاصل می‌شود. اما چون $b_{SC}^G(\lambda)$ و $b_{SC}^S(\lambda)$ مقادیری نامعلوم‌اند، به روش شبیه‌سازی و به ازای λ های معین $\{\lambda_1, \lambda_2, \dots, \lambda_k\}$ آنها را برآورد می‌کنیم. قبل از بیان نحوه شبیه‌سازی b_{SC}^S و b_{SC}^G تأکید می‌کنیم که مقدار دقیق λ^* به حجم نمونه‌های تصادفی (n) وابسته است. ابتدا N_1 نمونه تصادفی به حجم n از توزیع گوسین و N_2 نمونه تصادفی به حجم n توزیع اسلش به‌دست می‌آوریم. به منظور کاهش واریانس تفاوت بین متوسط درجات بدی در نمونه‌های گوسین و اسلش، بهتر است ابتدا نمونه‌های اسلش را تولید کرده و سپس با استفاده از رابطه‌ی زیر نمونه‌های گوسین را به‌دست آوریم.