



پایان نامه کارشناسی ارشد در رشته آمار ریاضی

عنوان :

آزمون نرمال بودن در حضور داده‌های پرت

استاد راهنما:

دکتر هادی جباری نوقابی

استاد مشاور:

دکتر ناصر رضا ارقامی

نگارنده:

حسین صابر جیدرق

بهمن ۱۳۸۹

اظہار نامہ

اینجانب حسین صابر جبرقی دانشجوی دورہ کارشناسی ارشد رشته آمار ریاضی دانشکده ریاضی دانشگاه فردوسی مشهد نویسنده رسالہ/پایان نامہ آزمون نرمال بودن در حضور داده‌های پرت تحت راهنمایی دکتر ہادی جباری نوقابی متعہد می‌شوم:

- تحقیقات در این رسالہ/پایان نامہ توسط اینجانب انجام شدہ است و از صحت و اصالت برخوردار است.
- در استفادہ از نتایج پژوهشہای محققان دیگر بہ مرجع مورد استفادہ استناد شدہ است.
- مطالب مندرج در رسالہ/پایان نامہ تاکنون توسط خود یا فرد دیگری برای دریافت ہیچ نوع مدرک یا امتیازی در ہیچ جا ارائه شدہ است.
- کلیہ حقوق معنوی این اثر متعلق بہ دانشگاه فردوسی مشهد می‌باشد و مقالات مستخرج با نام « دانشگاه فردوسی مشهد » و یا « Ferdowsi University of Mashhad » بہ چاپ خواہد رسید.
- حقوق معنوی تمام افرادی کہ در بہ دست آمدن نتایج اصلی رسالہ/پایان نامہ تأثیرگذار بودہ‌اند در مقالات مستخرج از رسالہ/پایان نامہ رعایت شدہ است.
- در کلیہ مراحل انجام این رسالہ/پایان نامہ، در مواردی کہ از موجود زندہ (یا بافتہای آنها) استفادہ شدہ است ضوابط و اصول اخلاقی رعایت شدہ است.
- در کلیہ مراحل انجام این رسالہ/پایان نامہ، در مواردی کہ بہ حوزہ اطلاعات شخصی افراد دسترسی یافتہ یا استفادہ شدہ است، اصل رازداری، ضوابط و اصول اخلاق انسانی رعایت شدہ است.

تاریخ امضای دانشجو

مالکیت نتایج و حق نشر

- کلیہ حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامہ های رایانہ‌ای، نرم افزارها و تجهیزات ساختہ شدہ) متعلق بہ دانشگاه فردوسی مشهد می‌باشد. این مطلب باید بہ نحو مقتضی در تولیدات علمی مربوطہ ذکر شود.
- استفادہ از اطلاعات و نتایج موجود در رسالہ/پایان نامہ بدون ذکر مرجع مجاز نمی‌باشد.



بِسْمِ تَعَالَى

Graduate Studies Thesis/Dissertation Information
Ferdowsi University of Mashhad

Author: Hossein Saber Jabdaragh
Supervisor(s): Dr. Hadi Jabbari Noughabi
Advisor(s): Prof. Naserreza Arghami

Faculty: Mathematical Sciences

Department: Statistics

Specialization: Mathematical Statistics

Approval Date:

Defence Date: 23/1/2011

M.Sc.

Ph.D.

Number of Pages: 107

Abstract:

Statistical models are often based on normal distributions and procedures for testing this distributional assumption are needed. Many goodness-of-fit tests suffer from the presence of outliers, in the sense that they may reject the null hypothesis even in the case of a single extreme observation. In this thesis, we show a possible extension of the Shapiro-Wilk test that is not affected by such a problem. The presented method is inspired by the forward search (FS), a new, recently proposed by Daniele Coin (2008), diagnostic tool. An application to univariate observations shows how the procedure is able to capture the structure of the data, even in the presence of outliers.

Signature of Supervisor:

Key Words:

Date:

1. Computational algorithm
2. Forward search
3. Least absolute deviation
4. Outlier
5. Shapiro-Wilk test

تقدیم به

پدر و مادر مهربان و عزیزم

قدردانی

بدین وسیله مراتب قدردانی و تشکر خود را نسبت به استاد راهنمای بزرگوام آقای دکتر هادی جبّاری نوقابی و استاد مشاورم آقای دکتر ناصر رضا ارقامی ابراز می‌دارم و برای ایشان توفیق روزافزون آرزومندم که بدون راهنمایی‌های ایشان در مراحل تحقیق پایان‌نامه، به انجام رساندن این پژوهش میسر نبود. همچنین سپاس و قدردانی خود را به اساتید محترم دکتر جعفر احمدی و دکتر غلام رضا محتشمی که قبول زحمت نموده و پایان‌نامه‌ام را مورد مطالعه و داوری قرار داده‌اند، تقدیم می‌نمایم.

از کارمندان دانشکده، واحد انتشارات، اداره آموزش و بخش کتابخانه، دوستان خوبم آقایان بابک نصیری، سعید امیرنژاد، عبدالسعید توماج، یاسر سنچولی، مجتبی اصفهانی و همه کسانی که به نوعی برگردن بنده حقی دارند، سپاسگزارم.



بسمه تعالی .
مشخصات رساله /پایان نامه تحصیلی دانشجویان .
دانشگاه فردوسی مشهد

عنوان رساله /پایان نامه: آزمون نرمال بودن در حضور داده‌های پرت

نام نویسنده: حسین صابر جیدرق

نام استاد(ان) راهنما: دکتر هادی جباری نوقابی

نام استاد(ان) مشاور: دکتر ناصررضا ارقامی

رشته تحصیلی: آمار ریاضی	گروه: آمار	دانشکده : علوم ریاضی
تاریخ دفاع: ۱۳۸۹/۱۱/۳	تاریخ تصویب:	
تعداد صفحات: ۱۰۷	<input type="radio"/> دکتری	<input checked="" type="radio"/> کارشناسی ارشد

چکیده رساله /پایان نامه : اغلب مدل‌های آماری روی توزیع‌های نرمال پایه ریزی می‌شوند و روش‌هایی برای آزمون کردن این توزیع مورد نیاز است. بسیاری از آزمون‌های نیکویی برازش از حضور داده‌های پرت تأثیر منفی می‌پذیرند، به این معنی که ممکن است آن‌ها فرضیه صفر را حتی به خاطر وجود یک مشاهده بسیار بزرگ یا بسیار کوچک رد کنند. در این مطالعه بسطی از آزمون شاپیرو-ویلک را ارائه می‌دهیم که از چنین مشکلی تأثیر نمی‌پذیرد. روش حاضر الهام گرفته از روش جستجوی پیشرو^۱ (FS) است که یک روش تشخیصی جدید است و اخیراً توسط دانیل کوین (۲۰۰۸) پیشنهاد شده است. با به کار گیری مشاهدات یک متغیر نشان داده می‌شود که این روش قادر است ساختار داده‌ها را حتی در حضور داده‌های پرت تشخیص دهد.

امضای استاد راهنما:	کلید واژه:
تاریخ:	۱. آگوریتم محاسبه ۲. جستجوی پیشرو ۳. کمترین انحراف مطلق ۴. داده‌ی پرت ۵. آزمون شاپیرو-ویلک

فهرست مندرجات

۶	فهرست جداول
۸	فهرست شکل‌ها
۹	نمادها
۱۰	پیش‌گفتار
۱۳	مقدمه و کلیات ۱
۱۴	تاریخچه ۱.۱
۱۷	آزمون‌های نیکویی برآزش ۲.۱

۱۷ آزمون کولموگروف - اسمیرنف	۱.۲.۱
۱۹ معیار V ی کوپر	۲.۲.۱
۱۹ آزمون کرامر-ون میسز	۳.۲.۱
۲۰ آزمون اندرسون - دارلینگ	۴.۲.۱
۲۰ آزمون های اصلاح شده مبتنی بر EDF	۵.۲.۱
۲۱ آزمون کی دو پیرسن	۶.۲.۱
۲۹ آزمون شاپیرو - ویلک	۷.۲.۱
۳۲ داده‌های پرت	۳.۱
۳۵ برخی از دلایل اتفاق افتادن داده های پرت	۱.۳.۱
۳۶ چرا باید داده‌های پرت شناسایی شوند	۲.۳.۱
۳۶ آزمون‌هایی برای k تا داده پرت	۳.۳.۱
۳۷ آزمون های تشخیص داده پرت یکتا	۴.۳.۱
۳۹ آزمون‌های تشخیص داده‌های پرت k تایی	۵.۳.۱
۴۰ تأثیر داده‌های پرت بر آزمون‌های نیکویی برازش	۴.۱
۴۱ برآوردگر رگرسیونی کمترین میانه مربعات (LMS)	۵.۱
۴۱ مروری بر برآوردگرهای رگرسیونی دیگر	۱.۵.۱
۴۴ نقطه از کار افتادگی و برآوردگر LMS	۲.۵.۱

۴۸	کشف داده‌های پرت در توزیع نرمال	۲
۵۰	کشف داده‌های پرت در توزیع نرمال با واریانس‌های نابرابر	۱.۲
۵۱	روش کمترین مربعات	۱.۱.۲
	الگوریتم روش کمترین انحراف مطلق (LAD) برای کشف داده‌های	۲.۲
۵۳	پرت تحت توزیع نرمال	
۵۵	روش LAD و الگوریتم محاسبه آن	۱.۲.۲
۵۸	الگوریتم روش LAD	۲.۲.۲
۶۵	روش پیشنهادی برای کشف داده‌های پرت در توزیع نرمال*	۳
	روش کلی حذف یکی، یکی ترکیباتی از داده‌ها و بررسی نرمال بودن	۱.۳
۶۶	باقی مانده‌ی داده‌ها	
۶۷	روش ماکسیمم تفاضل مرتبه اول آماره‌های ترتیبی	۲.۳
۶۸	آماره M	۱.۲.۳
۶۸	مقادیر بحرانی M	۲.۲.۳
۷۰	تشخیص داده‌های پرت	۳.۳

۷۳	استفاده از آماره M در حالت کلی	۴.۳
۷۳	ترکیب روش LAD با روش پیشنهادی	۵.۳
۷۸	بزرگترین زیر مجموعه دارای توزیع نرمال از یک نمونه تصادفی	۴
۸۰	آماره‌ی شاپیرو-ویلک	۱.۴
۸۲	نسخه جستجوی پیشرو آزمون شاپیرو-ویلک برای نرمال بودن	۲.۴
۸۴	اندازه W_F	۱.۲.۴
۸۶	به کارگیری و رفتار W_F	۳.۴
۸۸	توان تجربی W_F	۱.۳.۴
۹۳	نتیجه گیری و آینده تحقیق	
۹۴	ضمیمه	
۹۶	کتاب نامه	

۱۰۲	واژه‌نامه‌ی فارسی به انگلیسی
-----	-------	------------------------------

فهرست جداول

عنوان	صفحه
جدول ۱.۲: تعداد داده‌های پرت، Q_k و برآورد پارامترها	۵۳
جدول ۲.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 1$	۶۱
جدول ۳.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 2$	۶۱
جدول ۴.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 3$	۶۲
جدول ۵.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 4$	۶۲
جدول ۶.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 1$	۶۳
جدول ۷.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 2$	۶۳
جدول ۸.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 3$	۶۳
جدول ۹.۲: تعداد داده‌های پرت، $SMAD$ و برآورد پارامترها برای $j = 4$	۶۴
جدول ۱.۳: برخی مقادیر بحرانی برای M در سطح $\alpha = 0.05$	۶۹
جدول ۲.۳: نمونه تصادفی از توزیع نرمال	۷۲
جدول ۳.۳: نمونه تصادفی از توزیع نرمال	۷۶
جدول ۴.۳: تعداد داده‌های پرت، LAD و برآورد پارامترها	۷۷

-
- جدول ۵.۳: تعداد داده‌های پرت، LAD و برآورد پارامترها ۷۷
- جدول ۱.۴: کران تجربی ناحیه رد برای $n = 50$ ۸۴
- جدول ۲.۴: برآورد پارامترهای رابطه (۴-۹) ۸۵

فهرست شکل‌ها

صفحه	عنوان
۴۲.....	شکل ۱.۱: خط رگرسیونی کمترین مربعات خطا با یک داده پرت
۴۳.....	شکل ۲.۱: خط رگرسیونی کمترین مقدار مطلق با یک داده پرت
۴۶.....	شکل ۳.۱: خط رگرسیونی کمترین مربعات میانه با یک داده پرت
۸۸.....	شکل ۱.۴: نمودار نتایج نسخه پیشرو آزمون شاپیرو-ویلک
۹۰.....	شکل ۲.۴: توان W_F در مقابل توزیع‌های مقابل متقارن
۹۱.....	شکل ۳.۴: توان W_F در مقابل توزیع‌های مقابل غیرمتقارن

نمادها

LMS	برآوردگر رگرسیونی کمترین مربعات میانه
LS	برآوردگر رگرسیونی کمترین مربعات خطا
L_1	برآوردگر رگرسیونی کمترین مقدار مطلق
LAD	الگوریتم روش کمترین انحراف مطلق
EDF	تابع توزیع تجربی

پیش‌گفتار

در این تحقیق آزمون نرمال بودن نمونه‌ای از مشاهدات که دارای داده‌های پرت است، بررسی می‌شود. نرمال بودن یکی از عمومی‌ترین فرض‌ها در توسعه و استفاده از روش‌های آماری است. بنابراین روش‌هایی برای آزمون کردن این توزیع (توزیع نرمال) مورد نیاز است. این مسأله توسط محققین زیادی مورد مطالعه و بررسی قرار گرفته است و در طول این بررسی آزمون‌هایی را یافتیم که تصور نمی‌کردیم وجود داشته باشند! در این بررسی حدود چهار روش آزمون را یافتیم که برای نرمال بودن پیشنهاد شده است و به همان تعداد روش‌هایی مثل رسم نمودار، آزمون داده‌های پرت، آزمون‌های نیکویی برازش تعمیم یافته و آزمون‌های دیگر که در کشف غیر نرمال بودن داده‌ها مفید هستند، وجود دارد. بنابراین، این بررسی ذره‌ای از این دریای بی‌کران است که مطمئناً یک بررسی جامع با تمام جزئیات نیست.

روش‌های زیادی برای آزمون نرمال بودن نمونه‌های یک متغیره در ادبیات آماری پیشنهاد شده است. فقط با در نظر گرفتن آزمون‌های با فرضیه صفر مرکب، آزمون‌های نیکویی برازش برای نرمال بودن را می‌توان در چهار گروه طبقه بندی کرد. اولی شامل اندازه فاصله بین تابع توزیع نظری و تابع توزیع تجربی است. دومین دسته شامل آماره‌هایی مرکب از

ضرایب چولگی و کشیدگی است. سوّمین خانواده با تعمیم آزمون نیکویی برازش پیرسن^۱ پایه ریزی می‌شود. چهارمین کلاس مربوط به روش‌های رگرسیون است. بالاخره آخرین آن‌ها آزمون‌های نیکویی برازش مبتنی بر آنتروپی می‌باشد. برای اطلاعات بیشتر در مورد آزمون‌های نیکویی برازش مبتنی بر آنتروپی به پایان‌نامه خانم فاطمه یوسف‌زاده (۱۳۸۷) مراجعه کنید.

مشکل عمومی بسیاری از این آزمون‌ها حساسیت به حضور داده‌های پرت در نمونه است. در حقیقت، تنها یک مشاهده پرت می‌تواند منجر به رد فرضیه صفر (فرضیه نرمال بودن) شود حتی اگر اکثر داده‌ها از توزیع نرمال به دست آمده باشند.

در این تحقیق یک روش اصلی که با وفق دادن روش جستجوی پیشرو (FS) به آزمون‌های نیکویی برازش حاصل می‌شود، را معرفی می‌کنیم. بدین وسیله مشکل تأثیر مشاهدات پرت روی آزمون نیکویی برازش حل می‌شود. به طور اساسی FS از یک زیر مجموعه پرت از مشاهدات شروع می‌کند و مقیاس برای پیشرفت در جستجو (افزایش زیر مجموعه با استفاده از افزودن یک یا چند مشاهده در هر مرحله) مجموعه‌ای از روش‌های تشخیصی است که در طول جستجو به کار می‌روند (برای اطلاعات بیشتر آتکینسون و ریانی a ، b ، 2000 ، 2001 ، a ، 2002 ، b ، 2002 ، 2004 ، آتکینسون و همکاران 2004 را ببینید).

به منظور وفق دادن FS برای آزمون نرمال بودن، به یک روش برای انتخاب یک زیر مجموعه پرت از مشاهدات و روش دیگر برای آزمون نرمال بودن نیاز داریم.

^۱ Pearson Goodness of Fit Test

این تحقیق شامل چهار فصل می‌باشد. فصل اول مقدمه و کلیات می‌باشد. در فصل اول تاریخچه‌ای از آزمون نرمال بودن، آزمون‌های نیکویی برازش مهم از جمله آزمون شاپیروویلک، تعریف داده‌های پرت و روش‌های تشخیص این داده‌ها، برآوردگر رگرسیونی کمترین میانه مربعات و در کل مفاهیم و تعاریف مورد نیاز مطرح می‌شود. فصل دوم شامل دوروش است که محققین برای کشف داده‌های پرت در مشاهداتی که از جامعه نرمال به دست آمده‌اند، پیشنهاد کرده‌اند. در فصل سوم روش پیشنهادی برای کشف داده‌های پرت در توزیع نرمال و در صورت وجود، یافتن زیرمجموعه‌ای نرمال از نمونه‌ای تصادفی مورد بحث قرار می‌گیرد. این فصل را با * مشخص کردیم که یک روش جدید و پیشنهاد خودمان است. بالاخره، فصل چهارم به روش جستجوی پیشرو آزمون شاپیروویلک برای نرمال بودن اختصاص یافته است که ساختار داده‌ها را با وجود داده‌های پرت به خوبی تشخیص می‌دهد. در نهایت یاد آوری می‌شود که در تنظیم این پایان‌نامه برای برگردان انگلیسی به فارسی از واژه نامه آمار و ریاضی استفاده شده است. در صورتی که واژه‌ای در این واژه نامه وجود نداشته است، با نظر استاد راهنما برگردان فارسی آن به کار برده شده است.

فصل ۱

مقدمه و کلیات

۱.۱ تاریخچه

روش‌های آماری مثل آزمون استودنت، آزمون‌هایی برای ضرایب رگرسیونی، آنالیز واریانس و آزمون فیشرفرای همگنی واریانس، همگی دارای یک فرض اساسی هستند که باید نمونه از توزیع نرمال به دست آمده باشد. البته در روش‌های آماری یا باید فرض نرمال بودن را آزمون کنیم که در صورت رد فرضیه صفر از روش‌های ناپارامتری برای بررسی داده‌ها استفاده کنیم، یا این که نشان دهیم در صورتی که داده‌ها دارای توزیع نرمال نباشند، تأثیری بر روی روش‌های بررسی مشاهدات ندارد. به عبارت دیگر باید نشان دهیم که نرمال بودن یا نبودن مشاهدات هیچ فرقی ندارند. بسیاری از محققان آماری با ارزیابی دامنه تأثیر تخطی از این فرض روی سطح معنی داری آزمون یا اثربخشی برآورد کردن پارامتر، را بررسی کرده‌اند.

ارزیابی محققان در مورد تأثیرات تخطی از فرض نرمال بودن روی روش‌های استاندارد آماری به تاریخ قبل از مقاله بارتلت (۱۹۳۵) روی آزمون استودنت برمی‌گردد. فیشرفرای (۱۹۶۴) در مورد این مسأله تحقیق کرد و نتایج مربوط به انباشتگی آماره‌های چولگی و کشیدگی در آزمون‌های نرمال بودن را توسعه داد.

مسأله بررسی نرمال بودن توزیع برای استواری^۱ آزمون فرضیه‌ها در ادبیات آماری مورد توجه قرار گرفته است. مرور این ادبیات و آزمون‌های نرمال بودن، تلاش و همت بسیاری از نظریه پردازان و شاغلان آماری را نشان می‌دهد.

پیتمن (a و b ۱۹۳۷) مسأله حساسیت آزمون استودنت و آنالیز واریانس یک طرفه را بررسی کرد و بدین ترتیب با استفاده از یک تبدیل نظری به نتایج مهمی برای عدم حساسیت آزمون استودنت و آنالیز واریانس یک طرفه به تخطی از نرمال بودن را ارائه نمود.

^۱Robustness