

بسمه تعالی



دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران)

دانشکده ریاضی و علوم کامپیوتر

پایان نامه کارشناسی ارشد آمارمحض

عنوان:

رگرسیون حداقل قدرمطلق انحرافات وزنی استوار

دانشجو:

منیژه محمودی

استاد راهنما:

جناب آقای دکتر صادق رضایی

استاد مشاور:

جناب آقای دکتر سعید رضاخواه

مهر ۱۳۸۵



تاریخ:

پیوست:

فرم اطلاعات پایان نامه

کارشناسی ارشد و دکترا

معاونت پژوهشی

معادل

بورسیه

دانشجوی آزاد

منیزه محمودی

نام و نام خانوادگی:

رشته تحصیلی: آمار ریاضی

دانشگاه: ریاضی و علوم کامپیوتر

۸۳۱۱۳۱۰۹

شماره دانشجویی:

نام و نام خانوادگی استاد راهنما: جناب آقای دکتر صادق رضایی

عنوان پایان نامه به فارسی: رگرسیون حداقل قدر مطلق انحرافات وزنی استوار

عنوان پایان نامه به انگلیسی: Robust weighted least absolute deviations regression

نظری

توسعه ای

بنیادی

کاربردی

کارشناسی ارشد:

نوع پروژه:

دکتری

تعداد واحد: ۶

تاریخ خاتمه: ۸۵/۷/۵

۸۴/۸/۲۰

تاریخ شروع:

سازمان تأمین کننده اعتبار: ندارد

واژه های کلیدی به فارسی: رگرسیون استوار؛ رگرسیون حداقل قدر مطلق انحرافات؛ رگرسیون چندک

واژه های کلیدی به انگلیسی: Robust regression, Least absolute deviations regression, Quantile regression

نظرها و پیشنهادات به منظور بهبود فعالیتهای پژوهشی دانشگاه: دسترسی به منابع علمی از جمله مقالات و مجلات لاتین معتبر آسانتر شود.

استاد راهنما: جناب آقای دکتر صادق رضایی

دانشجو: منیزه محمودی

تاریخ:

امضاء استاد راهنما:

نسخه ۱: معاونت پژوهشی

نسخه ۲: کتابخانه و به انضمام دوجلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز اسناد و مدارک علمی

تقدیم به پدر و مادر مهربان و خواهر عزیزم

که سرچشمه های امید و مهربانی در زندگیم هستند و تشکر از

زحمات و مهر بی دریغ این فرشتگان که مرا در تمام مراحل زندگیم

یاری کردند.

تشکر و قدردانی

فرصت مناسبی است برای سپاسگزاری و قدردانی از زحمات ارزنده اساتید محترم، جناب آقای دکتر صادق رضایی و جناب آقای دکتر سعید رضاخواه که با راهنماییهای مفیدشان، همچون پدری دلسوز و مهربان مرا در اتمام این پایان نامه یاری رساندند. همچنین سپاس صمیمانه خود را به جناب آقای دکتر اسماعیل خرم و جناب آقای دکتر حمید پزشک که زحمت داوری این پایان نامه را برعهده گرفته اند، تقدیم می کنم.

فهرست مطالب:

| | |
|----|---------------|
| ۸ | چکیده |
| ۱۰ | علائم اختصاری |

فصل اول کلیات

| | |
|----|-------------|
| ۱۲ | ۱-۱ مقدمه |
| ۱۵ | ۲-۱ تاریخچه |
| ۱۸ | ۳-۱ ساختار |

فصل دوم مشاهدات مؤثر

| | |
|----|-------------------------------|
| ۲۰ | ۱-۲ مقدمه |
| ۲۱ | ۲-۲ مشاهدات مؤثر |
| ۲۳ | ۳-۲ راههای شناخت مشاهدات مؤثر |
| ۲۶ | ۴-۲ مثال |
| ۲۹ | ۵-۲ نتیجه گیری |

فصل سوم رگرسیون استوار و انواع برآوردهای آن

| | |
|----|--|
| ۳۱ | ۱-۳ مقدمه |
| ۳۳ | ۲-۳ ویژگی مهم برآوردهای استوار |
| ۳۴ | ۳-۳ برآوردهای استوار |
| ۳۴ | ۱-۳-۳ برآوردهای M |
| ۴۰ | ۲-۳-۳ برآوردهای حداقل میانگین مربعات خطا (LMS) |
| ۴۱ | ۳-۳-۳ برآوردهای حداقل مجموع مربعات پیراسته (LTS) |
| ۴۲ | ۴-۳-۳ برآوردهای S |
| ۴۳ | ۵-۳-۳ برآوردهای MM |
| ۴۳ | ۶-۳-۳ برآوردهای R |
| ۴۴ | ۴-۳ مثالهای کاربردی |
| ۵۷ | ۵-۳ نتیجه گیری |

فصل چهارم رگرسیون حداقل قدرمطلق انحرافات وزنی

| | |
|----|-----------|
| ۵۹ | ۱-۴ مقدمه |
|----|-----------|

- ۶۰ ۲-۴ رگرسیون حداقل قدرمطلق انحرافات (LAD)
- ۶۵ ۳-۴ نقطه فروریزش برآوردگر LAD
- ۶۸ ۴-۴ برآوردگر حداقل قدرمطلق انحرافات وزنی (WLAD)
- ۶۹ ۵-۴ شیوه های برآورد برآوردگر WLAD
- ۷۰ ۶-۴ نقطه فروریزش برآوردگر WLAD در یک نمونه متناهی
- ۷۰ ۷-۴ تعیین وزنهای رگرسیون WLAD
- ۷۵ ۸-۴ ویژگیهای مجانبی برآوردگر WLAD
- ۷۸ ۹-۴ مثال
- ۸۲ ۱۰-۴ نتیجه گیری
- ۸۳ ۱۱-۴ استنباط آماری روی رگرسیون WLAD

پیوستها

- ۸۵ ۱.A اثبات رابطه cook
- ۸۵ ۲.A اثبات رابطه Dffits
- ۸۷ B اثبات مسئله ثانویه
- ۸۹ C تعاریف همگراییها
- ۹۰ D اثبات قضیه ۱

۹۴ E برنامه صحیح مختلط

۹۵ F اثبات قضیه ۲

۹۶ G برنامه بوت استرپ مثال بخش ۴-۱۰

۹۸ منابع

۱۰۲ واژه نامه انگلیسی - فارسی

چکیده :

مقاله اصلی که در این پایان نامه مبنای کار در نظر گرفته شده است تحت عنوان

" Robust Weighted LAD Regression "

است که توسط Simonoff, Sengupta, Giloni در سال ۲۰۰۶ ارائه شده است.

زمانی که در مدل رگرسیون خطی نقاط پرت و دورافتاده وجود داشته باشد، یا مشاهدات از توزیع غیر نرمال تبعیت کنند؛ شیوه حداقل مربعات، دیگر شیوه خوبی برای برآورد پارامترها نیست؛ زیرا این برآوردگر نسبت به مشاهدات غیرمعمول بسیار حساس است. بنابراین شیوه رگرسیون استوار با تعداد زیادی برآوردگر پیشنهاد شده است.

یکی از قدیمی ترین پیشنهادات، شیوه حداقل قدرمطلق انحرافات (LAD) بوده است، که ضرایب رگرسیونی در آن از طریق مینیم کردن مجموع قدرمطلق باقیمانده ها برآورد می شوند.

البته از رگرسیون حداقل قدرمطلق انحرافات به عنوان یک جایگزین قوی برای شیوه حداقل مربعات چشم پوشی شده است، به این دلیل که این برآوردگر به شدت می تواند بوسیله یک تک مشاهده تحت تأثیر قرار گرفته شود، زیرا این برآوردگر دارای نقطه فروریزش برابر $1/n$ است. (n حجم نمونه است.)

هدف اصلی در این پایان نامه این است که نشان دهیم، با انتخاب وزنه های منطقی می توان برآوردگر حداقل قدرمطلق انحرافات وزنی را که دارای نقطه فروریزش بالاتری نسبت به برآوردگر LAD است، را بدست آورد. سپس ویژگی های این برآوردگر را مورد بررسی قرار داده و همچنین کاربرد این برآوردگر استوار را روی داده های واقعی نشان می دهیم.

لازم به ذکر است که ما مقاله ای تحت عنوان " رگرسیون حداقل قدرمطلق انحرافات وزنی استوار" را به مجله علمی- پژوهشی دانشگاه الزهراء (س) ارسال کرده ومنتظر داوری نهایی آن می باشیم، همچنین مقاله ای را تحت همان عنوان در هشتمین کنفرانس آمار ایران در شیراز در تاریخ ۸۵/۵/۲ ارائه دادیم.

برای اولین بار هم مقاله ای را تحت عنوان

" Statistical Inferences on Robust Weighted LAD Regression "

در مورد استنباط آماری رگرسیون حداقل قدرمطلق انحرافات وزنی استوار به یک ژورنال معتبر Statistica Sinica ارسال کرده و منتظر داوری آنها می باشیم.

علائم اختصاری

IRLS: Iteratively Reweighted Least Squares

LAD: Least Absolute Deviation

LMS: Least Median Squares

LTS: Least Trimmed Squares

M-estimator: Maximum likelihood estimator

MAD: median absolute deviations

OLS: Ordinary Least Squares

WLAD: Weighted Least Absolute Deviations

فصل ۱

کلیات

۱-۱ مقدمه

تحلیل رگرسیونی یکی از فنون آماری است که بیشترین کاربرد را در علوم اجتماعی، پزشکی، اقتصادی، کشاورزی و در بسیاری از زمینه های دیگر علمی و علوم کاربردی داراست. تحلیل رگرسیونی شامل مفاهیمی می شود که به سادگی قابل درک هستند و عملاً در اکثر نرم افزارهای آماری قابل اجرا می باشند و برای بدست آوردن ارتباط بین متغیرها کاربرد فراوانی دارد.

در اغلب برازشهای رگرسیونی از شیوه حداقل مربعات معمولی ($OLS^{(1)}$) استفاده می شود، اما زمانی که خطاها دارای توزیع نرمال نباشند و یا مجموعه داده ها شامل داده های پرت باشند؛ روش حداقل مربعات معمولی، دیگر کارا نیست زیرا این شیوه به داده های پرت حساس می باشد و بایستی از روشهای استوار برای برآورد پارامترها استفاده نمود. رگرسیون استوار معمولاً به روشی گفته می شود که نه تنها وقتی خطاها دارای توزیع نرمال است و مشاهدات پرت در مدل حضور ندارند، خوب عمل می کند بلکه نسبت به انحرافات کوچک از فرض نرمال بودن و نسبت به حضور نقاط پرت در مدل نیز حساس نمی باشد. رگرسیون استوار دارای تعداد زیادی برآوردگر است و تکنیکهای آنها مکمل تکنیک کمترین مربعات می باشند به طوریکه وقتی توزیع خطاها نرمال است و مشاهدات پرت در مدل حضور نداشته باشند، جوابهای آنها مشابه با جواب رگرسیون حداقل مربعات است.

یکی از این شیوه ها، شیوه حداقل قدر مطلق انحرافات ($LAD^{(2)}$) است که ضرایب رگرسیونی را از

۱- Ordinary Least Squares

۲- Least Absolute Deviations

از طریق مینیمم کردن مجموع قدرمطلق باقیمانده ها برآورد می کند. این روش تقریباً ۵۰ سال قبل از روش حداقل مربعات، در سال ۱۷۵۷ توسط Boscovich معرفی شد. پس از آن Laplace آنرا ۳۰ سال بعد پذیرفت ولی خیلی زود روش حداقل مربعات به دلیل سادگی محاسبات، بر آن سایه افکند. اما امروزه که برای محاسبات محدودیتی وجود ندارد و از طرف دیگر به دلیل حساسیت رگرسیون $OLS^{(۱)}$ نسبت به حضور نقاط پرت و فرض غیر نرمال خطاها و استواری رگرسیون $LAD^{(۲)}$ در این موارد، موجب شده است که بار دیگر این برآوردگر در زمره برآوردگرهای خوب و استوار قرار بگیرد؛ ولی به این دلیل که این شیوه نسبت به سایر برآوردگرهای استوار، دارای نقطه فروریزش کمتری است، (نقطه فروریزش یک ویژگی مهم برآوردگرهای استوار است که به عنوان ملاکی برای مقایسه آنها استفاده می شود)، از کاربرد آن کاسته است.

به همین دلیل در صدد برآمدیم تا در جهت بالابردن استواری برآوردگر LAD این تحقیق را انجام دهیم. در این پایان نامه، شیوه جدید رگرسیون استوار را به نام رگرسیون حداقل قدرمطلق انحرافات وزنی $(WLAD^{(۳)})$ را ارائه می دهیم که این شیوه دارای نقطه فروریزش بالاتری نسبت به رگرسیون LAD است و نشان می دهیم که با انتخاب وزنه‌های ساده و منطقی می توانیم استواری آنرا افزایش دهیم.

۱- Ordinary Least Squares

۲- Least Absolute Deviations

۳- Weighted Least Absolute Deviations

مقاله اصلی که در اینجا مبنای کار در نظر گرفته شده است تحت عنوان،

" Robust Weighted Least Absolute Deviations Regression "

است که توسط Giloni، Simonoff و Sengupta در سال ۲۰۰۶ برای اولین بار در این زمینه ارائه شده است.

از طرف دیگر چون استنباط آماری در مباحث رگرسیونی دارای اهمیت فراوانی است و تاکنون هیچ مقاله ای درباره استنباط آماری رگرسیون WLAD ارائه نشده است، طبق این ضرورت، ما برای اولین بار مطالعات و تحقیقاتی را در این زمینه انجام داده و مقاله ای را با عنوان

" Statistical Inference on Robust Weighted LAD Regression "

ارائه و به ژورنال معتبر علمی ارسال کردیم.

این بخش را به مروری بر تاریخچه مقالاتی که در این پایان نامه بکاربردیم، اختصاص می دهیم.
Huber در سال ۱۹۷۳ در مقاله ای تحت عنوان " Robust Regression " و در سال ۱۹۸۱ در
مقاله " Robust Statistics " ، Rousseeuw در سال ۱۹۸۱ در کتابی تحت عنوان:

" Robust Regression and Outlier Detection "

و Chatterjee و Hadi در سال ۱۹۸۸ در کتاب

" Sensitivity analysis in linear regression "

تحقیقات وسیعی را در زمینه رگرسیون استوار و ضرورت استفاده از آن ارائه کرده اند.

برای تکمیل تحقیقات رگرسیون $WLAD^{(۱)}$ به تاریخچه رگرسیون حداقل قدرمطلق انحرافات

نیازمندیم. رگرسیون $LAD^{(۲)}$ دارای قدمت بسیار زیادی است، Dielman در سال ۲۰۰۵ در مقاله

ای تحت عنوان " Least absolute value regression " یک مروری بر تحقیقات گسترده

ای که تا آن زمان در زمینه رگرسیون LAD انجام شده بود را ارائه کرده است. او به تاریخچه کامل

رگرسیون LAD ، شیوه های پیشنهادی برای برآورد آن و به محاسبات آنها اشاره کرده است.

Bosovich در سال ۱۷۵۷ برای اولین بار این شیوه را پیشنهاد کرد، در سالهای اخیر الگوریتمهای

بسیار زیادی برای برآورد LAD پیشنهاد شده است که یکی از آنها، الگوریتم سیمپلکس بر اساس

۱- Weighted Least Absolute Deviations

۲- Least Absolute Deviations

مسئله ثانویه است که توسط Barradale و Roberts در سال ۱۹۷۴ بیان شد و Dielman در سال ۱۹۹۲ آنرا اصلاح کرد. Bloomfield و Steiger در سال ۱۹۸۰ الگوریتم دیگری را به نام الگوریتم حداقل مربعات موزون تکراری (IRLS) بیان کردند که Leroy و Rousseeuw در سال ۱۹۸۷، Portnoy و Koenker در سال ۱۹۹۷، Olive و Hawkins در سال ۲۰۰۲، Koenker در سال ۲۰۰۴ به تصحیح شیوه های ذکر شده پرداختند.

همانطور که در مقدمه ذکر شد؛ مقدار فروریزش، یک ویژگی مهم برآوردگر استوار است که Donoho و Huber (۱۹۸۳)، He et al (۱۹۹۰)، Mizera و Muller (۲۰۰۱) شیوه هایی را برای محاسبه مقدار فروریزش رگرسیون LAD پیشنهاد کردند که Giloni و Padberg در سال ۲۰۰۴ در مقاله ای با عنوان

" The finite sample breakdown point of L_1 regression "

شیوه ساده تری را برای مقدار فروریزش رگرسیون LAD در یک نمونه متناهی ارائه کردند. برای افزایش مقدار فروریزش این برآوردگر، Rousseeuw و Zomren (۱۹۹۲) در مقاله ای با نام " A comparison of some quick algorithms " و همچنین Huber و Rousseeuw در سال ۱۹۹۷ در مقاله

" Robust regression with both continues and binary regressors "

طرحهای وزنی را پیشنهاد کردند که اولین طرح به دلیل دقت کم آن و دومی به دلیل محاسبات زیاد آن موجب شدند که Sengupta، Simonoff، Giloni (۲۰۰۶) یک طرح وزنی منطقی را پیشنهاد کنند که باعث افزایش مقدار فروریزش آن می شود و محاسبات ساده تری هم دارد و آنها این برآوردگر جدید را برآوردگر استوار WLAD معرفی کردند.

لازم به ذکر است که با استفاده از مقاله Dielman و Pfaffenberger (۱۹۹۲) با عنوان

"A further comparison of tests of hypotheses in LAV regression "

و مقاله Dielman (۱۹۹۶) تحت عنوان

" A note on hypothesis testing in LAV^(۱) multiple regression "

و مقاله Dielman (۲۰۰۵) با عنوان " Least absolute value regression " که آزمون

فرضهائی را برای رگرسیون LAD پیشنهاد کردند و با استفاده از مقاله Koenker و Bassett

(۱۹۷۸) با عنوان " Regression Quantile " و مقاله Zhou و Portnoy (۱۹۹۸) تحت عنوان "

" Statistical inferences on regression quantile " ، که فاصله اطمینانی را بر مبنای

رگرسیون چندک برای رگرسیون LAD ارائه دادند، توانستیم برای اولین بار فاصله اطمینان و آزمون

فرضی را برای رگرسیون WLAD تحت مقاله ای جدید پیشنهاد کنیم.

۱- Least Absolute Values

۳-۱ ساختار

ساختار پایان نامه به صورت زیر است:

فصل دوم به تعریف مشاهدات مؤثر و انواع آنها، شیوه های شناخت مشاهدات مؤثر همراه با مثال اختصاص دارد.

در فصل سوم به رگرسیون استوار، کاربرد آن و معرفی تعدادی از برآوردهای استوار M ، LMS ^(۱) و LTS ^(۲) می پردازیم و برنامه های نرم افزاری $Splus$ آنها را در مثالهایی ارائه می دهیم.

فصل چهارم به معرفی رگرسیون حداقل قدر مطلق انحرافات، ویژگیهای مجانبی برآوردهای آن، تعریف نقطه فروریزش آن و همچنین به معرفی رگرسیون حداقل قدر مطلق انحرافات وزنی، چگونگی انتخاب وزنها، چگونگی محاسبه نقطه فروریزش و ویژگیهای مجانبی آن اختصاص دارد. کاربرد مطالب این فصل را هم روی یک مثال نشان می دهیم. بخش آخر از این فصل را هم به اساس مقاله ای که برای اولین بار، ما در زمینه فاصله اطمینان و آزمون فرض برای ضرایب رگرسیون $WLAD$ ارائه دادیم، اختصاص می دهیم.

در پایان هر فصل هم نتیجه مربوط به همان فصل را بیان می کنیم.

۱) **Least Median Squares**

۲) **Least Trimmed Squares**

فصل دوم

مشاهدات مؤثر