

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



پایان نامه کارشناسی ارشد در رشته آمار ریاضی

چند روش برآوردیابی استوار در مدل های رگرسیونی چندکی با داده های
پرت

به وسیله‌ی

الهام یوسفی

استاد راهنما

دکتر امین قلم فرسای مستوفی

شهریور ۹۲

به نام خدا

اظہار نامہ

اینجانب الہام یوسفی (۹۰۰۳۷۵) دانشجوی رشته‌ی آمار گرایش آمار ریاضی دانشکده‌ی علوم، اظہار می‌کنم کہ این پایان نامہ حاصل پژوهش خودم بوده و در جاهایی کہ از منابع دیگران استفادہ کردہ ام نشانی دقیق و مشخصات کامل آن را نوشتہ ام. ہم چنین اظہار می‌کنم کہ تحقیق و موضوع پایان نامہ ام تکراری نیست و تعہد می‌نمایم کہ بدون مجوز از دانشگاه دستاوردهای آن را منتشر ننمودہ و یا در اختیار غیر قرار ندم. کلیہ حقوق این اثر مطابق با آیین نامہ مالکیت فکری و معنوی متعلق بہ دانشگاه شیراز است.

نام و نام خانوادگی: الہام یوسفی

تاریخ و امضاء: ۱۳۹۲/۶/۱۹

به نام خدا

چند روش برآوردیابی استوار در مدل های رگرسیونی چندکی با داده های

پرت

به وسیلهی

الهام یوسفی

پایان نامه

ارائه شده به تحصیلات تکمیلی دانشگاه شیراز به عنوان به عنوان بخشی از فعالیت-
های لازم برای اخذ درجهی کارشناسی ارشد

در رشتهی:

آمار ریاضی

از دانشگاه شیراز

شیراز

جمهوری اسلامی ایران

ارزیابی کمیتهی پایان نامه، با درجه : عالی

دکتر امین قلم فرسای مستوفی، استادیار بخش آمار (رئیس کمیته).....
دکتر احمدرضا سلطانی، استاد بخش آمار.....
دکتر علیرضا نعمت اللهی، دانشیار بخش آمار.....
دکتر عبدالرسول برهانی حقیقی، استادیار بخش آمار.....

شهریور ۱۳۹۲

تقدیم به دو مهربان زندگیم؛

به پدر عزیزم که عظمت هستی را در وجودش یافتم

به مادر مهربان و صبورم

که از هیچ کوششی برای سعادت من دریغ نکردند

سپاس گزاری

نخست، سپاس بی کران پروردگار یکتا را که هستیمان بخشید... بعد از حمد و سپاس به درگاه خداوند منان از نعمت‌های بیکرانی که بر ما ارزانی داشته است، از پدر و مادر صبور و مهربانم به خاطر تمام حمایت‌ها، تلاش‌ها و فداکاری‌هایشان که در تمام مراحل زندگی همراه من بوده اند سپاسگزارم. سپس بر خود لازم می‌دانم از زحمات استاد گرانقدرم جناب آقای دکتر قلم فرسا که پیوسته با راهنمایی‌هایشان راه گشای مشکلاتم در این مسیر بوده و از هیچ کوششی فروگذار ننمودند کمال تشکر و سپاس را داشته باشم. هم چنین از زحمات جناب آقای دکتر نعمت‌اللهی که با کمال صبر و حوصله در قسمت‌های مختلف پایان نامه در راه گشایی مشکلاتم کمک شایانی به من کرده و نکات ارزنده‌ای را به من آموختند، بی نهایت سپاس گزارم. از اساتید محترم جناب آقایان دکتر برهانی و دکتر سلطانی که با کمال دقت، زحمت مطالعه‌ی این پایان نامه را تقبل کرده و نکات ارزنده‌ای را به من گوشزد کردند، کمال تشکر و قدردانی را دارم. در انتها از امیدواری‌ها، همیاری‌ها و مهربانی‌های دوست همیشه همراهم... صحرا فرامرزی که پیوسته و همه جا همراهم بوده بسیار سپاس گزارم و امیدوارم روزی بتوانم قسمتی از مهربانی‌هایش را جبران کنم.

چکیده

چند روش برآوردیابی استوار در مدل های رگرسیونی چندکی با داده های پرت

به کوشش

الهام یوسفی

در آمار و روش های آماری در اقتصاد، اغلب با مجموعه داده هایی روبرو هستیم که شامل نقاط پرت هستند. رگرسیون چندکی خطی نیز به عنوان روشی پرکاربرد در آمار، نسبت به مشاهدات پرت، مخصوصا مشاهدات پرت موجود در متغیرهای مستقل، حساس است. در این پایان نامه ابتدا چند برآوردگر رگرسیونی را معرفی کرده و نشان می دهیم که چگونه می توان با استفاده از نقطه شکستگی استواری آنها را ارزیابی کرد. سپس درباره ی برآوردگرهای رگرسیون چندکی بحث می کنیم. به دلیل اینکه چنین برآوردگری در حضور نقاط پرت استوار نیست، برآوردگر استواری را که توسط سیزک و همکارانش در سال ۲۰۱۲ معرفی شد، مورد بررسی قرار می دهیم. در نهایت عملکرد رگرسیون چندکی را با رگرسیون چندکی حداقل پیراسته از طریق یک مثال عددی با هم مقایسه می کنیم.

واژه های کلیدی: برآوردگر استوار، نقطه شکستگی، رگرسیون چندکی، پیراستگی، نقطه

اهرمی

نام : الهام

نام خانوادگی: یوسفی

مقطع تحصیلی: کارشناسی ارشد

رشته و گرایش: آمار ریاضی

استاد راهنما: دکتر امین قلم فرسا

تاریخ دفاع: ۱۳۹۲/۶/۱۹

چکیده

چند روش برآوردیابی استوار در مدل های رگرسیونی چندکی با داده های پرت

در آمار و روش‌های آماری در اقتصاد، اغلب با مجموعه داده‌هایی روبرو هستیم که شامل نقاط پرت هستند. رگرسیون چندکی خطی نیز به عنوان روشی پرکاربرد در آمار، نسبت به مشاهدات پرت، مخصوصاً مشاهدات پرت موجود در متغیرهای مستقل، حساس است. در این پایان نامه ابتدا چند برآوردگر رگرسیونی را معرفی کرده و نشان می‌دهیم که چگونه می‌توان با استفاده از نقطه شکستگی استواری آنها را ارزیابی کرد. سپس درباره‌ی برآوردگرهای رگرسیون چندکی بحث می‌کنیم. به دلیل اینکه چنین برآوردگری در حضور نقاط پرت استوار نیست، برآوردگر استواری را که توسط سیزک و همکارانش در سال ۲۰۱۲ معرفی شد، مورد بررسی قرار می‌دهیم. در نهایت عملکرد رگرسیون چندکی را با رگرسیون چندکی حداقل پیراسته از طریق یک مثال عددی با هم مقایسه می‌کنیم.

فهرست مطالب

عنوان صفحه

فصل اول: مقدمه

- ۱-۱- مقدمه ۲
- ۱-۱-۱- اهداف کلی رگرسیون و ایده‌ی معرفی رگرسیون چندکی ۲
- ۱-۱-۲- اهداف ایجاد روش های رگرسیونی استوار ۳
- ۲-۱- مفاهیم اولیه ۷
- ۱-۲-۱- برآوردگر استوار ۷
- ۲-۲-۱- نمونه‌ی آلوده شده ۱۱
- ۳-۲-۱- وضعیت کلی مشاهدات ۱۱
- ۴-۲-۱- نقطه اهرمی ۱۲
- ۵-۲-۱- برآوردگر پایای رگرسیون ۱۲
- ۶-۲-۱- چند تعریف آنالیزی مورد نیاز ۱۲

فصل دوم: معرفی چند برآوردگر رگرسیونی و بررسی استواری آن ها

- ۱-۲- برآوردگر حداقل مربعات ۱۶
- ۲-۲- برآوردگر حداقل میانه مربعات ۱۸
- ۳-۲- روش برآوردیابی S ۲۳
- ۴-۲- برآوردگر حداقل وزن مربعات ۲۶
- ۵-۲- برآورد استوار و کارای روش حداقل مربعات وزنی ۲۹

- ۳۴ ۲-۵-۱- استواری برآوردگر REWLS
- ۳۷ ۲-۶- برآوردگر حداقل وزن مربعات دو مرحله ای (2S-LWS)
- ۴۰ ۲-۶-۱- نقطه شکستگی برآوردگر 2S-LWS

فصل سوم: معرفی رگرسیون چندکی

- ۴۴ ۳-۱- چندک ها و دست یابی به آن ها از طریق بهینه سازی
- ۴۷ ۳-۲- چندک های نمونه ای
- ۴۸ ۳-۳- دیدگاه کلی درباره ی رگرسیون چندکی
- ۴۹ ۳-۳-۱- برآوردگر رگرسیون چندکی

فصل چهارم: معرفی برآوردگر استوار رگرسیون چندکی حداقل

- ۵۴ ۴-۱- برآوردگر پیراسته ی تعمیم یافته
- ۵۸ ۴-۲- برآوردگر رگرسیون چندکی حداقل پیراسته
- ۵۹ ۴-۲-۱- نقطه شکستگی برآوردگر LTQR
- ۶۷ ۴-۲-۲- سازگاری برآوردگر LTQR

فصل پنجم: شبیه سازی و نتیجه گیری

- ۷۰ ۵-۱- خواص برآوردگرهای معرفی شده در فصل دوم، در حالت نمونه های متناهی
- ۷۲ ۵-۱-۱- کارایی برای نمونه های متناهی
- ۷۵ ۵-۱-۲- رفتار برای داده های غیر هم توزیع
- ۷۷ ۵-۲- خواص برآوردگرهای رگرسیون چندکی و رگرسیون چندکی حداقل پیراسته
- ۷۸ ۵-۲-۱- مجموعه داده های خوشه ی اختری CYB OB1
- ۸۱ ۵-۲-۲- آزمایش های شبیه سازی

منابع ۹۲

پیوست الف ۹۸

پیوست ب ۱۳۳

فهرست جدول‌ها

عنوان و شماره	صفحه
جدول ۱-۳- برآوردهای رگرسیون چندکی برای داده‌های engel	۵۲.....
جدول ۱-۵- کارایی نسبی MSE برای خطاهای نرمال، $\varepsilon_i \sim N(0,1)$	۷۳.....
جدول ۲-۵- کارایی نسبی MSE برای خطاهای دارای توزیع t	۷۴.....
جدول ۳-۵- متوسط مربع خطا (MSE) برای مدل‌های رگرسیونی مختلف	۷۶.....

فهرست شکل‌ها

صفحه	عنوان و شماره
۲۲.....	شکل ۱-۲- رسم برآوردگرهای LS و LMS برای داده‌های stars
۴۵.....	شکل ۱-۳- تابع زیان $\rho_{\tau}(u)$ در رابطه‌ی (۲-۳)
۵۱.....	شکل ۲-۳- نمودار برازش خطوط رگرسیون چندکی برای داده‌های engel
۷۸.....	شکل ۱-۵- برآوردگر رگرسیون چندکی برای داده‌های stars
۸۰.....	شکل ۲-۵- برآوردگر LTQR برای داده‌های stars
	شکل ۳-۵- نمودارهای جعبه‌ای حاصل از برآورد ضرایب مدل (۳-۵) در داده‌های
۸۲.....	حقیقی (مقدار آرایش ۰٪) و توزیع یکنواخت برای متغیرهای مستقل
	شکل ۴-۵- نمودارهای جعبه‌ای حاصل از برآورد ضرایب مدل (۳-۵) در داده‌های
۸۳.....	حقیقی (مقدار آرایش ۰٪) و توزیع نرمال برای متغیرهای مستقل
	شکل ۵-۵- نمودارهای جعبه‌ای حاصل از برآورد ضرایب مدل (۳-۵) در آزمایش
۸۵.....	شبیه‌سازی اول با ۱۰٪ آرایش و توزیع یکنواخت برای متغیرهای مستقل
	شکل ۶-۵- نمودارهای جعبه‌ای حاصل از برآورد ضرایب مدل (۳-۵) در آزمایش
۸۶.....	شبیه‌سازی دوم با ۱۰٪ آرایش و توزیع نرمال برای متغیرهای مستقل

- شکل ۵-۷- نمودارهای جعبه ای حاصل از برآورد ضرایب مدل (۳-۵) در آزمایش شبیه سازی اول با ۲۰٪ آلاینش و توزیع یکنواخت برای متغیرهای مستقل ۸۷
- شکل ۵-۸- نمودارهای جعبه ای حاصل از برآورد ضرایب مدل (۳-۵) در آزمایش شبیه سازی دوم با ۲۰٪ آلاینش و توزیع نرمال برای متغیرهای مستقل ۸۸
- شکل ۵-۹- نمودارهای جعبه ای حاصل از برآورد ضرایب مدل (۳-۵) در آزمایش شبیه سازی اول با ۳۰٪ آلاینش و توزیع یکنواخت برای متغیرهای مستقل ۸۹
- شکل ۵-۱۰- نمودارهای جعبه ای حاصل از برآورد ضرایب مدل (۳-۵) در آزمایش شبیه سازی دوم با ۳۰٪ آلاینش و توزیع نرمال برای متغیرهای مستقل ۹۰
- شکل ب-۱- تشریح مکان زیرفضاهای S و V در فضای E در حالت اول ۱۳۴
- شکل ب-۲- تشریح یک ساختار هندسی در اثبات لم در حالت دوم ۱۳۷

فصل اول

مقدمه

۱-۱- مقدمه

۱-۱-۱- اهداف کلی رگرسیون و ایده‌ی معرفی رگرسیون چندکی

یکی از اهداف مهم بررسی‌های آماری، یافتن روابطی است که به کمک آن بتوان اثر تغییرات یک یا چند متغیر را بر روی متغیرهای دیگر پیش بینی کرد. رگرسیون در واقع روشی برای استخراج روابط موجود بین متغیرها ارائه می‌دهد. رگرسیون چندکی نیز به عنوان زیر شاخه‌ای از رگرسیون به دنبال این هدف است که راهکار کاملی را برای تکمیل عملکرد رگرسیون میانگین ارائه دهد.

بسیاری از مصادیق آمار کاربردی را می‌توان مثلاً در پیچیدگی مدل رگرسیون خطی و روش‌های برآوردیابی مرتبط با حداقل مربعات دید. در این باره Mosteller and Tukey (1977, p. 266) گفته‌اند:

آنچه که خط یا منحنی رگرسیونی انجام می‌دهد، بررسی متغیر پاسخ از طریق میانگین مشاهدات با توجه به مجموعه‌ی x ها می‌باشد. همچنین می‌توانیم خط یا منحنی‌های رگرسیونی مختلف را با توجه به درصد نقاط توزیع‌ها محاسبه کنیم و بنابراین تصویر کامل‌تری از نقاط بدست بیاوریم. معمولاً این کار به طور متداول انجام نمی‌شود و بنابراین رگرسیون تصویر ناکاملی را ارائه می‌دهد. همانطور که میانگین تصویر ناکاملی را از یک توزیع ارائه می‌دهد، بنابراین منحنی رگرسیون نیز در این باره تصویر ناکاملی را از یک مجموعه توزیع‌ها ارائه می‌دهد.

هدف رگرسیون چندکی این است که راهکار کاملی را برای تکمیل عملکرد رگرسیون میانگین ارائه دهد.

چرا برآورد حداقل مربعات در مدل رگرسیونی خطی اینقدر در آمار کاربردی نفوذ کرد؟ چه چیزی آن را به ابزار قدرتمندی تبدیل کرد؟ ایده‌های اولیه درباره‌ی انجام رگرسیون به کارهای اولیه‌ی **Bosovich** در نیمه‌های قرن ۱۸م تا کارهای **Edgeworth** در پایان قرن ۱۹م برمی‌گردد. در ابتدا باید گفت که یک حقیقت مهم این است که آسانی محاسبه‌ی این برآوردگرها بسیار خوشایند است. به طور قطع این از انگیزه‌های اولیه برای موفقیت این روش بود. پس اگر خطاها در مدل به طور نرمال توزیع شده باشند، روش حداقل مربعات دارای مقدار بهینه (ماکزیمم یا مینیمم) می‌شود. همانطور که **Mosteller and Tukey** اظهار کرده اند، به ندرت دیده می‌شود که میانگین به تنهایی برای آنالیز آماری یک تک نمونه هدف رضایت بخشی باشد. ابزارهای پراکندگی مانند چولگی، برجستگی، نمودارهای جعبه‌ای، هیستوگرام ها و روش‌های برآوردیابی چگالی همگی برای اینکه به بینش بیشتری دست یابیم به کار گرفته می‌شود. آیا می‌توان چنین کاری را در رگرسیون انجام داد؟ پس می‌توان گفت که ایده‌ی ایجاد رگرسیون چندکی می‌تواند این باشد که به جای برآورد سطح میانگین شرطی (که شامل خط، صفحه و ... است) از طریق روش حداقل مربعات، سطح‌های مختلف چندک‌های شرطی را برآورد کنیم که باعث تکمیل روش قبل می‌شود.

۱-۱-۲- اهداف ایجاد روش‌های رگرسیونی استوار

اخیرا در آمار و روش‌های آماری در اقتصاد بیشتر به تکنیک‌هایی توجه شده است که با مشاهداتی سروکار دارد که غیرمعمول هستند. که این امر از توزیع‌هایی با دنباله پهن، کدگذاری اشتباه یا ناهمگنی یافته شده در مدل ناشی می‌شود و دارای اهمیت خیلی زیاد، مخصوصا در مدل‌های رگرسیونی خطی (یا غیرخطی) و سری‌های زمانی است. هدف این است که در رگرسیون، برآوردگرهایی ارائه دهیم که تحت تاثیر نقاط پرت یا انحراف از فرضیات مدل نیستند. این برآوردگرها، برآوردگرهای استوار^۱ نامیده می‌شوند. برای اثبات استواری برآوردگرهای معرفی شده از معیاری به نام نقطه شکستگی^۲ استفاده می‌کنیم که در ادامه

^۱ Robust estimators

^۲ Breakdown point

معرفی می‌شود (Rousseeuw (1997), Genton and Lucas (2003), Davies and Gather (2005). اولین برآوردگر رگرسیونی استوار با نقطه شکستگی مجانبی $1/2$ برآوردگر حداقل میانه‌ی مربعات بود (LMS; Rousseeuw, 1984). برآوردگرهای حداقل مربعات پیراسته (LTS; Rousseeuw, 1985) و برآوردگرهای S، Rousseeuw and Yohai, (1984) نمی‌توانند همزمان دارای نقطه شکستگی بالا و کارایی نسبی بالا باشند، Hössjer, (1992). از سوی دیگر بسیاری از روش‌های رگرسیونی استوار کلاسیکی، ظاهراً مشکل استوار بودن را، مخصوصاً برای داده‌هایی با توزیع نرمال، با کارایی نسبی پایین جبران می‌کنند. برای بهبود بیشتر کیفیت روش‌های با نقطه شکستگی بالا (Rousseeuw and Leroy (1987) پیشنهاد استفاده از روش حداقل وزن مربعات (LWS)³ را دادند که در آن به مشاهداتی که بیشتر از یک نقطه‌ی قطع⁴ ثابت (نقطه‌ای که مشاهدات بزرگتر از آن را در یک نمونه، پرت به حساب می‌آوریم) هستند، وزن صفر داده می‌شود. این کار تغییر پذیری برآورد را نسبت به برآورد استوار اولیه کم می‌کند؛ اما توزیع مجانبی آن به برآورد استوار اولیه بستگی پیدا خواهد کرد (Welch and Ronchetti, 2002). بنابراین (Gervini and Yohai (2002) پیشنهاد استفاده از حداقل مربعات وزنی با استفاده از نقطه قطع وابسته به داده‌ها را دادند که منجر به حاصل شدن برآورد استوار و کارایی حداقل مربعات وزنی (REWLS)⁵ شد. در این حالت باز هم توزیع مجانبی REWLS بستگی به برآورد اولیه دارد و برخی از مشاهدات از برآورد حذف می‌شوند. برای برطرف شدن این مشکلات کلاسی از روش‌های برآوردیابی دو مرحله‌ای را ارائه می‌دهیم که روش حداقل وزن مربعات دو مرحله‌ای (2S-LWS)⁶ نامیده می‌شود. در این روش به جای استفاده از وزن صفر برای مشاهدات پرت و به کارگیری روش حداقل مربعات وزنی (WLS)⁷، از وزن‌های کاملاً مثبت استفاده می‌شود و حداقل وزن مربعات LWS را به جای WLS به کار می‌بریم. که منجر به حاصل شدن 2S-LWS و خواص مطلوب آن می‌شود.

گاهی اوقات رگرسیون میانگین برای برازش به داده‌ها مناسب نیست و اطلاعات کاملی از نحوه‌ی پراکندگی داده‌ها ارائه نمی‌دهد. در این گونه موارد نیاز است از رگرسیون چندکی

³ Least Weighted Squares

⁴ cut off point

⁵ Robust and Efficient Weighted Least Squares

⁶ 2-Step Least Weighted Squares

⁷ Weighted least squares

(QR)⁸ به جای رگرسیون میانگین استفاده کنیم. رگرسیون چندکی توسط Koenker and Basset (1978) معرفی شد. برآوردگر رگرسیون چندکی نسبت به دم‌های چوله و اختلافات از توزیع نرمال، استوار است. به علاوه تحت شرایط خیلی کلی توزیع مجانبی بردار ضرایب برآورد شده نرمال چند گانه است که باعث ایجاد استنباط‌های استاندارد می‌شود (Koenker and Bassett, 1978). اما این برآوردگر نسبت به حضور نقاط پرت در متغیرهای مستقل حساس است (He et al, 1990). بنابراین در این گونه موارد نیاز است برآوردگرهای استوار دیگری جایگزین رگرسیون چندکی شود. بنابراین یک رویکرد پیشنهادی را برای برآوردیابی استوار در چارچوب رگرسیون چندکی بررسی می‌کنیم که رگرسیون چندکی حداقل پیراسته (LTQR)⁹ نام دارد. روش کار این برآوردگر بر اساس پیراسته کردن و حذف مشاهدات است که اثر نقاط پرت در متغیرهای مستقل را کاهش می‌دهد. روش ارائه شده برآوردگر استوار مکانی میانه را که توسط Tableman (1994 a,b) و برآوردگر حداقل انحراف مطلق پیراسته (LTA)¹⁰ که توسط Hawkins and Olive (1999) بررسی شد، گسترش می‌دهد. برآورد LTA و LTQR چندک‌های رگرسیونی را برای داده‌هایی که از تابع هدف قطع نشده‌اند، یعنی برای زیرمجموعه‌ای از داده‌ها، برآورد می‌کنند و باعث ایجاد خواص استواری می‌شوند.

در فصل اول این پایان نامه ابتدا ایده و انگیزه‌ی ایجاد رگرسیون چندکی که یکی از شاخه‌های پرکاربرد رگرسیون است را عنوان می‌کنیم و دلیل کاملتر بودن آن نسبت به رگرسیون میانگین را بحث می‌کنیم. سپس مفاهیم پایه برای استفاده در فصل‌های بعد را خواهیم آورد. در فصل دوم چند برآوردگر پرکاربرد رگرسیون میانگین را معرفی کرده و آنها را از نظر استواری مورد بررسی قرار می‌دهیم. برای بررسی استواری از معیار نقطه شکستگی استفاده می‌کنیم. در فصل سوم رگرسیون چندکی را معرفی کرده و برآوردگر پرکاربرد آن را ارائه می‌دهیم. به دلیل ضعف برآوردگر رگرسیون چندکی در حضور نقاط پرت موجود در متغیرهای مستقل، در فصل چهارم، رگرسیون چندکی حداقل پیراسته را معرفی می‌کنیم و استواری آن را از نظر معیار نقطه شکستگی مورد بررسی قرار می‌دهیم. ایده معرفی این برآوردگر بر اساس پیراسته سازی و کوتاه سازی چندک‌های رگرسیونی است. در انتها، در فصل پنجم، ابتدا کارایی نسبی برآوردگرهای معرفی شده در فصل دوم را در حالت نمونه‌های متناهی محاسبه کرده و متوسط مربع خطا (MSE) را برای آنها تحت مدل‌های توزیعی مختلف با هم مقایسه می‌کنیم. سپس

⁸ Quantile Regression

⁹ Least Trimmed Quantile Regression

¹⁰ Least Trimmed Absolute deviation

عملکرد برآوردگرهای رگرسیون چندکی و رگرسیون چندکی حداقل پیراسته را هم از طریق یک مجموعه داده واقعی و هم از طریق داده های شبیه سازی شده با هم مقایسه می کنیم.