

۱۳۷۷ / ۱۲ / ۱۹



۲۵۹۷

۲۵۹۷

۲۵۹۷

بیت‌الاحرام  
بیت‌الاحرام

۲۴۱۱۷



دانشگاه تربیت مدرس

دانشکده علوم پایه

پایان نامه دوره کارشناسی ارشد آمار

# مقایسه اسپلاین با رگرسیون چند جمله‌ای و کاربرد آن

توسط

روشنک علی محمدی

استاد راهنما

دکتر محسن محمدزاده

۱۳۱۹۷/۲

آذر ۱۳۷۷






۲۴۱۱۷

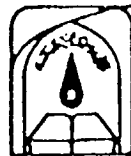
## تأییدیه اعضای هیأت داوران حاضر در جلسه دفاع از پایان نامه کارشناسی ارشد

اعضای هیئت داوران نسخه نهایی پایان نامه خانم/ آقای روشنگر علی محمدی

تحت عنوان: مقایسه اسپیلاین با رگرسیون چند جمله‌ای و کاربرد آن

را از نظر فرم و محتوی بررسی نموده و پذیرش آنرا برای تکمیل درجه کارشناسی ارشد پیشنهاد می‌کنند.

اعضای هیأت داوران	نام و نام خانوادگی	رتبه علمی	امضاء
۱- استاد راهنما	آقای دکتر محمدزاده	استادیار	
۲- نماینده شورای تحصیلات تکمیلی	آقای دکتر محمدزاده	استادیار	
۳- استاد ناظر	آقای دکتر اسماعیل خرم	استادیار	
۴- استاد ناظر	آقای دکتر عین‌الله پاشا	دانشیار	
۵- استاد ناظر	آقای دکتر عباس گرامی	استادیار	



شماره.....

تاریخ.....

پیوست.....

## آیین نامه چاپ پایان نامه (رساله) های دانشجویان دانشگاه تربیت مدرس

نظر به اینکه چاپ و انتشار پایان نامه (رساله) های تحصیلی دانشجویان دانشگاه تربیت مدرس مبین بخشی از فعالیتهای علمی - پژوهشی دانشگاه است بنابراین به منظور آگاهی و رعایت حقوق دانشگاه، دانش آموختگان این دانشگاه نسبت به رعایت موارد ذیل متعهد می شوند:

ماده ۱ در صورت اقدام به چاپ پایان نامه (رساله) ی خود، مراتب را قبلاً به طور کتبی به مرکز نشر دانشگاه اطلاع دهد.

ماده ۲ در صفحه سوم کتاب (پس از برگ شناسنامه)، عبارت ذیل را چاپ کند:  
کتاب حاضر، حاصل پایان نامه کارشناسی ارشد / رساله دکتری نگارنده در رشته ..... است  
که در سال ..... در دانشکده ..... دانشگاه تربیت مدرس به راهنمایی سرکار خانم / جناب آقای دکتر ..... و مشاوره سرکار خانم / جناب آقای دکتر ..... از آن دفاع شده است.

ماده ۳ به منظور جبران بخشی از هزینه های نشریات دانشگاه تعداد یک درصد شمارگان کتاب (در هر نوبت چاپ) را به مرکز نشر دانشگاه اهدا کند دانشگاه می تواند مازاد نیاز خود را به نفع مرکز نشر در معرض فروش قرار دهد.

ماده ۴ در صورت عدم رعایت ماده ۳، ۵۰٪ بهای شمارگان چاپ شده را به عنوان خسارت به دانشگاه تربیت مدرس، تأدیه کند.

ماده ۵ دانشجوی تعهد و قبول می کند در صورت خودداری از پرداخت بهای خسارت، دانشگاه می تواند خسارت مذکور را از طریق مراجع قضایی مطالبه و وصول کند؛ به علاوه به دانشگاه حق می دهد به منظور استیفای حقوق خود، از طریق دادگاه، معادل وجه مذکور در ماده ۴ را از محل توقیف کتابهای عرضه شده نگارنده برای فروش، تأمین نماید.

ماده ۶ اینجناب روستک علی محمدی دانشجوی رشته آمار مقطع کارشناسی ارشد تعهد فوق و ضمانت اجرایی آن را قبول کرده، به آن ملتزم می شوم.

## قدردانی

در این جا لازم است از کلیه افرادی که مرا در انجام این پایان نامه کمک نموده‌اند، خصوصاً استاد گرامی جناب آقای دکتر محسن محمدزاده که در تمام مراحل انجام کار با مساعدت‌ها و راهنمایی‌های بی‌دریغ خود مرا یاری کردند، تشکر کنم.

همچنین از اساتید محترم آقایان دکتر پاشا، دکتر گرامی و دکتر خزّم که حق تعلیم بر اینجانب دارند و به سبب حضور در جمع داوران سپاسگزاری و قدردانی می‌نمایم.

سلامت، سعادت و موفقیت روزافزون اساتید ارجمند را از درگاه خداوند متعال خواستارم.

## مقایسه اسپلاین با رگرسیون چندجمله‌ای و کاربرد آن

### چکیده

رگرسیون چندجمله‌ای روشی پارامتری برای برازش منحنی به داده‌ها می‌باشد که در آن ارتباط بین متغیر پاسخ و متغیرهای مستقل به صورت چندجمله‌ای برآورد می‌گردد. در بسیاری از مواقع در مورد نحوه ارتباط بین متغیرها اطلاع زیادی در دست نیست. در این صورت بهتر است به جای مفروض داشتن یک الگوی پارامتری خاص (مانند چندجمله‌ای) برای داده‌ها از روشی استفاده شود که داده‌ها ماهیت روند خود را بهتر نشان دهند. اسپلاین همواری روشی برای برآورد منحنی است که در آن در مورد شکل منحنی فرضهای قوی اعمال نمی‌شود و تنها فرض همواری منحنی در نظر گرفته می‌شود. در این پایان‌نامه برازش منحنی به دو روش رگرسیون چندجمله‌ای و اسپلاین همواری صورت گرفته و نهایتاً میزان دقت این دو روش بر اساس معیار مجموع مربعات باقیمانده‌ها مورد مقایسه عددی قرار گرفته‌اند. این مطالعه نشان می‌دهد که برای نمونه‌های باحجم کوچک و برای هر مقدار انحراف معیار، اسپلاین همواری منحنی بهتری از رگرسیون چندجمله‌ای به داده‌ها برازش می‌دهد، اما برای نمونه‌های باحجم بزرگ رگرسیون چندجمله‌ای روشی آسانتر و سریعتر است.

واژه‌های کلیدی: اسپلاین همواری، رگرسیون چندجمله‌ای و اعتبار متقابل تعمیم‌یافته

# فهرست مندرجات

۱	مقدمه	۱
۱	۱.۱ آنالیز رگرسیون .....	۱
۳	۲.۱ هموارسازی .....	۳
۸	۲ رگرسیون چند جمله‌ای	۸
۸	۱.۲ رگرسیون خطی .....	۸
۹	۱.۱.۲ روش کمترین مربعات .....	۹
۱۱	۲.۲ رگرسیون چند جمله‌ای .....	۱۱
۱۲	۱.۲.۲ مدل‌های چند جمله‌ای یک متغیره .....	۱۲
۱۵	۲.۲.۲ مدل‌های چند جمله‌ای دو یا چند متغیره .....	۱۵

۱۸ ..... چند جمله‌ای‌های متعامد ۳.۲

۲۵ ..... هموارسازها و اسپلاین ۳

۲۵ ..... مقدمه ۱.۳

۲۶ ..... هموارسازی ۲.۳

۲۶ ..... فنون هموارسازی ۱.۲.۳

۳۰ ..... رهیافت جریمه ناهمواری ۲.۲.۳

۳۳ ..... اسپلاین رگرسیونی ۳.۲.۳

۳۴ ..... اسپلاین درجه سه ۴.۲.۳

۳۵ ..... اسپلاین درجه سه طبیعی ۵.۲.۳

۳۵ ..... محاسبه  $g$  ۶.۲.۳

۳۸ ..... الگوریتم رینش ۷.۲.۳

۴۰ ..... هموارسازها برای بیش از یک متغیر مستقل ۸.۲.۳

۴۲ ..... تقابل اریبی-واریانس ۹.۲.۳

۴۳ ..... انتخاب پارامتر همواری ۱۰.۲.۳

۴۴ ..... روش اعتبار متقابل ۱۱.۲.۳

۴۵ ..... تقابل اریبی-واریانس برای هموارسازهای خطی ۱۲.۲.۳

۴۵ ..... اعتبار متقابل برای هموارسازهای خطی ۱۳.۲.۳

۴۷ ..... اعتبار متقابل تعمیم‌یافته ( $GCV$ ) ۱۴.۲.۳

۴۷ ..... درجه آزادی هموارساز ۱۵.۲.۳



۴۸ ..... ۳.۳ برآزش مدل به داده‌ها

۵۰ ..... ۴ مقایسه اسپلاین و رگرسیون چندجمله‌ای

۵۰ ..... ۱.۴ مقدمه

۵۰ ..... ۲.۴ روش محاسبه و کاربرد

۵۳ ..... ۳.۴ برآزش رگرسیون چندجمله‌ای

۵۳ ..... ۴.۴ شیوه شبیه‌سازی

۵۴ ..... ۵.۴ نتایج شبیه‌سازی

۶۴ ..... ۶.۴ بحث و نتیجه‌گیری

## فصل ۱

### مقدمه

در بسیاری از مسائل کاربردی، تعیین ارتباط بین عوامل مختلف از اهمیت خاصی برخوردار است. در آمار روشهای مختلفی برای تعیین نوع ارتباط و ارائه مدل ارتباطی بین عوامل وجود دارد. یکی از این فنون، آنالیز رگرسیون و دیگری روشهای هموارسازی است که در این فصل مختصراً مرور خواهد شد.

#### ۱.۱ آنالیز رگرسیون

آنالیز رگرسیون یک فن آماری برای تعیین نحوه ارتباط بین دو یا چند متغیر است به طوری که بتوان مقدار یک متغیر  $Y$  را با استفاده از مقادیر متغیرهای  $X_1, \dots, X_p$  برآورد یا پیش‌بینی کرد. بدلیل آنکه مقدار متغیر  $Y$  تحت تاثیر مقادیر متغیرهای  $X_1, \dots, X_p$  می‌باشد،  $Y$  متغیر وابسته (پاسخ) و  $X_1, \dots, X_p$  متغیرهای مستقل (پیشگو) نامیده می‌شوند.

معمولاً منظور از آنالیز رگرسیون، برازش یک مدل ریاضی به مجموعه داده‌ها به عنوان الگوی ارتباطی بین متغیرها است. به عبارت دیگر ابتدا خانواده‌ای از توابع در نظر گرفته شده و

سپس تعیین یک عضو خانواده مورد نظر، به عنوان مدل ارتباط بین متغیرها تعیین می‌شود. پس از تعیین مدل، با بررسی نمودار باقیمانده‌ها و برآورد، آزمون فرض در مورد پارامترهای مدل و نیکویی برازش مدل مورد تحقیق قرار می‌گیرد. این شیوه تحلیل رگرسیون، رگرسیون پارامتری نامیده می‌شود. جنبه دیگر رگرسیون، رگرسیون ناپارامتری است. ابتدا حالتی را که تنها یک متغیر مستقل (پیشگو) و یک متغیر وابسته (پاسخ) موجود است، در نظر می‌گیریم.

فرض کنید که مشاهدات متغیر تصادفی  $Y$  در  $n$  مقدار  $X$  به دست آمده است. می‌توان

مدل

$$y_i = g(x_i) + e_i \quad ; i = 1, \dots, n \quad (1.1.1)$$

را در نظر گرفت که در آن  $e_i$  ها متغیرهای تصادفی مستقل از توزیع  $N(0, \sigma^2)$  و  $g$  تابعی نامعلوم است که تابع یا منحنی رگرسیون نامیده می‌شود.

بکارگیری رگرسیون پارامتری یا ناپارامتری در مدل (۱.۱.۱) بستگی به فرضیهایی دارد که درباره تابع  $g$  می‌شود. در رگرسیون پارامتری شکل تابع  $g$  معلوم فرض می‌شود. به عبارت دیگر، خانواده‌ای از توابع را در نظر گرفته و با استفاده از فنون تحلیل رگرسیون، تابع مناسب برای برآورد منحنی تعیین می‌شود.

مدلهای پارامتری می‌توانند به طور خطی یا غیرخطی به پارامترها وابسته باشند. برای مدل‌هایی که نسبت به پارامترها خطی‌اند، می‌توان با استفاده از یک روش مناسب برآورد مانند حداقل مربعات خطا یا حداکثر درست‌نمایی، پارامترها را برآورد نموده و برآورد منحنی  $g$  را بدست آورد. این مدلها شامل رگرسیون چندجمله‌ای نیز می‌باشد. بعضی از مدل‌هایی که نسبت به پارامترها غیرخطی‌اند ممکن است با کمک یک تبدیل، خطی شوند. منحنی‌های لگاریتمی و توانی مثال‌هایی از این نوع هستند. تحلیل این نوع از مدلها را می‌توان مانند مدل‌های خطی انجام داد.

آنالیز آن دسته از مدلها که حتی توسط تبدیلات، خطی نمی‌شوند، با بکارگیری فنون

رگرسیون غیرخطی ممکن می‌گردد. برای مطالعه در این زمینه می‌توان به گلانت (۱۹۷۵)، دراپر و اسمیت (۱۹۶۶)، بیتس و واتس (۱۹۸۸)، مایرز (۱۹۹۰) و روتکوسکی (۱۹۹۰) مراجعه نمود.

رگرسیون ناپارامتری شیوه کلی‌تری برای استنباط درباره  $g$  است. برآوردهایی از منحنی رگرسیون که با این شیوه حاصل می‌شود، انعطاف پذیرند و به یک الگوی پارامتری خاص محدود نمی‌شوند. تفاوت عمده رگرسیون پارامتری و ناپارامتری درجه اتکای آنها به اطلاعاتی است که از داده‌ها و تجربیات آماردان برای برآورد  $g$  وجود دارد. در مدل ناپارامتری فرد بنا بر اعتقاد خود، یک فضای تابعی<sup>۱</sup> مناسب برای  $g$  برمی‌گزیند. سپس از داده‌ها جهت تعیین عضوی از این فضای تابعی برای برآورد منحنی نامعلوم رگرسیون استفاده می‌کند. اما در یک الگوی پارامتری، یک خانواده از منحنی‌ها انتخاب شده و آنالیز رگرسیونی بر مبنای این خانواده از منحنی‌ها قرار می‌گیرد. بنابراین فنون رگرسیون ناپارامتری بیشتر از رگرسیون پارامتری به داده‌ها وابسته هستند و باعث می‌شوند که برآورد حاصل روند داده‌ها را بهتر تبیین نماید. به طور کلی هنگامی که مدل پارامتری معتبر باشد باید از رگرسیون پارامتری بهره گرفت و تنها در صورتی که اطلاعات در مورد  $g$  ناچیز است، رگرسیون ناپارامتری مناسب می‌باشد. زیرا در صورتی که مدل پارامتری مناسب باشد کارایی برآوردهای ناپارامتری کمتر است و در ضمن برآوردهای ناپارامتری اریب هستند.

## ۲.۱ هموارسازی

برای بررسی نحوه ارتباط متغیر پاسخ  $Y$  و متغیر مستقل  $X$  مجدداً مدل (۱.۱.۱) را در نظر می‌گیریم. تابع  $g$  نامعلوم است و هدف برآورد این تابع می‌باشد. در حالت کلی برای

<sup>۱</sup> Functional Space

تعیین برآوردی منحصر بفرد از تابع  $g$  باید شرایطی را برای آن در نظر گرفت. به عنوان مثال، در روش رگرسیون خطی فرض می‌شود  $g$  تابعی خطی از  $X$  است و با استفاده از روش کمترین توانهای دوم<sup>۲</sup> ضرایب مدل خطی بطور یکتا برآورد می‌گردند. در مواردی که برازش مدل خطی به داده‌ها مناسب نباشد، ممکن است رگرسیون چندجمله‌ای برآورد بهتری از  $g$  فراهم نماید و بسته به اینکه یک یا چند متغیر مستقل موجود باشد، مدل چندجمله‌ای یک یا چند متغیره برازش داده می‌شود. برای تعیین درجه و ضرایب چندجمله‌ای از روش پیشرو<sup>۳</sup> یا پسرو<sup>۴</sup> استفاده می‌شود که برای مطالعه آن می‌توان به نتر (۱۹۷۴) و مونت گومری (۱۹۹۱) مراجعه نمود. ساویلین (۱۹۹۱) به طور هندسی به مطالعه مدل‌های چندجمله‌ای و چندجمله‌ای‌های متعامد پرداخته است.

اگر چه رگرسیون چندجمله‌ای برای برخی از مجموعه داده‌ها مناسب می‌باشد لیکن کاستی‌هایی نیز دارد. یکی از اشکالات این روش این است که تک مشاهدات دورافتاده به طور قابل توجهی روی برآورد تابع  $g$  تاثیر می‌گذارند. هر چند چندجمله‌ای برازش شده به داده‌ها از درجه بالاتری باشد به رابطه واقعی داده‌ها نزدیکتر می‌شود اما با افزایش درجه چندجمله‌ای برازاندن مدل پیچیده‌تر می‌شود. این اشکالات از مفروض داشتن الگویی مشخص برای داده‌ها ایجاد می‌شود. برای رهایی از مشکلات فوق می‌توان  $g$  را تابعی دلخواه در نظر گرفت و اقدام به برآورد آن نمود. اگر مدل (۱.۱.۱) با روش کمترین مربعات باقیمانده‌ها برازش داده شود حالات متفاوتی را می‌توان در نظر گرفت. اگر شرطی برای تابع نامعلوم  $g$  منظور نشود و نقاط با خطوط راست به هم وصل شوند، مجموع مربعات باقیمانده‌ها در نقاط مشاهده شده صفر می‌شود، منحنی حاصل ناهموار خواهد بود. برای رفع این نقیصه می‌توان شرط همواری  $g$  را روی مدل (۱.۱.۱) اعمال نمود. میزان همواری تابع  $g$  را می‌توان توسط  $\int_a^b \{g''(x)\}^2 dx$  اندازه‌گیری نمود. بنابراین برای برآورد منحنی  $g$ ، علاوه بر عامل مجموع مربعات باقیمانده‌ها،

---

Mean Square Error<sup>۱</sup>

forward<sup>۲</sup>

backward<sup>۳</sup>

میزان ناهمواری را نیز در نظر گرفته و ملاک مجموع مربعات جریمه‌ای بصورت

$$S(g, \lambda) = \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int_a^b \{g''(x)\}^2 dx$$

تعریف می‌شود که در آن  $\lambda > 0$  پارامتر همواری<sup>۵</sup> نامیده می‌شود. تابعی مانند  $\hat{g}$  که مجموع مربعات جریمه‌ای  $S(g, \lambda)$  را کمینه نماید، اسپلاین همواری<sup>۶</sup> نامیده می‌شود. این شیوه برآورد منحنی را رهیافت جریمه ناهمواری<sup>۷</sup> نامند. در مورد روشهای جریمه ناهمواری در اوبانک (۱۹۸۸)، واهبا (۱۹۹۰)، هاردل (۱۹۹۰) و راسنبلات (۱۹۹۱) مطالب مبسوطی آمده است.

تعمیمی از اسپلاینها حالتی است که دو متغیر مستقل وجود دارد که در اینصورت اسپلاین رویه نازک<sup>۸</sup> نامیده می‌شود. در این مورد می‌توان به فصل دوم واهبا (۱۹۹۰) مراجعه نمود. در روش اسپلاین همواری مقدار  $\lambda$  نقش مهمی در برآورد منحنی ایفا می‌نماید. اگر  $\lambda$  بزرگ باشد،  $\hat{g}$  خیلی کم‌انحنی خواهد شد (شکل ۳.۲.۳). در حد وقتی  $\lambda$  به سمت بی‌نهایت میل کند، برآورد منحنی همان رگرسیون خطی خواهد بود. از طرف دیگر، اگر  $\lambda$  نسبتاً کوچک باشد برآورد  $g$  تا حد زیادی روند داده‌ها را دنبال می‌کند، حتی اگر به بهای تغییرات سریع آن تمام شود (شکل ۴.۲.۳). در وضعیت حدی اگر  $\lambda$  به صفر نزدیک شود، برآورد حاصل به اسپلاین درونیاب<sup>۹</sup> میل می‌کند (شکل ۲.۲.۳). بنابراین محاسبه مقدار مناسب  $\lambda$  در اسپلاین همواری مساله مهمی است که در مورد آن می‌توان به هاردل، مارون و هال (۱۹۸۸) و گرین و سیلورمن (۱۹۹۴) مراجعه نمود. هاردل (۱۹۹۱ و ۱۹۹۰) چگونگی انتخاب پارامتر همواری بر اساس مینیمم کردن متوسط مربع خطای برآوردگر هسته‌ای را ارائه داده

Smoothing Parameter<sup>۵</sup>

Smoothig Spline<sup>۶</sup>

Roughness Penalty<sup>۷</sup>

Thin plate<sup>۸</sup>

Interpolating Spline<sup>۹</sup>

است. دو روش اعتبار متقابل<sup>۱۰</sup> و اعتبار متقابل تعمیم یافته<sup>۱۱</sup> برای تعیین مقدار  $\lambda$ ، در استون (۱۹۷۳)، کریون و واهبا (۱۹۷۹)، گرین و سیلورمن (۱۹۹۴) و محمدزاده (۱۹۹۸) آمده است.

از جمله کاربردهای روش اسپلاین همواری در مدل بندی نیمه پارامتری<sup>۱۲</sup> است. مدل های خطی تعمیم یافته (GLM)<sup>۱۳</sup> که اولین بار نلدر و ودربرن (۱۹۷۲) آنرا ارائه نمودند و بعدها توسط افراد دیگری از جمله مک کانگ و نلدر (۱۹۸۹) و اوسالیوان (۱۹۸۶) مورد مطالعه قرار گرفته است، از دیگر مواردی است که در آنها روش اسپلاین همواری بکار گرفته می شود.

هر دو روش اسپلاین همواری و رگرسیون چند جمله ای به منظور برآورد منحنی رگرسیون بکار می رود. هر یک از دو روش فوق الذکر ممکن است برای مجموعه داده های خاصی برآورد خوبی از منحنی ایجاد نمایند.

هدف ما از این تحقیق، مقایسه دو روش اسپلاین همواری و رگرسیون چند جمله ای در برآورد منحنی رگرسیون می باشد. نکته قابل توجه این است که لزوماً هیچیک از دو روش اسپلاین همواری یا رگرسیون چند جمله ای همواره نسبت به یکدیگر برآورد بهتری از منحنی را ایجاد نمی کنند و برای مجموعه داده های مختلف روش برآورد برتر قابل تغییر است. در این پژوهش برای مقایسه این دو روش از شیوه شبیه سازی استفاده می شود. در این شیوه به تولید داده های تصادفی و برازش منحنی به آنها می پردازیم. بدین منظور ابتدا یک نمونه تصادفی  $n$  تایی از توزیع یکنواخت استخراج نموده تا مشاهدات  $x_1, \dots, x_n$  بدست آیند و سپس از توزیع  $N(0, \sigma^2)$  نمونه تصادفی  $n$  تایی گرفته می شود که مقادیر  $e_1, \dots, e_n$  را تشکیل می دهند و از رابطه  $y = x + x^2 + x^3 + e$  مقادیر  $y$  از یک مدل درجه سه حاصل می شوند.

Cross Validation<sup>۱۰</sup>Generalized Cross Validation<sup>۱۱</sup>Semiparametric<sup>۱۲</sup>Generalized Linear Model<sup>۱۳</sup>