

حاشا  
الرحمن الرحيم



دانشکده علوم ریاضی

رساله دکتری آمار

# انتخاب متغیر و تشخیص ساختار در مدل‌های نیمه پارامتری

نگارنده: محمد کاظمی

اساتید راهنما

دکتر داود شاهسونی

دکتر محمد آرشی

بهمن ۱۳۹۷



فرم شماره ۱۲: صورت جلسه نهایی دفاع از رساله دکتری (Ph.D)  
(ویژه دانشجویان ورودی های ۹۴ و ما قبل)

بدینوسیله گواهی می شود آقای محمد کاظمی دانشجوی دکتری رشته آمار- استنباط به شماره دانشجویی ۹۳۰۱۱۸۵ ورودی مهر ماه سال ۱۳۹۳ در تاریخ ۹۷/۱۱/۲۸ از رساله نظری  عملی  خود با عنوان: انتخاب متغیر و تشخیص ساختار در مدل های نیمه پارامتری دفاع و با اخذ نمره ۱۹.۸۲ به درجه معالج نائل گردید.

<input type="checkbox"/> الف) درجه عالی: نمره ۱۹-۲۰	<input type="checkbox"/> ب) درجه بسیار خوب: نمره ۱۸/۹۹-۱۷
<input type="checkbox"/> ج) درجه خوب: نمره ۱۶/۹۹-۱۵	<input type="checkbox"/> د) غیر قابل قبول و نیاز به دفاع مجدد دارد
<input type="checkbox"/> ه) رساله نیاز به اصلاحات دارد	

ردیف	هیئت داوران	نام و نام خانوادگی	مرتبه علمی	امضاء
۱	دکتر داود شاهسونی	استاد راهنمای اول	دانشیار	
۲	دکتر محمد آرشی	استاد راهنمای دوم	دانشیار	
۳	دکتر مهدی روزبه	استاد مدعو خارجی	دانشیار	
۴	دکتر حسین باغیشنی	استاد مدعو داخلی	استادیار	
۵	دکتر محمدرضا ربیعی	استاد مدعو داخلی	استادیار	
۶	دکتر ابراهیم هاشمی	سرپرست ( نماینده ) تحصیلات تکمیلی دانشکده	استاد	

مدیر محترم تحصیلات تکمیلی دانشگاه:

ضمن تأیید مراتب فوق مقرر فرمائید اقدامات لازم در خصوص انجام مراحل دانش آموختگی آقای محمد کاظمی بعمل آید.

نام و نام خانوادگی: رئیس دانشکده: دکتر ابراهیم هاشمی

تاریخ و امضاء و مهر دانشکده: ۹۷/۱۱/۲۸



## تقدیم بہ:

پدرم، کہ ہموارہ حامی و پشتیبان من بودہ است  
مادرم، کہ بیچ واثرہ ای نمی تواند مہر او را معنی کند  
ہمسرم، کہ مسج وار با صبرش در تمامی سحطات رفیق راہ بود  
دختر نامزم، امید بخش جانم کہ آسایش او آرامش من است.

# سپاس‌گزاری

سپاس ایزد منان را که در تمام مراحل زندگی مرا قوت قلب بود و توفیق را رفیق راهم ساخت تا این رساله را به پایان برسانم. اکنون که با یاری او توانسته‌ام قدمی هرچند ناچیز در راه کسب علم و دانش بردارم، بر خود لازم می‌دانم از پدر و مادر عزیزم، که همواره بر کوتاهی و درستی من، قلم عفو کشیده و با عاطفه سرشار و گرمای امیدبخش وجودشان، در این سردترین روزگاران، بهترین پشتیبان من بوده‌اند، تشکر و قدردانی نمایم. از همسر عزیزم که همواره رفیق راه و مشوق اصلی من در امر تحصیل و تحقیق بوده و موفقیت خود را از آغاز این دوره تحصیلی تاکنون مرهون حمایت‌ها و صبر و شکیبایی بی‌دریغش هستم، بی‌نهایت سپاس‌گزارم.

تقدیر به گونه‌ای رقم خورد که مقطع دکتری را در دانشگاه صنعتی شاهرود به پایان برسانم و از محضر اساتید خوب گروه آمار این دانشگاه بهره‌جوییم. سپاس از استاد عزیزم جناب آقای دکتر داود شاهسونی، استادی دلسوز و با اخلاق، که درس ”فنون آماری“ را با سبکی جذاب و عشق و علاقه‌ای وصف‌ناپذیر تدریس می‌کرد و بارها از وی درس زندگی آموختم. بی‌شک کلاس درس ایشان از بهترین خاطرات من در این دوره تحصیلی است. کاش می‌شد دوباره دانشجوی ترم اول کارشناسی آمار بودم و آمار را جور دیگری یاد می‌گرفتم، اما افسوس نه راهی برای بازگشت مانده و نه پایی برای رفتن! سپاس از استاد خوبم جناب آقای دکتر محمد آرشی، که هنوز اولین جلسه کلاس درس ”استنباط آماری پیشرفته“، ایشان در خاطر هست. استادی جوان، با اخلاق و دانا که آن روز آرزو کردم روزی مثل او باشم و هنوز هم از تلاش برای تحقق این آرزویم دست نکشیده‌ام. انگار بخت با من یار بود که این دو استاد بزرگوار زحمت راهنمایی مرا در این دوره تحصیلی برعهده گیرند. بدون شک جایگاه و منزلت معلم، والاتر از آن است که در مقام قدردانی از زحمات بی‌شائبه او، با زبان قاصر و دست ناتوان، چیزی بنگارم، اما بر حسب وظیفه و از باب ”مَنْ لَمْ يَشْكُرِ الْمُنْعِمَ مِنَ الْمَخْلُوقِينَ لَمْ يَشْكُرِ اللَّهَ عَزَّ وَ جَلَّ“ از اساتید ارجمندم جناب آقای دکتر داود شاهسونی و جناب آقای دکتر محمد آرشی که با صبر و تلاش فراوان، از هیچ کمکی در این عرصه بر من دریغ نمودند و زحمت راهنمایی بنده را در تهیه و تنظیم این رساله برعهده گرفتند، کمال تشکر و قدردانی را دارم و از درگاه خداوند متعال برای آنها توفیق روزافزون آرزومندم.

از اساتید محترم، جناب آقای دکتر مهدی روزبه (دانشگاه سمنان)، جناب آقای دکتر حسین باغیشنی (دانشگاه صنعتی شاهرود) و جناب آقای دکتر محمدرضا ربیعی (دانشگاه صنعتی شاهرود)، که زحمت تصحیح و داوری این رساله را متقبل شدند، سپاس‌گزارم.

از استاد و دوست خوبم دکتر پائولو کنس رُدِری گوئز<sup>۱</sup>، که در طول دوره فرصت مطالعاتی من در دانشگاه تامپره فنلاند زحمت راهنمایی و میزبانی مرا برعهده گرفت و از هیچ کمکی به من دریغ نکرد، سپاس‌گزارم.

<sup>۱</sup>Paulo Canas Rodrigues

از مدیر دلسوز و پرتلاش گروه آمار دانشگاه صنعتی شاهرود، جناب آقای دکتر محمدرضا ربیعی به دلیل خدمات، راهنمایی‌ها و زحمات‌شان سپاس‌گزارم. از دوستان و همکلاسی‌های خوبم که مرا در رسیدن به اهدافم کمک کرده‌اند، تشکر کرده و برایشان آرزوی موفقیت و سربلندی می‌کنم. در نهایت، از همه عزیزانی که تا بدین‌جا مرا یاری کرده‌اند، به‌ویژه خانواده‌مهربانم که همواره پشتیبان من بوده‌اند، سپاس‌گزارم. امیدوارم بتوانم از عهده ادای حق این عزیزان برآیم. یقیناً این مجموعه بدون عیب و نقص نیست. لذا از خوانندگان عزیز تقاضا دارم عیب‌ها و کاستی‌های این رساله را از طریق پست الکترونیکی<sup>۲</sup> به اطلاع اینجانب برسانند.

محمد کاظمی

بهمن ۱۳۹۷

## تعهد نامه

اینجانب **محمد کاظمی** دانشجوی دکتری رشته **آمار دانشکده علوم ریاضی** دانشگاه صنعتی شاهرود، نویسنده رساله با عنوان **انتخاب متغیر و تشخیص ساختار در مدل‌های نیمه پارامتری**، تحت راهنمایی **دکتر داود شاهسونی و دکتر محمد آرشی** متعهد می‌شوم:

- تحقیقات در این رساله توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش‌های دیگر پژوهش‌گران، به مرجع مورد استفاده استناد شده است.
- مطالب این رساله، تا کنون توسط خود، یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ‌جا ارایه نشده است.
- حقوق معنوی این اثر، به دانشگاه صنعتی شاهرود تعلق دارد، و مقالات مستخرج با نام ”دانشگاه صنعتی شاهرود“ یا ”Shahrood University of Technology“ به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به‌دست آوردن نتایج اصلی رساله تاثیرگذار بوده‌اند، در مقالات مستخرج از رساله رعایت می‌گردد.
- در تمام مراحل انجام این رساله، در مواردی که از موجود زنده (یا بافت‌های آنها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در تمام مراحل انجام این رساله، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته (یا استفاده شده است)، اصل رازداری و اصول اخلاق انسانی رعایت شده است.

**محمد کاظمی**

**بهمن ۱۳۹۷**

### مالکیت نتایج و حق نشر

- تمام حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه‌های رایانه‌ای، نرم‌افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می‌باشد. این مطلب باید به نحو مقتضی، در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در این رساله بدون ذکر منبع مجاز نمی‌باشد.





## چکیده

مدل‌های نیمه‌پارامتری، به دلیل برخورداری از انعطاف‌پذیری مناسب و تفسیرپذیری، اغلب در تحلیل داده‌های واقعی مورد استفاده قرار می‌گیرند. دو مورد از رایج‌ترین مدل‌های نیمه‌پارامتری، مدل جمعی خطی-جزئی و مدل ضریب متغیر خطی-جزئی می‌باشند. استنباط آماری این مدل‌ها بر این شرط استوار است که اجزای پارامتری و ناپارامتری مدل از پیش تعیین شده باشند. اما در عمل چنین اطلاعات پیشین به ندرت در دسترس است، به ویژه هنگامی که تعداد متغیرهای توضیحی زیاد است. بنابراین لازم است روشی ارائه شود که بتواند مؤلفه‌های پارامتری و ناپارامتری را به درستی تشخیص دهد. در این رساله، به معرفی یک روش دو مرحله‌ای برای انتخاب متغیرهای مهم و تشخیص اثرات خطی و غیرخطی در مدل جمعی خطی-جزئی با بعد بسیار بالا پرداخته می‌شود. علاوه بر این، با استفاده از رگرسیون نما، یک روش نیرومند برای انتخاب متغیرهای مهم و تشخیص اثرات ثابت و متغیر در مدل ضریب متغیر خطی-جزئی ارائه می‌شود که نسبت به داده‌های پرت در متغیر پاسخ و توزیع غیر نرمال خطاها نیرومند است. کارایی روش‌های ارائه شده با مطالعه شبیه‌سازی و تحلیل داده‌های واقعی بررسی شده است. نتایج به دست آمده عملکرد مطلوب روش‌های مذکور را در مقایسه با برخی روش‌های موجود نشان می‌دهند.

**کلمات کلیدی:** انتخاب متغیر، تشخیص ساختار، رگرسیون نما، غربالگری مستقل مطمئن، کاهش بعد، مدل جمعی خطی-جزئی، مدل ضریب متغیر خطی-جزئی.

## لیست مقالات مستخرج از رساله

۱. کاظمی، م.، شاهسونی، د. و آرشی، م. (۱۳۹۷). انتخاب متغیر و تشخیص ساختار در بعد بالا برای مدل‌های جمعی خطی-جزئی، *مجله علوم آماری*، ۱۲(۲)، ۴۸۵-۵۱۲.
۲. کاظمی، م.، شاهسونی، د. و آرشی، م. (۱۳۹۶). غربالگری متغیرها در داده‌های بزرگ با استفاده از همبستگی فاصله‌ای، کنفرانس روش‌های مدرن در قیمت‌گذاری‌های بیمه‌ای و آمارهای صنعتی، دانشگاه بوعلی‌سینا، همدان.
۳. کاظمی، م.، شاهسونی، د. و آرشی، م. (۱۳۹۶). حل مسئله بهینه‌سازی برای انتخاب متغیرهای مهم در داده‌های با بعد بالا، دهمین کنفرانس بین‌المللی انجمن ایرانی تحقیق در عملیات، دانشگاه مازندران، بابلسر.
4. Kazemi, M., Shahsavani, D. and Arashi, M. (2019). A sure independence screening procedure for ultra-high dimensional partially linear additive models, *Journal of Applied Statistics*, 46(8), 1385-1403.
5. Kazemi, M., Shahsavani, D., Arashi, M. and Rodrigues, P. C. (2019). Estimation in partial linear model with spline modal function, *Communications in Statistics - Simulation and Computation* (Revised).
6. Kazemi, M., Shahsavani, D. and Arashi, M. (2018). Variable selection and structure identification for ultrahigh-dimensional partially linear additive models with application to cardiomyopathy microarray data, *Statistics, Optimization and Information Computing*, DOI: 10.19139/soic.v6i3.576.
7. Kazemi, M., Arashi, M. and Shahsavani, D. (2017), Gene identification using independence screening procedure, The 11th Seminar on Probability and Stochastic Processes, Imam Khomeini International University, Qazvin, Iran.
8. Kazemi, M., Shahsavani, D. and Arashi, M. (2017). Group LASSO for high-dimensional partially linear errors-in-variables models, Second National and First International Conference on Soft Computing, University of Guilan, Faculty of Engineering and Technology East of Guilan, Rudsar, Iran.

## پیش‌گفتار

مدل‌های رگرسیونی نیمه‌پارامتری در سال‌های اخیر مورد توجه بسیاری از محققین قرار گرفته‌اند. یک رده‌پرکاربرد از مدل‌های نیمه‌پارامتری، رده مدل‌های خطی - جزئی، مانند مدل جمعی خطی - جزئی و مدل ضریب متغیر خطی - جزئی، است. در مدل جمعی خطی - جزئی، برخی از متغیرهای توضیحی دارای اثر خطی و برخی دیگر دارای اثر غیر خطی بر متغیر پاسخ هستند و در مدل ضریب متغیر خطی - جزئی، اثرات بعضی از متغیرهای توضیحی ثابت و اثرات برخی دیگر از متغیرها تابعی از یک متغیر شاخص است. علاوه بر این برخی از متغیرهای توضیحی ارتباطی با متغیر پاسخ ندارند. در برازش مدل‌های خطی - جزئی، معمولاً فرض بر این است که مؤلفه‌های خطی و غیر خطی یا اثرات ثابت و متغیر از پیش تعیین شده‌اند. اما در عمل معمولاً تعیین این مؤلفه‌ها بسیار مشکل است. بنابراین موضوعاتی از قبیل انتخاب متغیرهای مهم، تشخیص مؤلفه‌های خطی و غیر خطی در مدل جمعی خطی - جزئی و تعیین اثرات ثابت و متغیر در مدل ضریب - متغیر خطی - جزئی بسیار حائز اهمیت است. از طرفی امروزه با افزایش توان محاسباتی کامپیوترها و روش‌های نوین جمع‌آوری داده‌ها، داده‌های با بعد بسیار بالا در زمینه‌های مختلف علمی و تحقیقاتی با هزینه نسبتاً کمی در اختیار محققین قرار می‌گیرد. در این نوع داده‌ها فقط تعداد اندکی از متغیرهای توضیحی واقعاً با متغیر پاسخ مرتبط هستند و سایر متغیرها تأثیری بر پاسخ ندارند. هنگامی که تعداد متغیرهای توضیحی در مدل زیاد است، تفسیر آن مدل بسیار مشکل است. همچنین حضور تعداد زیادی از متغیرهای توضیحی باعث ایجاد همخطی بین متغیرها شده و منجر به برآوردگری با واریانس بزرگ می‌شود. بنابراین در این نوع داده‌ها شناسایی متغیرهای مؤثر بر پاسخ اهمیت بسزایی دارد.

لیان<sup>۳</sup> (۲۰۱۲c)، لیان و همکاران (۲۰۱۲) و لیان و همکاران (۲۰۱۵) به موضوع انتخاب متغیر و تشخیص ساختار در مدل جمعی خطی - جزئی پرداختند. مطالعات مذکور برای انتخاب متغیر و تشخیص ساختار در داده‌های با بعد بسیار بالا، به دلایل هزینه محاسباتی، دقت آماری و ناپایداری الگوریتمی، مناسب نیستند. لذا شناسایی متغیرهای مهم در این نوع داده‌ها و تشخیص نوع تأثیر آنها بر متغیر پاسخ یک چالش است. در این حالت لازم است قبل از انتخاب متغیر و تشخیص ساختار، از یک روش غربالگری برای کاهش بعد استفاده شود. در این راستا، فن و لیو<sup>۴</sup> (۲۰۰۸) روش غربالگری مستقل مطمئن را برای کاهش بعد در مدل خطی پیشنهاد دادند. در این روش، متغیرهای توضیحی براساس همبستگی پیرسن آنها با متغیر پاسخ رتبه‌بندی شده و سپس تعدادی متغیر با بیشترین همبستگی انتخاب می‌شوند.

<sup>3</sup>Lian

<sup>4</sup>Fan and LV

برای انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی- جزئی، زی<sup>۵</sup> و همکاران (۲۰۰۴)، لنگ<sup>۶</sup> (۲۰۰۹)، نه و کیلگوم<sup>۷</sup> (۲۰۱۲)، هو<sup>۸</sup> و زی<sup>۹</sup> (۲۰۱۲) روش‌های مختلفی را برای تشخیص اثرات ثابت و متغیر در مدل ضریب متغیر خطی- جزئی معرفی کردند، اما این روش‌ها متغیرهای مهم را انتخاب نمی‌کنند. تانگ<sup>۹</sup> و همکاران (۲۰۱۲) یک روش واحد برای انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی- جزئی معرفی کردند، اما این روش مستلزم استفاده از یک الگوریتم دو مرحله‌ای تکراری است که در عمل چندان کارا نیست. علاوه بر این روش معرفی شده بر اساس رگرسیون میانگین نسبت به مشاهدات پرت و توزیع دم‌سنگین خطاها بسیار حساس است. بنابراین ارائه یک روش کارا و نیرومند برای انتخاب متغیر و تشخیص اثرات ثابت و متغیر در این مدل ضروری به نظر می‌رسد. در این رساله به مسئله انتخاب متغیر و تشخیص ساختار در مدل‌های نیمه‌پارامتری خطی- جزئی پرداخته می‌شود. این مجموعه شامل ۴ فصل و ۴ پیوست است:

- **فصل اول**، به معرفی مفاهیم مقدماتی در تحلیل رگرسیونی و مدل‌های مورد مطالعه اختصاص دارد.
- در **فصل دوم**، به بیان جزئیات روش‌های انتخاب متغیر، مانند روش‌های انتخاب متغیر کلاسیک، روش کمترین توان‌های دوم جریمه‌شده و روش غربالگری پرداخته می‌شود.
- در **فصل سوم**، یک روش دو مرحله‌ای برای انتخاب متغیر و تشخیص مؤلفه‌های خطی و غیرخطی در مدل جمعی با بعد بسیار بالا معرفی می‌شود. در مرحله اول، از یک روش غربالگری برای کاهش بعد فضای متغیرهای توضیحی استفاده می‌شود. این روش غربالگری بر اساس همبستگی فاصله‌ای بین متغیرهای توضیحی و تابع توزیع حاشیه‌ای متغیر پاسخ ساخته شده و زمانی که متغیر پاسخ دم‌سنگین یا دارای مقادیر فرین باشد، عملکرد خوبی را از خود نشان می‌دهد. در مرحله دوم، از روشی مبتنی بر دو تابع جریمه برای انتخاب همزمان مؤلفه‌های غیرصفر و خطی استفاده می‌شود.
- در **فصل چهارم**، یک روش نیرومند برای انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی- جزئی بر اساس رگرسیون نما<sup>۱۰</sup> ارائه می‌شود. رگرسیون نما توسط یائو<sup>۱۱</sup> و همکاران (۲۰۱۲) و یائو و لی<sup>۱۲</sup> (۲۰۱۴) بر اساس نمای شرطی متغیر پاسخ معرفی شده است.
- **پیوست آ**، اثبات قضایای بیان‌شده در فصل‌های ۳ و ۴ را در برمی‌گیرد.

<sup>5</sup>Xia

<sup>6</sup>Leng

<sup>7</sup>Noh and Keilegom

<sup>8</sup>Hu

<sup>9</sup>Tang

<sup>10</sup>Modal regression

<sup>11</sup>Yao

<sup>12</sup>Li

- پیوست ب، مرور مختصری بر توابع اسپلاین و B- اسپلاین را شامل می‌شود.
- پیوست پ، نامساوی‌ها و تعاریف مورد استفاده در این رساله را یادآوری می‌نماید.
- پیوست ت، شامل گزیده‌ای از برنامه‌های نوشته‌شده با نرم‌افزار R است.



# فهرست مطالب

ق	فهرست تصاویر
ش	فهرست جداول
ث	فهرست نمادها
ذ	فهرست کلمات اختصاری
۱	۱ مقدمات و مفاهیم اولیه
۱	۱.۱ مدل بندی رگرسیونی . . . . .
۳	۱.۱.۱ رگرسیون پارامتری . . . . .
۵	۲.۱.۱ رگرسیون ناپارامتری . . . . .
۷	۳.۱.۱ رگرسیون نیمه پارامتری . . . . .
۱۳	۲.۱ اهداف رساله . . . . .
۱۵	۳.۱ ساختار رساله . . . . .
۱۷	۲ انتخاب متغیر در مدل های رگرسیونی
۱۷	۱.۲ مقدمه . . . . .
۱۹	۲.۲ روش کمترین توان های دوم . . . . .
۲۰	۱.۲.۲ همخطی در رگرسیون چندگانه . . . . .
۲۳	۳.۲ انتخاب متغیر با معیارهای کلاسیک . . . . .
۲۵	۴.۲ انتخاب متغیر با کمترین توان های دوم جریمه شده . . . . .
۲۶	۱.۴.۲ انواع توابع جریمه . . . . .
۳۲	۲.۴.۲ ویژگی های یک تابع جریمه مناسب . . . . .
۳۵	۳.۴.۲ انتخاب پارامتر جریمه . . . . .
۳۷	۵.۲ انتخاب متغیر گروهی . . . . .
۳۷	۱.۵.۲ جریمه لاسوی گروهی . . . . .
۳۹	۲.۵.۲ لاسوی گروهی تطبیقی . . . . .

۳۹	سایر توابع جریمه گروهی	۳.۵.۲
۴۱	انتخاب متغیر در داده‌های با بعد بسیار بالا	۶.۲
۴۳	غربالگری مستقل مطمئن	۱.۶.۲
۴۷	غربالگری و رتبه بندی مستقل مطمئن	۲.۶.۲
۴۸	غربالگری با همبستگی فاصله‌ای	۳.۶.۲
۵۰	الگوریتم محاسباتی روش کمترین توان‌های دوم جریمه‌شده	۷.۲
۵۱	تقریب درجه دو موضعی	۱.۷.۲
۵۳	تقریب خطی موضعی	۲.۷.۲
<b>۵۵</b>	<b>انتخاب متغیر و تشخیص ساختار در مدل جمعی خطی - جزئی</b>	<b>۳</b>
۵۵	مقدمه	۱.۳
۵۷	غربالگری مستقل مطمئن نیرومند بر اساس همبستگی فاصله‌ای	۲.۳
۵۹	بررسی خواص نظری روش RDC-SIS	۱.۲.۳
۶۰	غربالگری مستقل مطمئن تکراری	۲.۲.۳
۶۲	انتخاب متغیر و تشخیص ساختار	۳.۳
۶۴	مطالعه شبیه‌سازی	۴.۳
۷۳	تحلیل داده‌های مسمومیت کبد	۵.۳
<b>۷۷</b>	<b>انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی - جزئی</b>	<b>۴</b>
۷۷	مقدمه	۱.۴
۷۹	رگرسیون نما	۱.۱.۴
۸۲	روش نیرومند انتخاب متغیر و تشخیص ساختار	۲.۴
۸۴	ویژگی‌های مجانبی	۳.۴
۸۶	انتخاب پهنای باند و الگوریتم پیاده‌سازی	۴.۴
۸۶	انتخاب پهنای باند بهینه	۱.۴.۴
۸۷	الگوریتم MEM	۲.۴.۴
۸۸	انتخاب پارامترهای تنظیم‌کننده	۳.۴.۴
۸۹	مطالعات شبیه‌سازی	۵.۴
۹۲	تحلیل داده‌های قیمت مسکن	۶.۴
<b>۹۹</b>	<b>اثبات قضایا</b>	<b>آ</b>
۹۹	قضایای فصل ۳	۱.آ
۱۰۲	قضایای فصل ۴	۲.آ
<b>۱۱۱</b>	<b>توابع اسپلاین و B-اسپلاین</b>	<b>ب</b>
۱۱۱	اسپلاین	۱.ب



- ۱۱۳ . . . . . ب.۲ رگرسیون اسپلین
- ۱۱۵ . . . . . ب.۳  $B$ -اسپلین
- ۱۲۲ . . . . . ب.۴ کمترین توان‌های دوم  $B$ -اسپلین

۱۲۳ پ تعاریف و نامساوی‌ها

۱۲۷ ت گزیده‌ای از برنامه‌های رایانه‌ای

۱۴۷ مراجع

۱۶۳ واژه‌نامه فارسی به انگلیسی

۱۶۷ واژه‌نامه انگلیسی به فارسی



# فهرست تصاویر

۴	۱.۱	قاب سمت چپ اثر متغیر تحصیل و قاب سمت راست اثر متغیر سابقه کار بر پاسخ را نشان می‌دهد.
۵	۲.۱	تابع رگرسیونی برآورد شده به صورت پارامتری. منحنی قرمز رنگ، نیمرخ دستمزد- سابقه کار را برای فارغ‌التحصیلان دبیرستان نشان می‌دهد.
۶	۳.۱	تابع رگرسیونی برآورد شده با روش ناپارامتری. منحنی قرمز رنگ، نیمرخ دستمزد- سابقه کار را برای فارغ‌التحصیلان دبیرستان نشان می‌دهد.
۸	۴.۱	برآزش مدل جمعی در مقابل برآزش مدل پارامتری. قاب سمت چپ اثر متغیر تحصیل و قاب سمت راست اثر متغیر سابقه کار بر پاسخ را نشان می‌دهد.
۹	۵.۱	تابع رگرسیونی برآورد شده با روش نیمه پارامتری. منحنی قرمز رنگ، نیمرخ دستمزد- سابقه کار را برای فارغ‌التحصیلان دبیرستان نشان می‌دهد.
۱۲	۶.۱	توابع ضرایب برآورد شده مثال ۷.۱.۱. نمودارها، به ترتیب، از چپ به راست $\hat{\alpha}_1(u)$ ، $\hat{\alpha}_2(u)$ و $\hat{\alpha}_3(u)$ را نشان می‌دهند.
۲۸	۱.۲	منحنی‌های تراز و نواحی محدودیت رگرسیون لاسو (قاب سمت چپ) و رگرسیون ریج (قاب سمت راست). ناحیه محدودیت لاسو به صورت $ \beta_1  +  \beta_2  \leq s$ و ناحیه محدودیت ریج به صورت $\beta_1^2 + \beta_2^2 \leq s$ است. بیضی‌ها، منحنی‌های تراز تابع خطای برآورد کمترین توان‌های دوم هستند.
۳۲	۲.۲	نمودار توابع جریمه لاسو، SCAD و MCP.
۳۳	۳.۲	نمودار برآوردهای کمترین توان‌های دوم جریمه شده (PLS) در مقابل برآورد کمترین توان‌های دوم معمولی (OLS) وقتی که ماتریس طرح متعامد است. محور افقی OLS و محور عمودی PLS است.
۳۵	۴.۲	نمودارهای مشتق توابع جریمه لاسو، SCAD و MCP.
۴۰	۵.۲	برآوردهای لاسوی گروهی، SCAD گروهی و MCP گروهی برای مقادیر مختلف $\lambda$ . خط عمودی نقطه چین مقدار $\lambda = 0.04$ را نشان می‌دهد. نمودارهای آبی و قرمز، به ترتیب، برآوردهای ضرایب گروه اول و دوم و بقیه نمودارها مربوط به گروه‌های صفر هستند.
۴۱	۶.۲	برآوردهای SCAD گروهی برای مقادیر مختلف $\lambda$ .

۴۴	..... روش‌های انتخاب متغیر در ابعاد خیلی بالا	۷.۲
	تقریب درجه دو موضعی و تقریب خطی موضعی برای توابع جریمه $L_{0.5}$ و SCAD.	۸.۲
	نمودار آبی تقریب درجه دو موضعی، نمودار قرمز تقریب خطی موضعی و نمودار مشکی تابع جریمه را نشان می‌دهد.	۵۴
۷۵	نمودارهای پراکنش متغیر پاسخ در مقابل دو ژن مهم و منحنی‌های اسپلاین مکعبی.	۱.۳
۷۶	نمودار مقادیر مشاهده شده متغیر پاسخ در مقابل مقادیر پیش بینی شده	۲.۳
۸۱	نمودار $-\phi_n(\cdot)$ به ازای پارامترهای مختلف.	۱.۴
۹۴	نمودار جعبه‌ای متغیر پاسخ MDEV.	۲.۴
	برآوردهای ضرایب متغیر مدل (۱۷.۴) و فواصل اطمینان نقطه‌ای ۹۵٪ این ضرایب	۳.۴
۹۶	با استفاده از رگرسیون نمای جریمه‌شده (PMR).	
	نمودار پایه‌های $B$ - اسپلاین مرتبه یک که به ترتیب از سمت چپ نمودارهای $B_{-1,1}(x)$ ,	۱.ب
۱۱۷	..... $B_{1,1}(x)$ و $B_{2,1}$ هستند.	
	نمودار پایه‌های $B$ - اسپلاین مرتبه دو که به ترتیب از سمت چپ نمودارهای $B_{-2,2}(x)$ ,	۲.ب
۱۱۹	..... $B_{1,2}(x)$ ، $B_{0,2}(x)$ ، $B_{2,2}(x)$ و $B_{3,2}(x)$ هستند.	
	نمودار پایه‌های $B$ - اسپلاین مرتبه سه که به ترتیب از سمت چپ نمودارهای $B_{-3,3}(x)$ ,	۳.ب
۱۲۲	..... $B_{1,3}(x)$ ، $B_{0,3}(x)$ ، $B_{2,3}(x)$ ، $B_{3,3}(x)$ و $B_{4,3}(x)$ هستند.	

# فهرست جداول

۴	نتایج برآورد کمترین توان‌های دوم در مثال ۱.۱.۱	۱.۱
۴۳	مقادیر $\frac{p}{n}$ برای داده‌های با ابعاد بالا و بسیار بالا	۱.۲
۶۶	احتمال تجربی $P_j$ و احتمال $S$ در مثال ۱.۴.۳	۱.۳
۶۷	چندک‌های $M$ ، احتمال تجربی $P_j$ و احتمال $S$ در مدل ۱	۲.۳
۶۸	چندک‌های $M$ ، احتمال تجربی $P_j$ و احتمال $S$ در مدل ۲	۳.۳
۶۹	احتمال تجربی $P_j$ و احتمال $S$ در مثال ۳	۴.۳
۷۱	احتمال تجربی $P_j$ و احتمال $S$ در مثال ۴	۵.۳
۷۲	نتایج تشخیص ساختار مدل در مثال ۵	۶.۳
۷۳	مجذور میانگین توان‌های دوم خطا برای پنج مؤلفه اول و تابع رگرسیون	۷.۳
۷۴	چهار ژن بسیار مهم انتخاب شده توسط روش‌های مختلف	۸.۳
۷۴	ژن‌های مهم انتخاب شده و نوع تأثیر آن‌ها	۹.۳
۹۱	نتایج تشخیص ساختار روش‌های مختلف در مثال ۱.۵.۴	۱.۴
۹۲	درصد انتخاب درست ضرایب توسط رگرسیون نما در مثال ۲.۵.۴	۲.۴
	نتایج انتخاب متغیر و برآورد ضرایب ثابت غیر صفر در داده‌های مسکن بوستون. نمادهای $\circ$ و $\vee$ ، به ترتیب، بیانگر اثر صفر و اثر متغیر است. مقادیر سطر آخر جدول، ضریب تعیین تعدیل شده و انحراف معیار باقیمانده‌ها را نشان می‌دهند.	۳.۴
۹۵		



## فهرست نمادها

$\mathbf{X}$	.....	ماتریس طرح
$\mathbf{Y}$	.....	بردار پاسخ
$\beta$	.....	بردار ضرایب
$\beta_0$	.....	مقدار درست بردار $\beta$
$\varepsilon$	.....	بردار خطا
$U$	.....	متغیر شاخص در مدل ضریب متغیر
$I(A)$	.....	تابع نشانگر مجموعه $A$
$\mathbf{I}_m$	.....	ماتریس واحد با بعد $m$
$\mathbf{A}^T$	.....	ترانزاده ماتریس $\mathbf{A}$
$\mathbf{A}^{-1}$	.....	معکوس ماتریس $\mathbf{A}$
$\mathbf{C}$	.....	معکوس ماتریس $\mathbf{X}^T \mathbf{X}$
$R^2$	.....	ضریب تعیین
$R^2_{adj}$	.....	ضریب تعیین تعدیل شده
$\ \cdot\ $	.....	نرم اقلیدسی
$Var$	.....	واریانس
$Cov$	.....	کوواریانس
$dcov$	.....	کوواریانس فاصله‌ای
$dcorr$	.....	ضریب همبستگی فاصله‌ای
$\xrightarrow{D}$	.....	همگرایی در توزیع
$\xrightarrow{P}$	.....	همگرایی در احتمال
$E[\cdot]$	.....	امید ریاضی
$\hat{\beta}^{OLS}$	.....	برآوردگر کمترین توان‌های دوم معمولی
$\hat{\beta}^{Ridge}$	.....	برآوردگر ریج
$\hat{\beta}^{Lasso}$	.....	برآوردگر لاسو
$\hat{\beta}^{aLasso}$	.....	برآوردگر لاسوی تطبیقی
$\hat{\beta}^{SCAD}$	.....	برآوردگر SCAD
$\hat{\beta}^{MCP}$	.....	برآوردگر MCP

---



---

$p_\lambda(\cdot)$	تابع جریمه
$\mathcal{K}$	معیار ماکزیمم تقعر تابع جریمه MCP
$\text{diag}\{a_1, \dots, a_k\}$	ماتریس قطری با اعضای قطر اصلی $a_1, \dots, a_k$
$\text{diag}\{\mathbf{A}_1, \dots, \mathbf{A}_k\}$	ماتریس بلوکی با ماتریس‌های $\mathbf{A}_1, \dots, \mathbf{A}_k$
$\mathbb{M}$	مجموعه اندیس متغیرهای مهم
$K_h(\cdot)$	تابع هسته
$\mathbf{B}(\cdot)$	پایه‌های B – اسپلاین
$\mathcal{N}_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	توزیع نرمال $k$ متغیره با بردار میانگین $\boldsymbol{\mu}$ و ماتریس کوواریانس $\boldsymbol{\Sigma}$



# فهرست کلمات اختصاری

aLasso(Adaptive Least Absolute Shrinkage and Selection Operator)  
AIC(Akaike's Information Criterion)  
BIC(Bayesian Information Criterion)  
DC-SIS(Sure Independence Screening Procedure based on the Distance Correlation)  
GCV(Generalized Cross-Validation)  
ISIS(Iterative Sure Independence Screening)  
Lasso(Least Absolute Shrinkage and Selection Operator)  
LLA(Local Linear Approximation)  
LQA(Local Quadratic Approximation)  
MARS(Multivariate Adaptive Regression Splines)  
MCP(Minimax Concave Penalty )  
MEM(Modal Expectation-Maximization)  
NIS(Nonparametric Independence Screening)  
OLS(Ordinary Least Squares)  
PMR(Penalized Modal Regression)  
PQR(Penalized Quantile Regression)  
RDC-SIS(Robust Sure Independence Screening Procedure based on Distance Correlation)  
RMSE(Square Root of Mean Squared Errors)  
RSS(Residual Sum of Squares)  
SCAD( Smoothly Clipped Absolute Deviation)  
SIRS(Sure Independent Ranking and Screening)  
SIS(Sure Independence Screening)  
VIF(Variance Inflation Factor)



# فصل ۱

## مقدمات و مفاهیم اولیه

در این فصل، انواع مدل‌های رگرسیونی و مدل‌های مورد بررسی در این رساله، به اختصار بیان می‌شوند. مطالب عمده این فصل از هاردل و همکاران (۲۰۰۰)، هاردل و همکاران (۲۰۰۴) و هستی و همکاران (۲۰۰۹) گرفته شده است.

### ۱.۱ مدل بندی رگرسیونی

تحلیل رگرسیونی روشی ساده برای بررسی روابط تابعی میان متغیرهاست. به عنوان مثال، فرض کنید بخواهیم بررسی کنیم که آیا میزان دستمزد با تحصیلات و سابقه کار در بازار مرتبط است یا خیر. این رابطه به شکل یک معادله یا الگویی است که متغیر وابسته یا پاسخ را به یک یا چند متغیر توضیحی یا پیشگو مرتبط می‌کند. در این مثال، میزان حقوق متغیر پاسخ و متغیرهای تحصیل و سابقه کار متغیرهای توضیحی (پیشگو یا کمکی) هستند.

در ابتدا، رگرسیون خطی چندگانه را در نظر بگیرید. در این مدل، میانگین متغیر پاسخ  $Y$  به صورت

$$E(Y|\mathbf{X}) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \quad (1.1)$$

با متغیرهای توضیحی  $\mathbf{X} = (X_1, \dots, X_p)^T$  در ارتباط است که  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی نامعلوم است.

با تعریف  $\varepsilon$  به عنوان انحراف  $Y$  از میانگین شرطی  $E(Y|\mathbf{X})$ ، یعنی  $\varepsilon = Y - E(Y|\mathbf{X})$ ، مدل خطی را می‌توان به صورت زیر تعریف کرد

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon \quad (2.1)$$

که در آن  $\varepsilon$  می‌تواند نشان‌دهنده عواملی باشد که در مدل بندی تصادفی  $Y$  تأثیر داشته، اما در  $\mathbf{X}$  در نظر گرفته نشده‌اند.

**مثال ۱.۱.۱.** (هاردل<sup>۱</sup> و همکاران، ۲۰۰۴) فرض کنید  $Y$  لگاریتم دستمزد باشد. متغیرهای توضیحی تحصیل ( $School$ ) برحسب سال، سابقه کار در بازار ( $Exp$ ) و توان دوم متغیر  $Exp$  را در نظر بگیرید که متغیر  $Exp$  از رابطه  $Exp = Age - School - 6$  محاسبه می‌شود و  $Age$  بیانگر سن است. اگر فرض کنیم لگاریتم دستمزد به صورت خطی با این متغیرهای توضیحی مرتبط است، در این صورت مدل رگرسیون خطی عبارتست از

$$E(Y|School, Exp) = \beta_0 + \beta_1 School + \beta_2 Exp + \beta_3 Exp^2. \quad (3.1)$$

معادله (۳.۱) به معادله درآمد سرمایه انسانی<sup>۲</sup> یا معادله درآمد مینسر<sup>۳</sup> معروف است و نقش مهمی را در اقتصاد تجربی بازی می‌کند. ویژگی مهم مدل (۳.۱) صورت پارامتری آن است، یعنی شکل تابع رگرسیونی به ضرایب رگرسیونی نامعلوم  $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)^T$  بستگی دارد. بنابراین، برای تعیین تابع رگرسیونی، لازم است این ضرایب مجهول برآورد شوند. این مدل، روابط غیر خطی احتمالی بین متغیرهای توضیحی و پاسخ را نادیده می‌گیرد. لذا هنگامی که رابطه میان متغیر پاسخ و متغیرهای توضیحی نامعلوم است، استفاده از رگرسیون خطی ممکن است منجر به نتایج نادرست شود. حال فرض کنید شکل درست (اما نامعلوم) تابع رگرسیونی لگاریتم دستمزد روی متغیرهای توضیحی تحصیل و سابقه کار به صورت

$$E(Y|School, Exp) = m(School, Exp) \quad (4.1)$$

باشد و شما بخواهید تابع رگرسیونی  $m(School, Exp)$  را تنها در یک بار کوشش و تا حد ممکن دقیق برآورد کنید. بدین معنی که اگر پی بردید مدل اولیه داده‌ها را به خوبی برازش نمی‌کند، دیگر مجاز به تغییر مدل نیستید. در ابتدا ممکن است فرض کنید، تابع رگرسیونی که مسئولیت برآورد آن را به عهده گرفته‌اید، دارای صورت مشخص (۳.۱) است، یعنی فرض می‌کنید که

$$m(School, Exp) = \beta_0 + \beta_1 School + \beta_2 Exp + \beta_3 Exp^2$$

و سپس ضرایب نامعلوم را با استفاده از روش کمترین توان‌های دوم برآورد کنید. اما اگر به شما گفته شود که راه‌های دیگری برای برآورد تابع رگرسیونی فوق وجود دارند که به‌جز پذیره هموار بودن، هیچ

<sup>1</sup>Härdle

<sup>2</sup>Human capital earnings equation

<sup>3</sup>Mincer earnings equation

پذیره پیشین دیگری برای شکل تابع رگرسیونی در نظر نمی‌گیرند، شاید دیگر مدل پارامتری فوق را برازش نکنید. لازم به ذکر است که در اینجا منظور از هموار بودن، مشتق‌پذیر بودن تابع  $m(\cdot)$  تا یک مرتبه خاص و کوچک است. به خاطر داشته‌باشید که در این مسئله، شما تنها یک فرصت دارید و اگر شکل درست  $m(School, Exp)$  خیلی متفاوت از (۳.۱) باشد، برآورد مدل پارامتری ممکن است نتایج گمراه‌کننده‌ای را به دنبال داشته باشد. روش‌های متعددی برای برآورد تابع  $m(\cdot)$  وجود دارند که تنها فرض می‌شود  $m(\cdot)$  یک تابع هموار است. به این روش‌ها، رگرسیون ناپارامتری گویند.

برآوردگرهای رگرسیون ناپارامتری بسیار انعطاف‌پذیر هستند، اما با افزایش تعداد متغیرهای توضیحی در مدل، دقت آماری این برآوردگرها به شدت کاهش می‌یابد (هاردل و همکاران، ۲۰۰۴). این مشکل به مشقت بُعد<sup>۴</sup> معروف است. در نتیجه، محققان سعی کرده‌اند به مدل‌ها و برآوردگرهایی بپردازند که انعطاف‌پذیری بیشتری نسبت به رگرسیون پارامتری استاندارد دارند، و در عین حال بر مشکل مشقت بعد غلبه می‌کنند. چنین روش‌هایی معمولاً خصوصیات روش‌های پارامتری و ناپارامتری را با هم ترکیب می‌کنند و لذا به روش‌های نیمه‌پارامتری معروف هستند. مزیت‌های دیگر روش‌های نیمه‌پارامتری عبارتند از: امکان ورود متغیرهای طبقه‌ای به مدل، سهولت تفسیر نتایج و امکان مشخص کردن بخشی از مدل با استفاده از اطلاعات پیشین. بنابراین، رگرسیون را بر حسب نوع تابعی که رابطه بین متغیرهای توضیحی و متغیر پاسخ را تعیین می‌کند، می‌توان به سه نوع پارامتری، ناپارامتری و نیمه‌پارامتری تقسیم کرد.

در ادامه، برای درک بهتر مطلب و پیش‌درآمدی بر مطالب کلیدی این رساله، از معادله درآمد و مثال‌های دیگر برای نشان دادن تفاوت‌های بین روش‌های پارامتری، ناپارامتری و نیمه‌پارامتری استفاده می‌کنیم.

### ۱.۱.۱ رگرسیون پارامتری

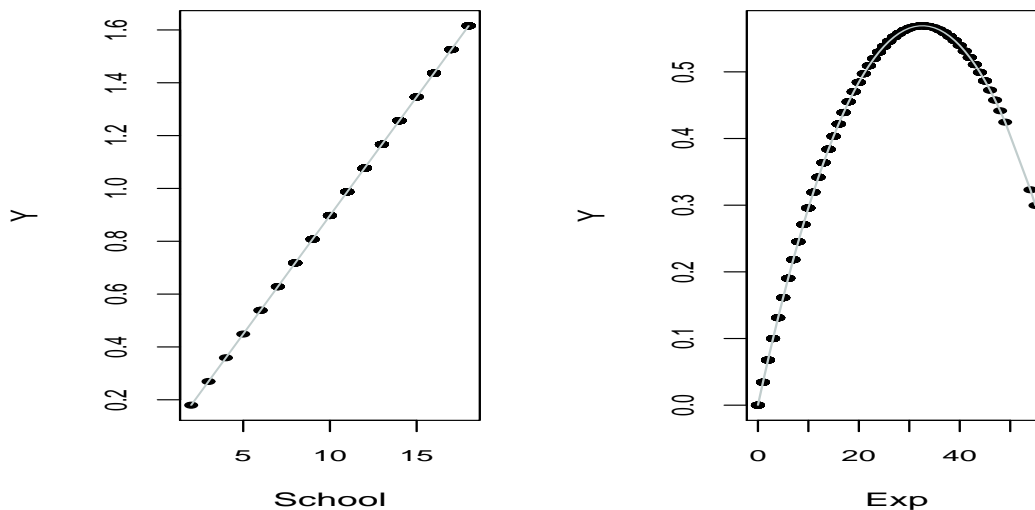
در این زیر بخش، برای بررسی روابط بین لگاریتم درآمد و متغیرهای توضیحی از مدل پارامتری (۳.۱) استفاده می‌کنیم. این مدل، بیش از هر مدل دیگری در اقتصاد تجربی توسط محققین استفاده شده است. طبق نظریه سرمایه انسانی، اثر سابقه کار بر دستمزد، یا نیمرخ دستمزد- سابقه کار، یک تابع مقعر است، چون مراحل اولیه شروع کار حرفه‌ای در بازار با افزایش دستمزد و انباشت سریع سرمایه انسانی همراه است و در یک نقطه‌ای در میانسالی به اوج می‌رسد. بعد از آن، میزان درآمد به دلیل کاهش ساعات کار و انگیزه برای سرمایه‌گذاری کاهش می‌یابد. این دلیل استفاده از  $Exp$  و  $Exp^2$  در مدل است. برای به‌دست آوردن مدلی که با نظریه سرمایه انسانی مطابقت داشته باشد، باید علامت  $\beta_2$  مثبت و  $\beta_3$  منفی باشد.

ضرایب مدل (۳.۱) را با استفاده از روش کمترین توان‌های دوم و براساس زیرنمونه‌ای به حجم  $n = 534$  از نظرسنجی جمعیت جاری سال ۱۹۸۵ که توسط اداره آمار ایالات متحده امریکا ارائه شده است، برآورد شده و نتایج در جدول ۱.۱ ارائه شده است. این داده‌ها در بسته "AER" نرم افزار R با نام "CPS۱۹۸۵" قابل دسترسی است.

<sup>4</sup>Curse of dimensionality

جدول ۱.۱: نتایج برآورد کمترین توان‌های دوم در مثال ۱.۱.۱

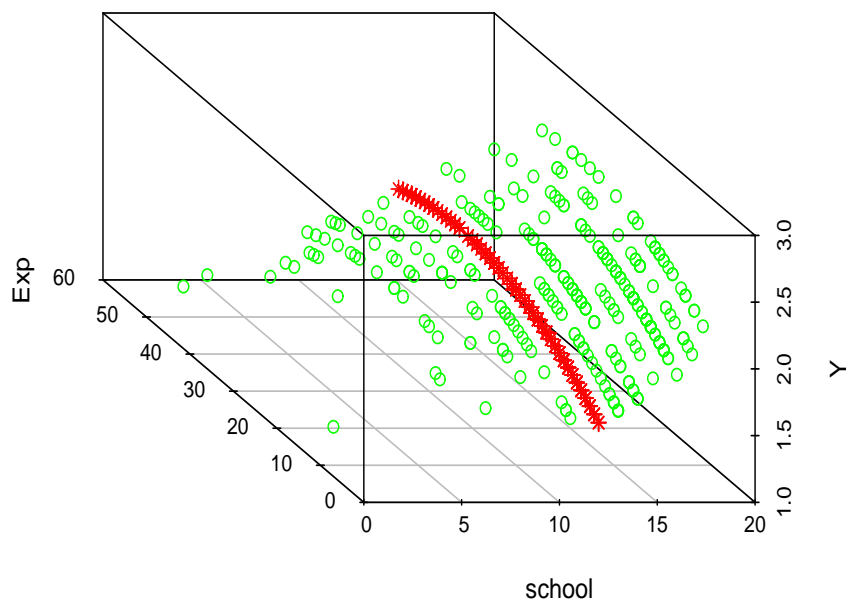
متغیر وابسته: لگاریتم درآمد			
متغیر	ضرایب	خطای معیار	آماره t
ثابت	۰/۵۲۰۲	۰/۱۲۳۶	۴/۲۰۹
<i>School</i>	۰/۰۸۹۸	۰/۰۰۸۳	۱۰/۷۸۸
<i>Exp</i>	۰/۰۳۴۹	۰/۰۰۵۶	۶/۱۸۵
<i>Exp</i> <sup>۲</sup>	-۰/۰۰۰۵	۰/۰۰۰۱	-۴/۳۰۷



شکل ۱.۱: قاب سمت چپ اثر متغیر تحصیل و قاب سمت راست اثر متغیر سابقه کار بر پاسخ را نشان می‌دهد.

توجه داشته باشید که ضرایب برآوردشده *Exp* و *Exp*<sup>۲</sup> دارای همان علامت پیش‌بینی‌شده توسط نظریه سرمایه انسانی است. شکل ۱.۱ نیمرخ دستمزد-تحصیل (*School* در مقابل *School* ۰/۰۸۹۸) و نیمرخ دستمزد-سابقه کار (نمودار *Exp* در مقابل *Exp* ۰/۰۰۰۵ - ۰/۰۳۴۹) را نشان می‌دهد. شکل ۱.۱ پیش‌بینی نظریه سرمایه انسانی درباره مقعر بودن نیمرخ دستمزد-سابقه کار را تایید می‌کند.

رویه رگرسیونی برآوردشده در شکل ۲.۱ نشان داده شده است. در این شکل، محور عمودی مقادیر تابع رگرسیونی برآوردشده و محورهای افقی متغیرهای *School* و *Exp* را نشان می‌دهند. همه منحنی‌های عضو این رویه در جهت *Exp* شبیه نمودار ۱.۱ (قاب سمت راست) و در جهت *School* شبیه نمودار ۱.۱ (قاب سمت چپ) هستند. برای درک بهتر این نمودار سه بعدی، نیمرخ دستمزد-سابقه کار را، با ثابت گرفتن تحصیل در ۱۲ سالگی، نشان داده‌ایم، یعنی منحنی قرمز رنگ نیمرخ دستمزد-سابقه کار را برای فارغ‌التحصیلان دبیرستان نشان می‌دهد. نتایج شکل ۲.۱ با شکل ۱.۱ همخوانی دارند.



شکل ۲.۱: تابع رگرسیونی برآورد شده به صورت پارامتری. منحنی قرمز رنگ، نیمرخ دستمزد- سابقه کار را برای فارغ التحصیلان دبیرستان نشان می دهد.

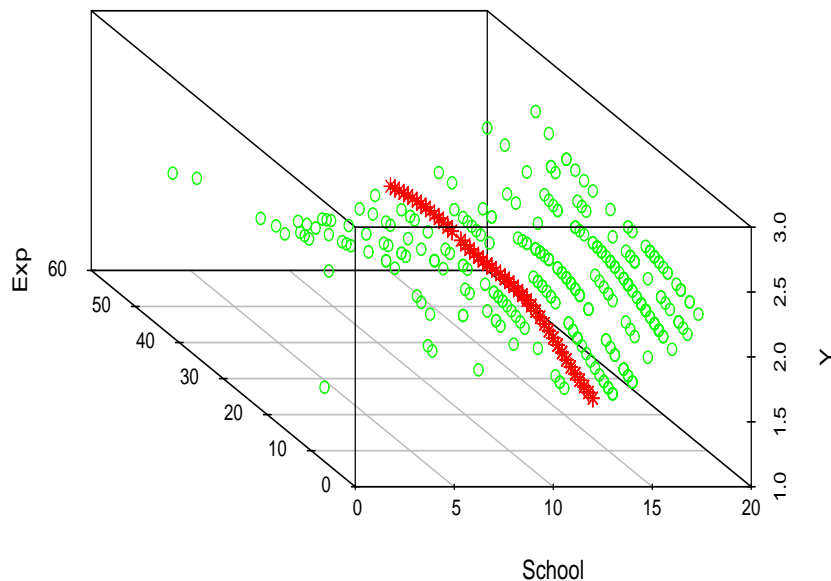
### ۲.۱.۱ رگرسیون ناپارامتری

در رگرسیون پارامتری رابطه بین متغیرهای پیشگو و متغیر پاسخ از طریق یک تابع پارامتری معین تعریف می شود. بنابراین برای برآورد یک مدل پارامتری نیاز به اطلاعات پیشین درباره ساختار داده ها است. با افزایش بعد داده ها و پیچیدگی مسائل مورد مطالعه، تعیین شکل درست مدل های پارامتری چندان ساده نیست. تعیین نادرست شکل مدل های پارامتری ممکن است منجر به آریبی برآورد و استنباط نادرست شود. در مقابل، مدل های رگرسیونی ناپارامتری، بدون نیاز به اطلاعات پیشین، انعطاف پذیری زیادی را برای شناخت ساختار داده ها فراهم آورده اند.

فرض کنید

$$E(Y|School, Exp) = m(School, Exp) \quad (۵.۱)$$

که  $m(\cdot)$  یک تابع هموار است. هدف برآورد تابع  $m(\cdot)$  بر اساس داده ها است. یک روش رایج در رگرسیون ناپارامتری میانگین وزنی موضعی است. فرض کنید  $x = (s, e)$  مقدار دلخواهی باشد که  $(School, Exp)$  می تواند اختیار کند. در روش میانگین وزنی موضعی، برای برآورد تابع  $m(\cdot)$  در نقطه  $x_0 = (s_0, e_0)$ ، به آن مقادیری از  $Y$  که متناظر با مقادیر  $x$  نزدیک تر به نقطه  $x_0$  باشند، وزن بیشتر و



شکل ۳.۱: تابع رگرسیونی برآوردشده با روش ناپارامتری. منحنی قرمز رنگ، نیمرخ دستمزد-سابقه کار را برای فارغ‌التحصیلان دبیرستان نشان می‌دهد.

به مقادیر  $Y$  متناظر با نقاط دورتر وزن کمتری نسبت داده می‌شود. نادارایا و واتسون<sup>۵</sup> (۱۹۶۴) پیشنهاد دادند که وزن‌ها در رگرسیون ناپارامتری می‌توانند به وسیله یک تابع متقارن تک‌مدی حول صفر که مقادیر آن در دو طرف صفر توسط یک پارامتر مقیاس کنترل می‌شود، انتخاب شود. نخستین نامزد آن‌ها برای چنین توابعی که هسته نامیده شدند، توابع چگالی احتمال بودند. سپس با استفاده از توابع کراندار هسته، برآوردگر زیر را پیشنهاد دادند

$$\hat{m}(x_0) = \sum_{i=1}^n W_i Y_i, \quad W_i = \frac{K_h(x_0, x_i)}{\sum_{i=1}^n K_h(x_0, x_i)} \quad (۶.۱)$$

که در آن  $K_h(\cdot)$  یک تابع هسته بوده و پارامتر  $h$  پهنای باند<sup>۶</sup> است. برآوردگر (۶.۱) به هموارساز هسته<sup>۷</sup> معروف است. شکل وزن‌ها توسط تابع هسته و مقادیر آن‌ها بر اساس فاصله از نقطه  $x_0$  و پهنای باند  $h$  تعیین می‌شود. بزرگ بودن پهنای باند باعث می‌شود که به نقاط نزدیک‌تر وزن‌های بزرگتری داده شود. ساده‌ترین هسته، هسته یکنواخت و پرکاربردترین آن‌ها، هسته نرمال می‌باشد. برآورد ناپارامتری تابع  $m(\cdot)$  با استفاده از هموارساز (۶.۱) با هسته نرمال در شکل ۳.۱ آمده است.

<sup>5</sup>Nadaraya and Watson

<sup>6</sup>Bandwidth

<sup>7</sup>Kernel smoothing



سایر روش‌ها برای برازش مدل ناپارامتری عبارتند از: رگرسیون اسپلاین، هموارسازی اسپلاین، اسپلاین رگرسیونی تطبیقی چندمتغیره (MARS) و رگرسیون موضعی. برای مطالعه بیشتر به هستی<sup>۸</sup> و همکاران (۲۰۰۹) و گرین و سیلورمن<sup>۹</sup> (۱۹۹۳) مراجعه کنید.

### ۳.۱.۱ رگرسیون نیمه پارامتری

اگرچه رگرسیون ناپارامتری از انعطاف‌پذیری بالایی برخوردار است، اما افزایش بعد فضای متغیرهای توضیحی، استفاده از مدل‌های ناپارامتری را با مشکل مواجه می‌کند. استون<sup>۱۰</sup> (۱۹۸۲) نشان داد که اگر مشتق تابع  $m(\cdot)$  تا مرتبه  $r$  موجود باشد، سریع‌ترین نرخ همگرایی قابل حصول برآوردگر  $\hat{m}(\cdot)$  به  $m(\cdot)$  برابر با  $(\log n/n)^{r/(2r+p)}$  است که در آن  $p$  بعد متغیرهای توضیحی است. این نرخ همگرایی نقش بعد متغیرهای توضیحی را در کارایی برآوردگر نشان می‌دهد، به عبارت دیگر، افزایش  $p$  منجر به کاهش عملکرد برآوردگر  $\hat{m}(\cdot)$  می‌شود. این مسئله به "مشقت بعد" معروف است و به دلیل تنگ بودن داده‌ها در فضای با ابعاد بالاتر اتفاق می‌افتد. برای اطلاعات بیشتر به گیننس<sup>۱۱</sup> (۲۰۱۱) مراجعه کنید. برای اجتناب از اریبی زیاد مدل‌های پارامتری و مشقت بعد مدل‌های ناپارامتری می‌توان از مدل‌های نیمه پارامتری استفاده کرد. برای شرح روش نیمه پارامتری، دوباره مدل درآمد سرمایه انسانی (۳.۱) را در نظر بگیرید. فرض کنید تابع رگرسیونی به شکل مدل جمعی زیر باشد

$$E(Y|School, Exp) = \alpha + g_1(School) + g_2(Exp) \quad (۷.۱)$$

که  $g_1(\cdot)$  و  $g_2(\cdot)$  توابع هموار نامعلوم و  $\alpha$  پارامتر نامعلوم است. این مدل ساختار جمعی رگرسیون پارامتری را با انعطاف‌پذیری رگرسیون ناپارامتری ترکیب می‌کند. در این مدل هیچ فرض محدود کننده قوی برای شکل توابع تعیین کننده اثر متغیرهای توضیحی  $School$  و  $Exp$  بر لگاریتم درآمد در نظر گرفته نمی‌شود. بنابراین برای برآورد توابع نامعلوم  $g_1(\cdot)$  و  $g_2(\cdot)$  باید از روش‌های ناپارامتری استفاده شود. شکل ۴.۱ برآوردهای توابع  $g_1(\cdot)$  و  $g_2(\cdot)$  را در مقابل نیمرخ دستمزد-تحصیل و دستمزد-سابقه کار برآورد شده با روش پارامتری نشان می‌دهد. در این شکل، نیمرخ‌های برآورد شده با روش پارامتری با خطوط نازک و برآوردهای  $g_1(\cdot)$  و  $g_2(\cdot)$  با خطوط ضخیم نشان داده شده است. با توجه به نمودار، شباهت زیادی بین این دو برآورد وجود دارد، اما در نقاط پایانی تفاوت مشاهده می‌شود. همچنین، رویه رگرسیونی برآورد شده با روش نیمه پارامتری در شکل ۵.۱ نشان داده شده است. این رویه برآورد شده همخوانی خوبی بین برآوردهای پارامتری و برازش مدل جمعی نشان می‌دهد.

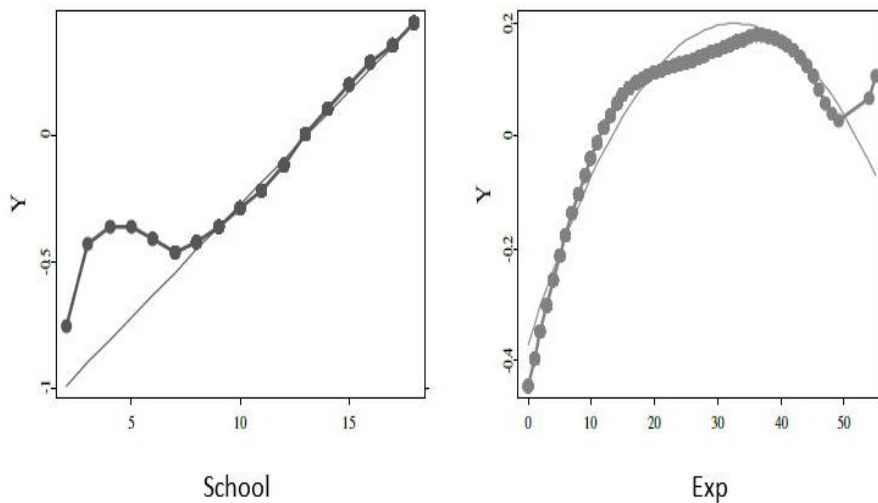
لذا می‌توانیم نتیجه بگیریم که در این مثال خاص، مدل پارامتری توسط روش‌های ناپارامتری و نیمه پارامتری تایید می‌شود. بنابراین روش‌های ناپارامتری و نیمه پارامتری را می‌توان برای بررسی کفایت مدل پارامتری نیز به کار برد.

<sup>۸</sup>Hasti

<sup>۹</sup>Green and Silverman

<sup>۱۰</sup>Stone

<sup>۱۱</sup>Geenens



شکل ۴.۱: برازش مدل جمعی در مقابل برازش مدل پارامتری. قاب سمت چپ اثر متغیر تحصیل و قاب سمت راست اثر متغیر سابقه کار بر پاسخ را نشان می‌دهد.

حال با دقت بیشتری به مدل‌های (۵.۱) و (۷.۱) نگاه کنید. مشاهده می‌کنید که در مدل (۵.۱) لازم است یک تابع نامعلوم شامل دو متغیر را برآورد کنیم، اما در مدل (۷.۱) باید دو تابع نامعلوم برآورد شود که هریک از این دو تابع شامل یک متغیر است. این بدین معنی است که با برازش مدل (۷.۱) توانسته‌ایم بعد مسئله برآورد را کاهش دهیم. اگرچه همه محققین بر این باورند که برازش مدل‌های جمعی مانند (۷.۱) باعث اعمال کاهش بعد بر مدل‌های کاملاً ناپارامتری می‌شود، اما برخی از آن‌ها مدل (۷.۱) را به‌عنوان مدل نیمه‌پارامتری قبول ندارند. زیرا در این مدل هیچ پارامتری، به جز عرض از مبدأ، برای برآورد وجود ندارد. تاکنون مدل‌های رگرسیونی نیمه‌پارامتری متعددی توسط افراد مختلف معرفی شده‌اند. رایج‌ترین مدل‌های نیمه‌پارامتری عبارتند از: مدل جمعی خطی-جزئی<sup>۱۲</sup>، مدل ضریب متغیر خطی-جزئی<sup>۱۳</sup> و مدل تک‌شاخص<sup>۱۴</sup> که در ادامه به معرفی این مدل‌ها خواهیم پرداخت.

### مدل جمعی خطی-جزئی

مدل جمعی خطی-جزئی به صورت

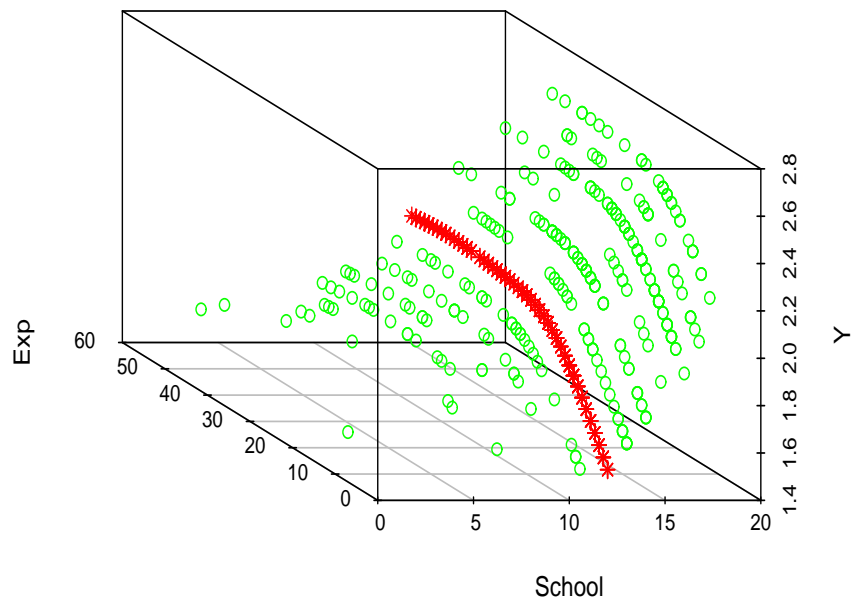
$$Y = \sum_{j \in S_1} \beta_j X_j + \sum_{j \in S_2} f_j(X_j) + \varepsilon \quad (۸.۱)$$

تعریف می‌شود که در آن  $f_j(\cdot)$  ها توابع هموار نامعلوم و  $S_1$  و  $S_2$  دو زیر مجموعه مجزا از  $\{1, \dots, p\}$  می‌باشند. این مدل تفسیرپذیری مدل‌های خطی و انعطاف‌پذیری مدل‌های ناپارامتری را با هم ترکیب می‌کند و یکی از رایج‌ترین مدل‌های نیمه‌پارامتری است. هاردل و همکاران (۲۰۰۰) بررسی دقیقی از این

<sup>12</sup>Partially linear additive models

<sup>13</sup>Partially linear varying coefficient model

<sup>14</sup>Single index model



شکل ۵.۱: تابع رگرسیونی برآوردشده با روش نیمه پارامتری. منحنی قرمز رنگ، نیمرخ دستمزد- سابقه کار را برای فارغ التحصیلان دبیرستان نشان می دهد.

مدل نیمه پارامتری ارائه داده اند. مدل جمعی خطی- جزئی دارای کاربردهای متعددی است. انگل<sup>۱۵</sup> و همکاران (۱۹۸۶) از نخستین افرادی بودند که این مدل را در نظر گرفتند. آن ها از این مدل برای بررسی رابطه بین دما و مصرف برق استفاده کردند. در ادامه به بیان چند مثال واقعی می پردازیم که محققین مدل جمعی خطی- جزئی را به کار برده اند.

**مثال ۲.۱.۱.** هاردل و همکاران (۲۰۰۴) از مدل جمعی خطی- جزئی برای بررسی تأثیر عوامل مختلف بر عرضه نیروی کار زن در شرق آلمان استفاده کردند. در این داده ها، عرضه نیروی کار زن با تعداد ساعت های کاری هفتگی ( $Y$ ) اندازه گیری می شود. مجموعه داده های مورد بررسی یک زیر نمونه از پنل اقتصادی اجتماعی آلمان (GSOEP) شامل  $n = 607$  زن است که دارای شغل و در حال زندگی با همسر خود هستند. متغیرهای توضیحی عبارتند از: نرخ بیکاری در ایالتی که زن در آن زندگی می کند ( $U_1$ )، زن دارای فرزند کمتر از ۱۶ ساله است، بله=۱ و خیر=۰ ( $U_2$ )، سن برحسب سال ( $T_1$ )، دستمزد هر ساعت ( $T_2$ )، شاخص اعتبار شغل ( $T_3$ )، اجاره ماهانه آپارتمان یا محل اقامت ( $T_4$ )، تعداد سال های تحصیل ( $T_5$ )، درآمد خالص ماهانه شوهر ( $T_6$ ).

<sup>15</sup>Engle

هاردل و همکاران (۲۰۰۴) مدل جمعی خطی-جزئی زیر را برای تحلیل این مجموعه داده به کار بردند

$$Y = \alpha + \sum_{j=1}^2 \beta_j U_j + \sum_{j=1}^6 g_j(T_j) + \varepsilon.$$

آن‌ها نشان دادند که میانگین متغیر پاسخ به‌طور خطی با متغیرهای  $U_1$  و  $U_2$  در ارتباط است، در حالی که نحوه ارتباط آن با سایر متغیرهای توضیحی نامعلوم است.

**مثال ۳.۱.۱.** انگل و همکاران (۱۹۸۶) از مدل جمعی خطی-جزئی برای بررسی رابطه بین میزان برق مصرفی ماهیانه ( $Y$ ) با متغیرهای قیمت برق ( $X_1$ )، درآمد ( $X_2$ ) و متوسط دمای هوا ( $t$ ) استفاده کردند. آن‌ها متغیر پاسخ  $Y$  را به صورت مجموع یک تابع هموار  $g$  از دمای ماهیانه  $t$  و متغیرهای  $X_1$  و  $X_2$  و همچنین متغیرهای نشانگر  $X_3, \dots, X_{13}$ ، که معرف ماه‌های میلادی هستند، به صورت

$$Y = \sum_{j=1}^{13} \beta_j X_j + g(t) + \varepsilon$$

مدل‌سازی کردند. آن‌ها دریافتند که متغیر پاسخ  $Y$  دارای رابطه خطی با متغیرهای  $X_1$  و  $X_2$  و غیر خطی با متغیر  $t$  است.

**مثال ۴.۱.۱.** اشمالنسی و استوکر<sup>۱۶</sup> (۱۹۹۹) مدل جمعی خطی-جزئی را برای بررسی رابطه بین مصرف گازوئیل خانواده با متغیرهایی مانند درآمد، سن، تعداد افراد خانوار، نوع خانه و محل اقامت در ایالات متحده آمریکا به کار بردند. آن‌ها دریافتند که لگاریتم مصرف گازوئیل بر حسب گالن دارای رابطه غیرخطی با لگاریتم متغیرهای درآمد و سن، و رابطه خطی با لگاریتم تعداد افراد خانوار، نوع خانه و محل اقامت است.

**مثال ۵.۱.۱.** هوانگ<sup>۱۷</sup> و همکاران (۲۰۱۲a) یک مجموعه داده مربوط به دیابت را، که توسط ویلمز<sup>۱۸</sup> (۱۹۹۷) گزارش شده است، مورد بررسی قرار دادند. در این داده‌ها، ۱۶ متغیر بر روی ۴۰۳ نفر اندازه‌گیری شده است تا میزان شیوع چاقی، دیابت و سایر عوامل خطر بیماری قلبی و عروقی را در ویرجینیای مرکزی بررسی کنند. دیابت نوع دوم با چاقی همراه است. مقدار هموگلوبین گلیکوزیله<sup>۱۹</sup> بیشتر از ۷/۰ معمولاً به‌عنوان یک تشخیص مثبت از این بیماری در نظر گرفته می‌شود.

هوانگ و همکاران (۲۰۱۲a) هموگلوبین گلیکوزیله را به‌عنوان متغیر پاسخ و ۱۵ متغیر دیگر را به‌عنوان متغیرهای توضیحی در نظر گرفتند. این ۱۵ متغیر شامل کلسترول ( $X_1$ )، گلوکز تثبیت شده ( $X_2$ )، لیپوپروتئین با چگالی بالا ( $X_3$ )، نسبت کلسترول به لیپوپروتئین با چگالی بالا ( $X_4$ )، مکان ( $X_5$ )، سن ( $X_6$ )، جنسیت ( $X_7$ )، قد ( $X_8$ )، وزن ( $X_9$ )، اسکلت ( $X_{10}$ )، اولین فشار خون سیستولیک ( $X_{11}$ )، اولین فشار خون دیاستولیک ( $X_{12}$ )، دور کمر ( $X_{13}$ )، دور ران ( $X_{14}$ )، و زمان طراحی آزمایشگاه ( $X_{15}$ )

<sup>16</sup>Schmalensee and Stoker

<sup>17</sup>Huang

<sup>18</sup>Willems

<sup>19</sup>Glycosylated hemoglobin

است. در بین ۱۵ متغیر توضیحی، ۳ متغیر مکان، جنسیت و اسکلت متغیرهای طبقه‌ای و ۱۲ متغیر دیگر پیوسته هستند. آن‌ها با به‌کارگیری روشی موسوم به رگرسیون تعقیب نیمه‌پارامتری<sup>۲۰</sup> دریافتند که علاوه بر ۳ متغیر طبقه‌ای، متغیرهای پیوسته  $X_2, X_9, X_{11}, X_{12}, X_{13}, X_{14}$  دارای اثر خطی و سایر متغیرها دارای رابطه غیرخطی با متغیر پاسخ هستند.

### مدل ضریب متغیر خطی - جزئی

مدل ضریب متغیر تعمیمی از مدل خطی است که برای افزایش انعطاف‌پذیری مدل، ضرایب رگرسیونی با توابع ناپارامتری هموار جایگزین می‌شوند (کولند<sup>۲۱</sup> و همکاران، ۱۹۹۱؛ هستی و تیبشیرانی<sup>۲۲</sup>، ۱۹۹۳). صورت کلی این مدل به صورت

$$Y = \sum_{j=1}^p \beta_j(U) X_j + \varepsilon \quad (9.1)$$

است که  $U$  متغیر شاخص<sup>۲۳</sup> یا تعدیل‌کننده<sup>۲۴</sup>، و  $\beta(U) = (\beta_1(U), \dots, \beta_p(U))^T$  بردار توابع ضرایب، شامل توابع هموار نامشخص از  $U$ ، است و به کمک نمونه برآورد می‌شوند. در این مدل، متغیر شاخص معمولاً سن یا زمان در نظر گرفته می‌شود و هر یک از متغیرهای پیشگو، برای مقادیر داده‌شده متغیر  $U$ ، دارای رابطه خطی با میانگین متغیر پاسخ می‌باشند. در عمل، برخی از توابع ضرایب در مدل (۹.۱) ثابت و بقیه ضرایب تابعی از متغیر شاخص  $U$  هستند. در این حالت خاص می‌توان مدل را به صورت

$$Y = \sum_{j \in S_1} \beta_j X_j + \sum_{j \in S_2} \beta_j(U) X_j + \varepsilon \quad (10.1)$$

نوشت. مدل (۱۰.۱) را مدل ضریب متغیر خطی - جزئی گویند. در نظر گرفتن ضرایب ثابت به صورت تابعی از متغیر شاخص باعث افزایش آریبی و کاهش دقت برآورد می‌شود. بنابراین، مدل ضریب متغیر خطی - جزئی در تحلیل داده‌های واقعی بسیار مفید است.

**مثال ۶.۱.۱.** ژائو<sup>۲۵</sup> و همکاران (۲۰۱۴) با استفاده از مدل ضریب متغیر به بررسی رابطه بین سطح بتاکاروتن پلاسما و متغیرهایی مانند جنسیت، وضعیت استعمال سیگار، شاخص توده بدنی، وضعیت مصرف ویتامین، تعداد کالری، میزان چربی، میزان فیبر، تعداد نوشیدنی‌های الکلی، کلسترول و سن پرداختند. ابتدا متغیرهای طبقه‌ای وضعیت استعمال سیگار و مصرف ویتامین را با متغیرهای ساختگی<sup>۲۶</sup> جایگزین کردند. سپس با در نظر گرفتن سن به‌عنوان متغیر شاخص نشان دادند که متغیر سن دارای اثر غیر خطی و متغیرهای جنسیت، وضعیت استعمال سیگار، وضعیت مصرف ویتامین و نوشیدنی‌های الکلی دارای اثر خطی روی پاسخ هستند و اثر سایر متغیرهای توضیحی تابعی از سن است.

<sup>20</sup>Semiparametric regression pursuit

<sup>21</sup>Cleveland

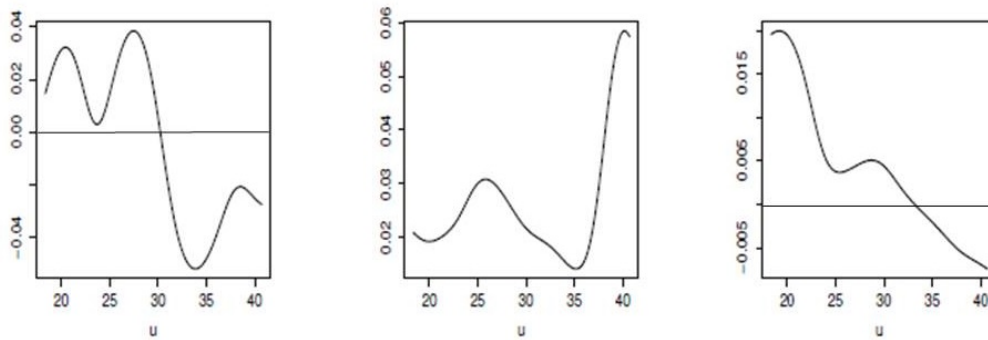
<sup>22</sup>Tibshirani

<sup>23</sup>Index variable

<sup>24</sup>Modifying predictor

<sup>25</sup>Zhao

<sup>26</sup>Dummy variables



شکل ۶.۱: توابع ضرایب برآورد شده مثال ۷.۱.۱. نمودارها، به ترتیب، از چپ به راست  $\hat{\alpha}_1(u)$ ،  $\hat{\alpha}_2(u)$  و  $\hat{\alpha}_3(u)$  را نشان می‌دهند.

**مثال ۷.۱.۱.** لی<sup>۲۷</sup> و همکاران (۲۰۱۶) از مدل ضریب متغیر برای مدل‌بندی داده‌های وزن نوزاد هنگام تولد استفاده کردند. بدین منظور، ۶۴ مادر که در یک بیمارستان خاص از جمهوری چک در فاصله زمانی نوامبر ۲۰۰۸ تا مارس ۲۰۰۹ زایمان کرده بودند، مورد بررسی قرار گرفتند. در این راستا، وزن نوزاد برحسب کیلوگرم (BW) همراه با متغیرهایی مانند سن مادر (MOA)، مدت زمان بارداری (GEA) و شاخص توده بدنی مادر (BMI) ثبت شد. علاوه بر این، سطح بیان ۵۸۶۹ ژن این مادران نیز در دسترس است. هدف شناسایی ژن‌های مؤثر بر وزن نوزاد است. آن‌ها با در نظر گرفتن سن مادر به‌عنوان متغیر شاخص، مدل ضریب متغیر خطی-جزئی زیر را برای برازش داده‌های وزن هنگام تولد به صورت زیر به کار بردند

$$BW = \alpha_1(U) + \alpha_2(U)GEA + \alpha_3(U)BMI + \sum_{j=1}^{5869} jGE_j + \varepsilon$$

که  $U = MOA$  و  $jGE_j$  بیان  $j$ -امین ژن است. پس از برازش مدل فوق، ۱۰ ژن مؤثر شناسایی شدند. توابع ضرایب برآوردشده در شکل ۶.۱ نشان داده شده‌اند. با توجه به این شکل، می‌توان دید که تا قبل از ۳۰ سالگی مادران، سن مادر دارای تاثیر مثبت بر وزن نوزاد و برای سن بیش از ۳۰ سال دارای اثر منفی است. مدت زمان بارداری دارای اثر مثبت بر وزن نوزاد است و میزان تاثیر بعد از ۳۵ سالگی به سرعت افزایش می‌یابد. همچنین، متغیر BMI دارای اثر مثبت تا سن ۳۵ سالگی، و بعد از ۳۵ سالگی دارای اثر منفی بر وزن نوزاد است.

### مدل تک شاخص

مدل رگرسیونی تک‌شاخص به شکل

$$Y = g(\mathbf{X}^T \beta) + \varepsilon \quad (11.1)$$

تعریف می‌شود که  $g$  یک تابع نامعلوم است. تابع  $g$  را معمولاً تابع پیوند و کمیت  $\mathbf{X}^T \beta$  را شاخص گویند. مدل (۱۱.۱) بسیاری از مدل‌های پارامتری پرکاربرد را به‌عنوان یک حالت خاص شامل می‌شود. اگر  $g$

تابع همانی باشد، پس (۱۱.۱) یک مدل خطی است. همچنین اگر  $g$  تابع توزیع تجمعی نرمال یا تابع توزیع لجستیک باشد، مدل فوق به ترتیب به مدل پروبیت یا مدل لوجیت تبدیل می‌شود. علاوه بر این، مدل‌های توبیت و پواسن نیز حالت خاصی از مدل تک‌شاخص می‌باشند. برای اطلاعات بیشتر به هاردل و همکاران (۱۹۹۳) و هاردل و استوکر (۱۹۸۹) مراجعه کنید.

## ۲.۱ اهداف رساله

همان‌طور که قبلاً اشاره شد، مدل‌های نیمه‌پارامتری علاوه بر کاهش آریبی مدل‌های پارامتری، با مشکل مشقت بعد مدل‌های ناپارامتری نیز مقابله می‌کنند. یکی از این مدل‌های پرکاربرد، مدل جمعی خطی-جزئی

$$Y = \sum_{j \in S_1} \beta_j X_j + \sum_{j \in S_2} f_j(X_j) + \varepsilon \quad (12.1)$$

است که در آن  $S_1$  و  $S_2$  دو زیر مجموعه مجزا از  $\{1, \dots, p\}$  است. فرض کنید تعداد متغیرهای توضیحی مؤثر بر پاسخ برابر با  $p$  باشد. در این صورت، مجموعه  $S_1 \cup S_2$  دارای اندازه  $p_1$  است. به عنوان مثال، اگر  $p_1 = 5$ ،  $p = 10$ ،  $S_1 = \{1, 2, 3\}$  و  $S_2 = \{4, 5\}$  باشد، آنگاه ۳ متغیر توضیحی اثر خطی و ۲ متغیر اثر غیرخطی در مدل بندی متغیر پاسخ دارند و ۵ متغیر دیگر بی‌اثر هستند. برای استفاده از مدل (۱۲.۱) لازم است اجزای خطی و غیرخطی مدل مشخص باشند. در تحلیل داده‌های واقعی، از پیش تعیین کردن مؤلفه‌های خطی و غیرخطی چندان ساده به نظر نمی‌رسد، به‌ویژه زمانی که تعداد متغیرهای توضیحی زیاد است. بنابراین، علاوه بر تشخیص مؤلفه‌های غیرصفر، شناسایی مؤلفه‌های خطی و غیرخطی نیز بسیار حائز اهمیت است. تاکنون بسیاری از محققین، با تکیه بر روش‌های جریمه، به موضوع انتخاب متغیر و تشخیص ساختار<sup>۲۸</sup> در مدل‌های جمعی خطی-جزئی پرداخته‌اند. از آن جمله می‌توان به ژانگ<sup>۲۹</sup> و همکاران (۲۰۱۱)، هوآنگ و همکاران (۲۰۱۲)، لیان<sup>۳۰</sup> (۲۰۱۲b)، لیان (۲۰۱۲c)، لیان و همکاران (۲۰۱۲)، لیان و همکاران (۲۰۱۵) اشاره کرد. لازم به ذکر است که منظور از تشخیص ساختار در مدل (۱۲.۱)، شناسایی متغیرهایی است که به صورت خطی یا غیرخطی وارد مدل می‌شوند. به عبارت دیگر، هدف تشخیص مجموعه‌های  $S_1$  و  $S_2$  است.

از طرفی، با پیشرفت فناوری و افزایش توان ذخیره‌سازی کامپیوترها، محققین می‌توانند با هزینه نسبتاً کمی داده‌هایی با "بعد بسیار بالا"<sup>۳۱</sup> در اختیار داشته باشند. تحلیل این نوع داده‌ها فرصت‌ها و چالش‌های جدیدی را پیش‌روی آماردانان قرار داده است. به دلیل وجود تعداد زیادی از متغیرها در مدل‌های با بعد بسیار بالا، تفسیر این مدل‌ها بسیار مشکل است. لذا مسئله انتخاب متغیر نقش بسیار مهمی را در تحلیل این نوع داده‌ها ایفا می‌کند.

<sup>28</sup>Structure identification

<sup>29</sup>Zhang

<sup>30</sup>Lian

<sup>31</sup>Ultrahigh dimensional

در این خصوص، مطالعات انجام شده فوق، که مبتنی بر توابع جریمه هستند، به دلایل هزینه محاسباتی، دقت آماری و ناپایداری الگوریتمی از کارایی کافی برخوردار نیستند (فن<sup>۳۲</sup> و همکاران، ۲۰۰۹). لذا شناسایی متغیرهای مهم و تشخیص نوع تاثیر آن‌ها در متغیر پاسخ به لحاظ خطی یا غیرخطی بودن، کماکان مورد توجه متخصصین علوم مذکور و آماردانان است.

در این رساله با تمرکز بر انتخاب متغیر و تشخیص ساختار در داده‌های با بعد بسیار بالا، ابتدا از یک روش غربالگری مستقل مطمئن برای کاهش بعد مدل استفاده می‌شود. این روش بر اساس همبستگی فاصله‌ای بین متغیرهای توضیحی و تابع توزیع حاشیه‌ای متغیر پاسخ تعریف می‌شود و انتظار می‌رود وقتی که متغیر پاسخ دم سنگین، چوله یا دارای مقادیر فرین باشد، عملکرد خوبی داشته باشد. سپس از دو تابع جریمه برای تشخیص همزمان متغیرهای مهم و مؤلفه‌های خطی در زیر مدل غربال شده در مرحله قبل استفاده می‌شود. علاوه بر این ویژگی غربالگری مطمئن را برای روش غربالگری معرفی شده اثبات کرده و با استفاده از شبیه‌سازی و تحلیل داده‌های واقعی عملکرد این روش را با روش‌های موجود مقایسه می‌کنیم.

در بعضی از مسائل، اثر برخی از متغیرهای توضیحی بر پاسخ به صورت خطی و اثر متغیرهای دیگر تابعی از متغیر شاخص است. بنابراین لازم است به جای مدل (۱۲.۱) از مدل ضریب متغیر خطی-جزئی (۱۰.۱) استفاده شود. مشابه مدل جمعی خطی-جزئی، استفاده از مدل (۱۰.۱) بر فرض معلوم بودن اثرات ثابت و متغیر استوار است. در عمل، افزاز این ضرایب، به دو جزء ثابت و متغیر، بسیار مشکل است. بنابراین لازم است روشی ارائه شود تا ضرایب صفر، ثابت غیر صفر و متغیر شناسایی شوند. مسئله انتخاب متغیر و تشخیص ساختار در مدل‌های ضریب متغیر خطی-جزئی تاکنون توسط افراد مختلف مورد بررسی قرار گرفته است. به عنوان مثال می‌توان به لنگ<sup>۳۳</sup> (۲۰۰۹)، هو و زیا<sup>۳۴</sup> (۲۰۱۲)، نه و کیلگوم<sup>۳۵</sup> (۲۰۱۲)، وانگ و سانگ<sup>۳۶</sup> (۲۰۱۳) اشاره کرد. با این حال، همه این روش‌ها بر رگرسیون میانگین یا تابع درست‌نمایی استوار هستند، که به شدت تحت تاثیر داده‌های پرت یا توزیع دم‌سنگین خطاها قرار می‌گیرند. بنابراین لازم است یک روش نیرومند برای انتخاب متغیر و تشخیص ساختار در مدل‌های ضریب متغیر خطی-جزئی ارائه شود. از طرفی، یائو<sup>۳۷</sup> و همکاران (۲۰۱۲) رگرسیون نما را معرفی کردند که، به جای میانگین، از نمای متغیر پاسخ به منظور مدل‌سازی استفاده می‌کند. بنابراین با توجه به ویژگی‌های شاخص نما، این روش در مقایسه با رگرسیون میانگین از نیرومندی بیشتری نسبت به داده‌های پرت و توزیع غیر نرمال خطاها برخوردار است.

در این رساله با استفاده از رگرسیون نما یک روش نیرومند برای انتخاب متغیر و تشخیص ساختار در مدل‌های ضریب متغیر خطی-جزئی معرفی می‌شود. عملکرد این روش با استفاده از مطالعه شبیه‌سازی و مثال واقعی مورد ارزیابی قرار می‌گیرد. علاوه بر این، ویژگی‌های مجانبی برآوردگر ارائه شده نیز ثابت می‌شود.

<sup>32</sup>Fan

<sup>33</sup>Leng

<sup>34</sup>Hu and Xia

<sup>35</sup>Noh and Keilegom

<sup>36</sup>Wang and Song

<sup>37</sup>Yao



## ۳.۱ ساختار رساله

- در فصل ۲ به معرفی برخی از روش‌های انتخاب متغیر مبتنی بر جریمه می‌پردازیم. در ادامه، با توجه به ناکارآمدی روش‌های جریمه برای انتخاب متغیر در داده‌های با بعد بسیار بالا، به توصیف برخی از روش‌های غربالگری برای کاهش بعد در این نوع داده‌ها خواهیم پرداخت. مسئله انتخاب متغیر با استفاده از توابع جریمه مستلزم حل یک مسئله بهینه‌سازی است. در انتهای این فصل حل این مسئله بهینه‌سازی مورد بحث قرار می‌گیرد.
- در فصل ۳، یک روش دو مرحله‌ای برای انتخاب متغیر و تشخیص ساختار در مدل‌های جمعی خطی-جزئی معرفی می‌شود. در مرحله اول، از یک روش غربالگری آزاد-مدل برای کاهش بعد فضای متغیرهای توضیحی استفاده می‌شود. این روش براساس ضریب همبستگی فاصله‌ای بین تابع توزیع متغیر پاسخ و متغیرهای توضیحی تعریف می‌شود. سپس در مرحله دوم، از دو تابع جریمه برای انتخاب متغیر و تشخیص ساختار استفاده می‌شود.
- در فصل ۴، مسئله انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی-جزئی مورد بررسی قرار می‌گیرد. در این فصل با ترکیب رگرسیون نما و توابع جریمه، یک روش نیرومند برای برازش مدل ضریب متغیر خطی-جزئی ارائه می‌شود.
- در بخش ضمیمه، اثبات قضایای مطرح‌شده در فصل‌های ۳ و ۴، مروری بر توابع اسپلاین و B-اسپلاین، گزیده‌ای از کدهای R و برخی از تعاریف و مفاهیم مورد نیاز در این رساله بیان می‌شوند.



## فصل ۲

# انتخاب متغیر در مدل‌های رگرسیونی

در این فصل، به معرفی رایج‌ترین روش‌های انتخاب متغیر پرداخته می‌شود. مطالب این فصل از منابع متعددی گرفته شده‌است که از آن جمله می‌توان به مونت‌گومری و همکاران (۲۰۱۲)، تیبشیرانی (۱۹۹۶)، فن و لی (۲۰۰۱)، هوانگ و همکاران (۲۰۱۲b) و فن و لیو (۲۰۰۸) اشاره کرد.

### ۱.۲ مقدمه

مدل رگرسیونی چندگانه زیر را در نظر بگیرید

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon \quad (1.2)$$

که در آن بردار متغیرهای توضیحی،  $\mathbf{X} = (X_1, \dots, X_p)^T$  بردار ضرایب  $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی و  $\varepsilon$  خطای تصادفی با میانگین صفر و واریانس ثابت است. در این فصل، بدون از دست دادن کلیت، فرض می‌شود عرض از مبدا برابر با صفر است. یک روش متداول برآورد ضرایب رگرسیونی در مدل (۱.۲)، روش کمترین توان‌های دوم است. هنگامی که بین متغیرهای توضیحی همخطی وجود دارد، استنباط بر اساس این روش می‌تواند گمراه‌کننده باشد. در صورت وجود همخطی، برآوردگر کمترین توان‌های دوم منجر به برآوردهایی با واریانس بزرگ می‌شود که دقت پیشگویی را کاهش می‌دهد. یکی از عوامل ایجاد همخطی، زیاد بودن تعداد متغیرهای توضیحی در مدل است. در مدل‌های با بعد بالا، مسئله همخطی بسیار شدید است، یعنی هر یک از متغیرهای توضیحی را می‌توان به صورت ترکیب خطی

از سایر متغیرها نوشت (جیمز<sup>۱</sup> و همکاران، ۲۰۱۳؛ فن و لی، ۲۰۰۸). روش پیشنهادی در این حالت، انتخاب زیر مجموعه‌ای از متغیرهای توضیحی برای ورود به مدل است. در این رهیافت علاقه‌مندیم زیرمجموعه‌ای از متغیرها را طوری انتخاب کنیم که رابطه بین متغیر پاسخ و متغیرهای توضیحی را به خوبی توصیف کند. یافتن زیر مجموعه‌ای مناسب از متغیرهای توضیحی برای ورود به مدل "انتخاب متغیر" نامیده می‌شود.

ساختن یک مدل رگرسیونی که تنها شامل زیرمجموعه‌ای از متغیرهای توضیحی است، دو واقعیت ناسازگار را در بر دارد: الف) علاقه‌مندیم که مدل، تا حد امکان، متغیرهای توضیحی بیشتری را شامل شود، به نحوی که اطلاعات موجود در این متغیرها بتواند در پیشگویی متغیر پاسخ مؤثر باشد. ب) می‌خواهیم مدل حتی الامکان دارای متغیرهای رگرسیونی کمتری باشد، چون با افزایش تعداد متغیرهای توضیحی در مدل، هزینه محاسبات افزایش می‌یابد و تفسیر مدل دشوارتر می‌شود. علاوه بر این، وجود تعداد زیادی متغیر توضیحی ممکن است باعث بیش‌برازش<sup>۲</sup> شود، یعنی مدل برازش شده برای داده‌های آموزشی<sup>۳</sup>، که برای ساخت مدل از آن‌ها استفاده شده است، مناسب باشد، اما برای پیشگویی مشاهدات آینده بسیار ضعیف باشد. بنابراین مسئله انتخاب متغیر امری بسیار مهم در تحلیل داده‌هاست. فرآیند به‌دست آوردن یک مدل که سازگاری و توافق بین این دو واقعیت است، انتخاب بهترین معادله رگرسیونی نامیده می‌شود.

به‌عنوان رهیافت دیگر، می‌توان از روش‌هایی مانند رگرسیون مؤلفه‌های اصلی<sup>۴</sup> یا رگرسیون کمترین توان‌های دوم جزئی<sup>۵</sup> برای کاهش بعد استفاده کرد. در این روش‌ها، ابتدا  $p$  ترکیب خطی ناهمبسته از متغیرهای توضیحی تولید کرده و سپس  $M$  تا از این ترکیب‌ها به عنوان جایگزین متغیرهای اصلی استفاده می‌شود که  $M < p$ . به‌وضوح، اگر  $M = p$ ، آنگاه برآورد به‌دست آمده با برآورد کمترین توان‌های دوم معادل است. استفاده از این مدل‌ها به منظور پیشگویی منطقی است، اما معمولاً تفسیر این مدل‌ها دشوار است (هستی و همکاران، ۲۰۰۹).

برای تعیین مدل پیشگو، لازم است پارامترهای مدل (ضرایب رگرسیونی) برآورد شوند. لذا فرآیندی که علاوه بر انتخاب متغیر و بهترین معادله رگرسیونی، برآورد ضرایب را نیز انجام دهد، مورد نظر است. روش کمترین توان‌های دوم جریمه‌شده<sup>۶</sup> یک روش کارا برای برآورد پارامترها است که همزمان انتخاب متغیر را نیز انجام می‌دهد.

در این فصل، ابتدا خواص برآوردگر کمترین توان‌های دوم و تاثیر همخطی بر عملکرد این برآوردگر بررسی می‌شود. سپس به معرفی رایج‌ترین روش‌های انتخاب متغیر می‌پردازیم. این روش‌ها شامل روش‌های انتخاب متغیر کلاسیک، روش کمترین توان‌های دوم جریمه‌شده و روش غربالگری<sup>۷</sup> هستند.

<sup>1</sup>James

<sup>2</sup>Overfitting

<sup>3</sup>Training data set

<sup>4</sup>Principal component regression

<sup>5</sup>Partial least squares regression

<sup>6</sup>Penalized least squares

<sup>7</sup>Screening

## ۲.۲ روش کمترین توان‌های دوم

فرض کنید  $\{(x_i, y_i)\}_{i=1}^n$  یک نمونه تصادفی از مدل (۱.۲) باشند. بنابراین مدل رگرسیون خطی چندگانه را می‌توان با نماد ماتریسی

$$Y = X\beta + \varepsilon \quad (۲.۲)$$

نوشت که در آن  $Y = (y_1, \dots, y_n)^T$  بردار متغیر پاسخ،  $X = (x_1, \dots, x_n)^T$  ماتریس طرح در اندازه  $n \times p$  با رتبه کامل ستونی،  $\beta$  بردار ضرایب رگرسیونی و  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$  بردار خطا است.

در رگرسیون خطی چندگانه فرض بر این است که رابطه بین متغیر پاسخ و متغیرهای توضیحی خطی است. بنابراین برای پیشگویی، لازم است ضرایب رگرسیونی برآورد شوند. این روش پارامتری هنگامی مناسب است که فرض خطی بودن رابطه میان متغیر پاسخ و متغیرهای توضیحی لااقل به صورت تقریبی برقرار باشد. اما هنگامی که رابطه میان متغیر پاسخ و متغیرهای توضیحی نامعلوم است، استفاده از رگرسیون خطی ممکن است منجر به نتایج گمراه کننده‌ای شود.

روش کمترین توان‌های دوم متداول‌ترین روش برآورد برای ضرایب رگرسیونی است که در سال ۱۸۰۹ توسط گاوس معرفی شد. به منظور یافتن برآوردگر کمترین توان‌های دوم معمولی (OLS)، لازم است برخی پذیره‌های زیربنایی در مدل (۲.۲) برقرار باشند که عبارتند از: الف) خطاها ناهمبسته‌اند. ب)  $E(\varepsilon) = 0$  (ج)  $E(\varepsilon^T \varepsilon) = \sigma^2 I_n$  که  $\sigma$  مجهول، مثبت و متناهی است. برآوردگر OLS بردار ضرایب  $\beta$  به صورت

$$\begin{aligned} \hat{\beta}^{OLS} &= \arg \min_{\beta \in \mathbb{R}^p} \{ \| Y - X\beta \|^2 \} \\ &= (X^T X)^{-1} X^T Y \end{aligned} \quad (۳.۲)$$

است. به سادگی می‌توان نشان داد که اگر  $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I_n)$ ، آنگاه  $\hat{\beta}^{OLS} \sim \mathcal{N}_p(\beta, (X^T X)^{-1} \sigma^2)$  بنا به قضیه گاوس-مارکوف،  $\hat{\beta}^{OLS}$  بهترین برآوردگر نارایب خطی (BLUE) است. همچنین برآوردگر نارایب واریانس  $\sigma^2$  بر اساس این برآوردگر به صورت

$$S^2 = \frac{1}{m} \| Y - X\hat{\beta}^{OLS} \|^2, \quad m = n - p$$

به دست می‌آید. با داشتن برآورد ضرایب رگرسیونی، بردار مقادیر برازش شده به صورت

$$\hat{Y} = X\hat{\beta}^{OLS} = HY \quad (۴.۲)$$

محاسبه می‌شود که در آن  $H = X(X^T X)^{-1} X^T$  ماتریس هت نامیده می‌شود و یک ماتریس متقارن و خودتوان است. بنابراین  $Var(\hat{Y}) = \sigma^2 H$  و بردار باقیمانده‌ها و واریانس آن، به ترتیب، به صورت  $e = (I_n - H)Y$  و  $Var(e) = \sigma^2 (I_n - H)$  به دست می‌آیند. همان‌طور که ملاحظه می‌شود ماتریس  $H$  نقش بسیار مهمی در نتایج روش کمترین توان‌های دوم بازی می‌کند. برای اطلاعات بیشتر در این باره به نوروزی‌راد (۱۳۹۶) مراجعه کنید.

اگر داده‌ها از یک مدل خطی با پذیره‌های مطرح‌شده پیروی کنند، روش کمترین توان‌های دوم برآوردهای خوبی را ارائه می‌دهد. هنگامی که حداقل یکی از این پذیره‌ها برقرار نباشد، برآوردهای OLS بهینه نیستند.

### ۱.۲.۲ همخطی در رگرسیون چندگانه

یک مسئله جدی که کاربرد برآوردگر OLS را با مشکل مواجه می‌کند، وجود همخطی بین متغیرهای توضیحی است، بدین معنی که متغیرهای توضیحی به صورت کامل یا تقریباً کامل وابسته خطی باشند. اگر یکی از متغیرهای توضیحی تابعی خطی از یک یا چند متغیر توضیحی دیگر باشد، گوییم رگرسیون دارای همخطی کامل است. همخطی ناقص هنگامی اتفاق می‌افتد که یکی از متغیرهای توضیحی به‌طور تقریبی یک تابع خطی از یک یا چند متغیر توضیحی دیگر باشد. به‌عبارت دیگر، همخطی کامل هنگامی رخ می‌دهد که حداقل به‌ازای یک  $j$  ( $j = 1, \dots, p$ ) داشته باشیم  $R_j^y = 1$  و همخطی ناقص زمانی پیش می‌آید که حداقل به‌ازای یک  $j$  داشته باشیم  $R_j^y \approx 1$ ، که  $R_j^y$  ضریب تعیین در رگرسیون خطی متغیر  $X_j$  بر سایر متغیرهای توضیحی است. طبق تعریف می‌دانیم که ضریب تعیین ( $R^2$ ) برابر با نسبت تغییرات بیان‌شده به تغییرات کل است. لذا در مدل رگرسیونی (۲.۲) ضریب تعیین برابر است با

$$R^2 = \frac{\hat{\beta}^T \mathbf{X}^T \mathbf{Y} - n\bar{Y}^2}{\mathbf{Y}^T \mathbf{Y} - n\bar{Y}^2}, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i. \quad (5.2)$$

وجود همخطی، تاثیرات جدی بر برآوردهای OLS دارد. برای درک بهتر موضوع، مدل رگرسیون خطی با دو متغیر توضیحی را در نظر بگیرید

$$Y = \beta_1 X_1 + \beta_2 X_2 + \varepsilon. \quad (6.2)$$

فرض کنید  $X_1$ ،  $X_2$  و  $Y$  استاندارد شده‌اند، به‌طوری‌که  $\bar{X}_j = 0$ ،  $\sum_{i=1}^n (X_{ij} - \bar{X}_j)^2 = 1$ ،  $\bar{Y} = 0$  و  $\sum_{i=1}^n (Y_i - \bar{Y})^2 = 1$ . در این صورت معادلات نرمال روش کمترین توان‌های دوم به‌صورت زیر است

$$(\mathbf{X}^T \mathbf{X}) \hat{\beta}^{OLS} = \mathbf{X}^T \mathbf{Y}$$

یا به‌طور معادل

$$\begin{bmatrix} 1 & r_{12} \\ r_{12} & 1 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} r_{1y} \\ r_{2y} \end{bmatrix}$$

که  $r_{12}$  همبستگی ساده بین  $X_1$  و  $X_2$  و  $r_{jy}$  همبستگی ساده بین  $X_j$  و  $Y$  به‌ازای  $j = 1, 2$  می‌باشند. معکوس  $\mathbf{X}^T \mathbf{X}$  عبارتست از

$$C = (\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{1-r_{12}^2} & \frac{-r_{12}}{1-r_{12}^2} \\ \frac{-r_{12}}{1-r_{12}^2} & \frac{1}{1-r_{12}^2} \end{bmatrix} \quad (7.2)$$

و برآوردهای ضرایب رگرسیونی به صورت

$$\hat{\beta}_1 = \frac{r_{1y} - r_{12}r_{2y}}{1 - r_{12}^2}, \quad \hat{\beta}_2 = \frac{r_{2y} - r_{12}r_{1y}}{1 - r_{12}^2} \quad (۸.۲)$$

است. اگر همخطی ناقص بین  $X_1$  و  $X_2$  به گونه‌ای باشد که  $r_{12}$  بزرگ شود، آنگاه با توجه به روابط (۷.۲) و (۸.۲)، برآورد ضرایب و واریانس برآوردگرها بسیار بزرگ خواهد بود. به عبارت دیگر، اگر  $|r_{12}| \rightarrow 1$ ، آنگاه، به ازای  $j = 1, 2$ ،  $Var(\hat{\beta}_j) = C_{jj}\sigma^2 \rightarrow \infty$  و  $Cov(\hat{\beta}_1, \hat{\beta}_2) = C_{12}\sigma^2 \rightarrow \infty$  که  $C_{jj}$  مؤلفه زام روی قطر اصلی ماتریس  $C$  است. بنابراین همخطی زیاد بین  $X_1$  و  $X_2$  منجر به واریانس‌ها و کوواریانس‌های بزرگ برای برآوردگرهای OLS خواهد شد. برای بیش از دو متغیر توضیحی نیز همخطی اثرات مشابهی ایجاد می‌کند. می‌توان نشان داد که در این حالت اعضای قطر اصلی ماتریس  $C = (X^T X)^{-1}$  عبارتند از

$$C_{jj} = \frac{1}{1 - R_j^2}, \quad j = 1, \dots, p \quad (۹.۲)$$

که در آن  $R_j^2$  ضریب تعیین حاصل از رگرسیون خطی  $X_j$  بر  $p - 1$  متغیر توضیحی دیگر است. اگر همخطی شدید بین  $X_j$  و هر زیر مجموعه‌ای از  $p - 1$  متغیر دیگر وجود داشته باشد، در این صورت مقدار  $R_j^2$  نزدیک به یک خواهد بود. چون  $Var(\hat{\beta}_j) = C_{jj}\sigma^2 = \sigma^2(1 - R_j^2)^{-1}$ ، همخطی شدید موجب می‌شود که واریانس برآوردگر OLS ضرایب رگرسیونی  $\beta$  بسیار بزرگ شود. در این صورت  $\hat{\beta}^{OLS}$  برآورد دقیقی از  $\beta$  نیست. به عبارت دیگر، دقت برآورد ضرایب رگرسیونی بسیار کم می‌شود. همچنین همخطی باعث می‌شود که بردار  $\hat{\beta}^{OLS}$  دارای اعضایی باشد که از لحاظ قدر مطلق بسیار بزرگ هستند. علاوه بر این،  $\hat{\beta}^{OLS}$  ناپایدار است، یعنی با حذف یا اضافه کردن یک متغیر توضیحی، یا با یک تغییر کوچک در مقادیر متغیرهای توضیحی، مقادیر  $\hat{\beta}^{OLS}$  به طور قابل توجهی تغییر می‌کنند.

در صورت وجود همخطی کامل، ماتریس  $X^T X$  وارون پذیر نیست، زیرا در این حالت پررتبه نیست. بنابراین نمی‌توان با استفاده از معادلات نرمال، جواب یکتا برای برآورد ضرایب به دست آورد. لازم به ذکر است که به دلیل وجود خطای اندازه‌گیری، همخطی کامل معمولاً در عمل رخ نمی‌دهد، مگر هنگامی که کاربر به اشتباه در مدل رگرسیونی یکی از متغیرها را تابعی خطی از متغیر دیگر قرار دهد یا تعداد متغیرهای توضیحی بیشتر از تعداد مشاهدات باشد. برای اطلاعات بیشتر درباره منابع پیدایش همخطی به مونت‌گومری<sup>۸</sup> و همکاران (۲۰۱۲) مراجعه کنید.

بنابراین لازم است بدانیم که چگونه وجود همخطی را کشف کنیم تا مواظب پیامدهای ممکن باشیم. روش‌های مختلفی برای تشخیص همخطی وجود دارند. در اینجا به برخی از این روش‌ها اشاره می‌کنیم.

۱. عامل تورم واریانس (VIF): یکی از روش‌های معروف برای تشخیص همخطی استفاده از کمیت VIF است که به صورت

$$VIF_j = \frac{1}{1 - R_j^2}, \quad j = 1, \dots, p \quad (۱۰.۲)$$

<sup>۸</sup>Montgomery

محاسبه می‌شود. واضح است که اگر  $X_j$  رابطه خطی قوی با سایر متغیرهای توضیحی داشته باشد، آنگاه  $R_j^2$  نزدیک به یک بوده و  $VIF_j$  بزرگ می‌شود. اگر  $X_j$  بر  $p - 1$  متغیر توضیحی دیگر عمود باشد، آنگاه  $R_j^2$  صفر بوده و  $VIF_j$  برابر با یک خواهد شد. بنابراین انحراف  $VIF_j$  از مقدار یک، انحراف از متعامد بودن و گرایش به همخطی را نشان می‌دهد. اگر هر یک از مقادیر  $VIF$  بیشتر از  $10$  باشد، همخطی شدید بین متغیرهای توضیحی وجود دارد و ممکن است موجب بروز مشکلاتی در برآورد شود. از آنجا که  $Var(\hat{\beta}_j) = VIF_j \sigma^2$ ، لذا  $VIF_j$  میزان افزایش در واریانس  $\hat{\beta}_j$  به واسطه ارتباط خطی  $X_j$  با بقیه متغیرهای توضیحی را نسبت به واریانس آن در صورت ناهمبسته بودن متغیر  $X_j$  با سایر متغیرهای توضیحی اندازه می‌گیرد. این مطلب دلیل نامگذاری این کمیت خاص را نشان می‌دهد.

۲. عدد شرطی<sup>۹</sup>: مقادیر ویژه ماتریس  $X^T X$  می‌توانند برای اندازه‌گیری میزان همخطی مورد استفاده قرار گیرند. یک روش متداول برای کشف وجود همخطی، محاسبه عدد شرطی به صورت

$$\kappa = \sqrt{\frac{\max_{1 \leq j \leq p} \lambda_j}{\min_{1 \leq j \leq p} \lambda_j}} \quad (11.2)$$

است که در آن  $\lambda_1, \dots, \lambda_p$  مقادیر ویژه ماتریس  $X^T X$  هستند.

بلسلی<sup>۱۰</sup> و همکاران (۱۹۸۰) پیشنهاد کردند که اگر عدد شرطی برای ماتریس  $X^T X$  بزرگتر از  $10$  باشد، آنگاه در ماتریس طرح همخطی وجود دارد. اگر  $100 < \kappa < 300$  باشد، همخطی شدید و چنانچه  $\kappa > 1000$  باشد همخطی بسیار جدی است.

تاکنون روش‌های متعددی برای مقابله با همخطی بین متغیرهای توضیحی معرفی شده است. برای اطلاعات بیشتر در این باره به مونت‌گومری و همکاران (۲۰۱۲) و چاترجی و هادی<sup>۱۱</sup> (۲۰۱۲) مراجعه کنید.

یک رهیافت مناسب برای مقابله با همخطی، استفاده از برآوردگرهای اریب مانند ریدج<sup>۱۲</sup> (هورل و کنارد<sup>۱۳</sup>، ۱۹۷۰) است. در این روش با پذیرفتن مقدار کمی اریبی، برآوردگری حاصل می‌شود که واریانس آن در مقایسه با برآوردگر OLS کمتر است. این کاهش واریانس، کاهش MSE را به دنبال خواهد داشت. در نتیجه برآوردگر ریدج در مقایسه با برآوردگر OLS پایدارتر است.

هنگامی که تعداد متغیرهای توضیحی زیاد است، همخطی امری اجتناب‌ناپذیر است. امروزه مجموعه داده‌هایی با تعداد زیادی از متغیرهای توضیحی در زمینه‌های مختلف علمی و تحقیقاتی مانند پزشکی، ژنتیک، بیوانفورماتیک، پردازش تصویر و مالی به وفور مورد استفاده قرار می‌گیرند. در این نوع داده‌ها، اطلاعات تعداد زیادی از متغیرهای توضیحی روی تعداد اندکی از واحدهای نمونه ثبت می‌شود ( $p > n$ ). ورود تعداد زیادی متغیر توضیحی به مدل اگرچه انعطاف‌پذیری مدل را افزایش می‌دهد، اما علاوه بر ایجاد همخطی شدید، باعث کاهش توان پیشگویی و مشکل‌شدن تفسیر مدل می‌شود. برآوردگر ریدج اگرچه

<sup>9</sup>Condition number

<sup>10</sup>Belsley

<sup>11</sup>Chatterjee and Hadi

<sup>12</sup>Ridge

<sup>13</sup>Hoerl and Kennard



توان پیشگویی بالایی دارد، اما به دلیل تعدد ضرایب، مدل پیچیده‌ای را نتیجه می‌دهد که تفسیر آن بسیار مشکل است. این در حالی است که در مدل‌های با بعد بالا، معمولاً تنها تعداد کمی از متغیرهای توضیحی با متغیر پاسخ مرتبط هستند. بهترین روش در این حالت، حذف متغیرهای زائد و مدل‌بندی با زیرمجموعه‌ای از مؤثرترین متغیرهای توضیحی است. تاکنون روش‌های مختلفی برای انتخاب متغیر معرفی شده‌اند که در ادامه به معرفی مهمترین آنها می‌پردازیم.

## ۳.۲ انتخاب متغیر با معیارهای کلاسیک

تاکنون الگوریتم‌های متعددی برای انتخاب متغیر در مدل خطی (۳.۲) معرفی شده‌اند. این الگوریتم‌ها بر اساس معیارهایی مانند معیار اطلاع آکائیک (AIC)، معیار اطلاع بیزی (BIC)، اعتبارسنجی متقابل تعمیم‌یافته (GCV)، ضریب تعیین تعدیل‌شده ( $R_{adj}^2$ )، مجموع توان‌های دوم باقیمانده‌ها (RSS) و  $C_p$  - مالو<sup>۱۴</sup> زیرمجموعه‌های متفاوتی از متغیرهای توضیحی را به‌عنوان بهترین زیرمجموعه انتخاب می‌کنند. از جمله این روش‌ها می‌توان انتخاب بهترین زیرمجموعه<sup>۱۵</sup>، انتخاب پیش‌رو<sup>۱۶</sup>، روش حذفی پس‌رو<sup>۱۷</sup> و روش گام به گام<sup>۱۸</sup> را نام برد.

در روش انتخاب بهترین زیرمجموعه، تمام زیرمجموعه‌های ممکن را در نظر گرفته و سپس بر اساس معیارهای فوق، بهترین زیرمجموعه از مدل کامل انتخاب می‌شود. در این حالت، اگر  $p$  تعداد کل متغیرهای توضیحی باشد، تعداد  $2^p$  زیرمجموعه باید مورد بررسی قرار گیرند. به‌عنوان مثال، اگر  $p = 10$ ، آنگاه باید ۱۰۲۴ مدل ممکن در نظر گرفته شوند و برای  $p = 20$  باید بیش از یک میلیون مدل ممکن بررسی شوند. بنابراین با افزایش  $p$ ، تعداد زیرمجموعه‌های ممکن به سرعت افزایش می‌یابد. لذا روش انتخاب بهترین زیرمجموعه از نظر محاسباتی، حتی در صورت استفاده از کامپیوترهای با سرعت بالا، برای  $p > 40$  قابل استفاده نیست. در این راستا، دراپر و اسمیت<sup>۱۹</sup> (۱۹۸۱) روش‌های انتخاب پیش‌رو، حذفی پس‌رو و گام به گام را معرفی کردند. روش انتخاب پیش‌رو با مدلی که تنها شامل عرض از مبدا است، شروع می‌شود و سپس در هر گام متغیری که بیشترین تأثیر را در بهبود مدل داشته باشد، انتخاب می‌شود. این رویه تا زمانی ادامه می‌یابد که تأثیر هیچ یک از متغیرهای باقیمانده در بهبود مدل معنی‌دار نباشد. برای اندازه‌گیری میزان تأثیر هر یک از متغیرها در بهبود مدل، از آماره‌هایی مانند  $R^2$ ، مجموع توان‌های دوم باقیمانده‌ها (RSS)، آماره F و آماره t استفاده می‌شود. به‌طور دقیق‌تر، در روش انتخاب پیش‌رو، برای انتخاب بهترین مدل بر اساس آماره F، به‌صورت زیر عمل می‌کنیم. اولین متغیر توضیحی که برای ورود به مدل نامزد می‌شود، متغیری است که دارای بیشترین همبستگی ساده با متغیر پاسخ است. این متغیر توضیحی همان متغیری است که بیشترین مقدار آماره F را برای آزمون معنی‌داری رگرسیون نتیجه می‌دهد. اگر آماره F متناظر با این متغیر از مقدار F از پیش تعیین‌شده  $F^*$

<sup>14</sup>Mallows's  $C_p$

<sup>15</sup>Best subset selection

<sup>16</sup>Forward selection

<sup>17</sup>Backward elimination

<sup>18</sup>Stepwise

<sup>19</sup>Draper and Smith

تجاوز کند، این متغیر توضیحی وارد مدل می‌شود. فرض کنیم اولین متغیر توضیحی وارد شده به مدل  $X_1$  باشد. دومین متغیر توضیحی منتخب برای ورود به مدل آن است که پس از تعدیل اثر  $X_1$  روی  $Y$ ، اکنون بیشتر همبستگی را با  $Y$  دارد. به عبارت دیگر، دومین متغیر منتخب متغیری است که دارای بیشترین همبستگی جزئی با  $Y$  باشد. فرض کنیم دومین متغیر انتخاب شده  $X_2$  باشد. در این صورت بزرگترین آماره  $F$  جزئی عبارتست از

$$F = \frac{SSR(X_2|X_1)}{MSE(X_1, X_2)}$$

که  $MSE(X_1, X_2)$  میانگین مربعات خطا و  $SSR(X_2|X_1) = SSR(X_1, X_2) - SSR(X_1)$  نشان‌دهنده میزان افزایش در مجموع مربعات رگرسیونی است که با افزودن متغیر توضیحی  $X_2$  به مدل حاصل می‌شود. این افزایش در مجموع مربعات رگرسیونی سهم اثرگذاری متغیر توضیحی  $X_2$  را به شرط حضور  $X_1$  در مدل اندازه‌گیری می‌کند. اگر آماره  $F$  جزئی از  $F^*$  بزرگتر باشد،  $X_2$  به مدل اضافه می‌شود. به طور کلی، در هر گام، متغیر توضیحی که دارای بیشترین همبستگی جزئی با  $Y$  است، یا به طور معادل دارای بیشترین آماره  $F$  جزئی به شرط حضور دیگر متغیرهای توضیحی در مدل است، برای ورود به مدل نامزد می‌شود. اگر آماره  $F$  جزئی آن متغیر از مقدار  $F^*$  تجاوز کند، متغیر مورد نظر وارد مدل می‌شود. این رویه هنگامی پایان می‌یابد که بزرگترین آماره  $F$  از مقدار  $F^*$  تجاوز نکند یا همه متغیرهای توضیحی وارد مدل شده باشند. در عمل معمولاً مقدار  $F^*$  برابر  $F_{\alpha, 1, n-p}$  در نظر گرفته می‌شود که  $\alpha$  بیانگر خطای نوع اول است.

در روش حذفی پس‌رو، فرآیند با مدل کامل شروع می‌شود و هر بار یک متغیر از مدل حذف می‌شود. متغیر حذف شده همان متغیری است که کمترین تأثیر را در بهبود مدل دارد. به عبارت دیگر، در هر گام، متغیر توضیحی با کوچکترین آماره  $F$  جزئی برای حذف از مدل نامزد می‌شود. اگر آماره  $F$  جزئی متغیر منتخب کمتر از مقدار  $F^*$  باشد، متغیر مورد نظر از مدل حذف می‌شود. فرآیند حذف متغیرها زمانی متوقف می‌شود که کوچکترین مقدار آماره  $F$  جزئی کوچکتر از مقدار از پیش تعیین شده  $F^*$  نباشد یا همه متغیرهای توضیحی از مدل حذف شده باشند.

در روش‌های انتخاب پیش‌رو و حذفی پس‌رو، هنگامی که یک متغیر حذف یا اضافه می‌شود، در مراحل بعدی معنی‌داری این متغیرها مجدداً بررسی نمی‌شود. روش گام به گام تعدیل روش انتخاب پیش‌رو می‌باشد که در هر گام همه متغیرهای توضیحی وارد شده به مدل در گام‌های قبلی، مجدداً ارزیابی می‌شوند. بنابراین یک متغیر توضیحی اضافه شده در گام قبلی ممکن است در گام بعدی از مدل حذف شود.

برای اطلاعات بیشتر درباره سایر روش‌های انتخاب متغیر می‌توان به دزیاک<sup>۲۰</sup> و همکاران (۲۰۰۵)، هاینز<sup>۲۱</sup> و همکاران (۲۰۱۸)، دسبولتس<sup>۲۲</sup> (۲۰۱۸) و دینگ<sup>۲۳</sup> (۲۰۱۸) مراجعه کرد. اگرچه روش‌های انتخاب زیرمجموعه با معیارهای کلاسیک از ویژگی‌های نمونه‌ای مطلوبی برخوردارند، اما این روش‌ها پایدار نیستند، بدین معنی که یک تغییر بسیار کوچک در داده‌ها ممکن است مدل‌های

<sup>20</sup>Dziak

<sup>21</sup>Heinze

<sup>22</sup>Desboulets

<sup>23</sup>Ding

کاملاً متفاوت را نتیجه دهد. لذا ممکن است این روش‌ها پیشگویی نامناسبی را ایجاد کنند (بريمن، ۱۹۹۶). مشکل دیگر این روش‌ها این است که تمامی مدل‌های ممکن را در نظر نمی‌گیرند و این موضوع ممکن است منجر به از دست دادن یک مدل خوب شود. در روش‌های انتخاب پیش‌رو و حذفی پس‌رو، حداکثر تعداد  $1 + \dots + (p-1) + p$  مدل مورد بررسی قرار می‌گیرند. به‌عنوان مثال، برای  $p = 15$ ، نسبت ماکزیمم تعداد مدل‌های مورد بررسی به تعداد کل مدل‌های ممکن برابر با

$$\frac{15 + 14 + \dots + 1}{2^{15} - 1} = \frac{120}{32767} \approx 0.4\% \quad (12.2)$$

است، یعنی  $99.6\%$  از مدل‌های ممکن برای انتخاب بهترین مدل در نظر گرفته نمی‌شوند. بنابراین مدل نهایی انتخاب شده ممکن است بهترین مدل نباشد. از طرفی در صورت وجود همخطی بین متغیرهای توضیحی، استفاده از روش‌های مذکور ممکن است منجر به انتخاب مدل مناسبی نشود (دزیاک و همکاران، ۲۰۰۵). علاوه بر این، زمانی که  $p$  بزرگ باشد، هزینه محاسباتی این روش‌ها بسیار بالاست. برای رفع این چالش‌ها، محققین رگرسیون جریمه‌شده را معرفی کردند که، با انتخاب یک تابع جریمه مناسب، دو عمل انتخاب متغیر و برآورد ضرایب را به‌طور همزمان انجام می‌دهد.

## ۴.۲ انتخاب متغیر با کمترین توان‌های دوم جریمه‌شده

در این بخش، به توصیف روش کمترین توان‌های دوم جریمه‌شده می‌پردازیم. بدیهی است وقتی تعداد متغیرهای توضیحی  $p$  زیاد است، اندازه مجموع قدر مطلق ضرایب  $\beta_j$ ،  $j = 1, \dots, p$ ، افزایش می‌یابد. از آنجایی که در روش کمترین توان‌های دوم به دنبال مینیمم کردن  $\|\varepsilon\|^2 = \varepsilon^T \varepsilon$  نسبت به  $\beta$  هستیم، پس چنانچه بخواهیم علاوه بر برآورد ضرایب رگرسیونی، تعداد متغیرهای توضیحی مدل را نیز کاهش دهیم، کافی است بزرگی قدرمطلق ضرایب یا تابعی از آن را نیز به‌طور همزمان مینیمم کنیم. بنابراین، در رگرسیون کمترین توان‌های دوم جریمه‌شده، محک  $\|\varepsilon\|^2$  را تحت این قید که تابعی از قدرمطلق ضرایب رگرسیونی مانند  $\sum_{j=1}^p p(|\beta_j|)$  نباید بزرگ باشد، مینیمم می‌کنند. این یک مسئله بهینه‌سازی مقید<sup>۲۴</sup> است که به‌صورت زیر بیان می‌شود

$$\min_{\beta} \|\varepsilon\|^2 \quad \text{s.t} \quad \sum_{j=1}^p p(|\beta_j|) \leq s. \quad (13.2)$$

یکی از راه‌های حل مسئله بهینه‌سازی مقید، استفاده از روش لاگرانژ<sup>۲۵</sup> است که با استفاده از این روش، مسئله بهینه‌سازی (۱۳.۲) را می‌توان به صورت زیر بیان کرد

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} \left\{ \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p p(|\beta_j|) \right\} \\ &= \arg \min_{\beta} \left\{ \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\}. \end{aligned} \quad (14.2)$$

<sup>24</sup>Constraint optimization

<sup>25</sup>Lagrange method

با دقت در ساختار مسئله (۱۴.۲) دیده می‌شود که به تابع مجموع توان‌های دوم خطا در رگرسیون معمولی، عبارت جریمه  $\sum_{j=1}^p p_\lambda(|\beta_j|)$  اضافه شده است. در اینجا، تابع جریمه و  $\lambda \geq 0$  پارامتر جریمه است که میزان جریمه را برحسب بزرگی ضرایب رگرسیونی کنترل می‌کند. در این حالت بین  $\lambda$  در (۱۴.۲) و  $s$  در (۱۳.۲) رابطه‌ای معکوس وجود دارد. به‌طور دقیق‌تر، اگر بخواهیم تعداد متغیرهای بیشتری از مدل حذف شوند، کافی است  $s$  را کوچک یا  $\lambda$  را بزرگ انتخاب کنیم و برعکس. پارامتر تنظیم‌کننده  $\lambda$  موازنه‌ای بین خطای پیشگویی و پیچیدگی مدل برقرار می‌کند. جواب مسئله بهینه‌سازی (۱۴.۲) به شدت به مقدار  $\lambda$  وابسته است و لذا تعیین مقدار مناسب برای  $\lambda$  یک مسئله چالش‌برانگیز است. روش‌های مختلف انتخاب مقدار بهینه  $\lambda$  در زیربخش ۳.۴.۲ آمده‌اند.

بنابراین با انتخاب یک مقدار بهینه برای  $\lambda$ ، مقادیر برخی از پارامترها ممکن است صفر شوند که در این صورت متغیر توضیحی متناظر با آن پارامتر، به عنوان متغیر بی‌اثر شناخته شده و بدین ترتیب مجموعه متغیرهای مهم و بی‌اهمیت از یکدیگر تشخیص داده می‌شوند و در اصطلاح "انتخاب متغیر" صورت می‌گیرد.

## ۱.۴.۲ انواع توابع جریمه

تا کنون توابع جریمه زیادی توسط افراد مختلف معرفی شده‌اند که برخی از آن‌ها دو عمل انتخاب متغیر و برآورد را به‌طور همزمان انجام می‌دهند. در این بخش به معرفی برخی از توابع جریمه معروف می‌پردازیم.

۱. جریمه  $L_q$ : این تابع جریمه، رگرسیون بریج را نتیجه می‌دهد که توسط فرانک و فریدمن<sup>۲۶</sup> در سال ۱۹۹۳ معرفی شد. برآوردگر بریج به صورت

$$\hat{\beta}^{Bridge} = \arg \min_{\beta} \left\{ \frac{1}{2} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\}, \quad 0 < q \leq 2 \quad (15.2)$$

تعریف می‌شود. این تابع جریمه، برخی از توابع جریمه مهم را به عنوان یک حالت خاص شامل می‌شود.

### حالت خاص $q = 2$

مدل رگرسیونی حاصل را در این حالت رگرسیون ریج گویند و برآورد پارامترها به صورت

$$\begin{aligned} \hat{\beta}^{Ridge} &= \arg \min_{\beta} \left\{ \frac{1}{2} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \lambda \sum_{j=1}^p |\beta_j|^2 \right\} \\ &= (\mathbf{X}^T \mathbf{X} + K \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{Y} \end{aligned}$$

به دست می‌آید که  $K = 2\lambda$  را پارامتر تنظیم‌کننده یا پارامتر ریج گویند. یک مزیت رگرسیون ریج پیاده‌سازی آسان و شکل صریح جواب آن است.

<sup>26</sup>Frank and Friedman

در برآورد کمترین توان‌های دوم معمولی، اگر ماتریس  $\mathbf{X}$  دارای رتبه کامل ستونی نباشد، ماتریس  $\mathbf{X}^T \mathbf{X}$  وارون‌پذیر نخواهد بود و لذا جواب یکتایی برای برآوردگر OLS وجود نخواهد داشت. اما برای هر ماتریس دلخواه  $\mathbf{X}$  و  $K > 0$ ، ماتریس  $\mathbf{X}^T \mathbf{X} + K \mathbf{I}_p$  همواره وارون‌پذیر است و در نتیجه پاسخ یکتا برای برآوردگر ریج وجود دارد. برآوردگر ریج اریب است، اما در مقایسه با برآوردگر OLS دارای واریانس کوچکتری است. همچنین همواره  $\lambda$  ای وجود دارد به طوری که  $MSE(\hat{\beta}^{Ridge}) < MSE(\hat{\beta}^{OLS})$  (برهنی و هوانگ<sup>۲۷</sup>، ۲۰۱۹). بنابراین رگرسیون ریج پیشگویی بهتری از رگرسیون کمترین توان‌های دوم معمولی ارائه می‌کند. با این وجود، برآوردگر ریج اگرچه ضرایب رگرسیونی را منقبض می‌کند، اما هیچ یک از برآوردهای ضرایب دقیقاً برابر صفر نیستند، یعنی این روش توانایی انتخاب متغیر ندارد. این موضوع خللی در پیشگویی ایجاد نمی‌کند، اما با توجه به حضور همه متغیرها در مدل، تفسیر مدل برازش‌شده به سادگی امکان‌پذیر نیست.

در حالتی که  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$  ماتریس طرح متعامد خواهد بود. در این حالت، برآوردگر ریج به صورت

$$\hat{\beta}^{Ridge} = \frac{1}{1+K} \hat{\beta}^{OLS} \quad (16.2)$$

در می‌آید. این رابطه ماهیت انقباضی رگرسیون جریمه‌شده را نشان می‌دهد، چون با افزایش  $K$  ضرایب به سمت صفر منقبض می‌شوند. هنگامی که  $K = 0$ ، برآوردهای OLS و ریج معادل‌اند.

#### حالت خاص $q = 1$

مدل رگرسیونی حاصل در این حالت را رگرسیون لاسو (Lasso) گویند که توسط تیبشیرانی<sup>۲۸</sup> در سال ۱۹۹۶ معرفی شد. از دیدگاه آماری، لاسو به معنی عملگر انتخاب‌کننده، منقبض‌کننده و مینیمم‌کننده از جنس قدرمطلق است. در رگرسیون لاسو برآورد ضرایب رگرسیونی به صورت

$$\hat{\beta}^{Lasso} = \arg \min_{\beta} \left\{ \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p |\beta_j| \right\} \quad (17.2)$$

به دست می‌آید که در آن  $\lambda$  پارامتر جریمه است و سطح تنگی (تعداد ضرایب صفر) را کنترل می‌کند. وقتی  $\lambda \rightarrow 0$ ، آنگاه  $\hat{\beta}^{Lasso} \rightarrow \hat{\beta}^{OLS}$  و وقتی  $\lambda \rightarrow \infty$ ، آنگاه  $\hat{\beta}^{Lasso} \rightarrow 0$ . به عبارت دیگر، هرچه جریمه بزرگتری به کار رود، تعداد بیشتری از ضرایب به سمت صفر منقبض می‌شوند. برآورد به دست آمده با رگرسیون لاسو تنگ است، یعنی برآوردهای برخی از ضرایب برابر صفر هستند. به عبارت دیگر، این روش، انتخاب متغیر و برآورد پارامتر را به طور همزمان انجام می‌دهد.

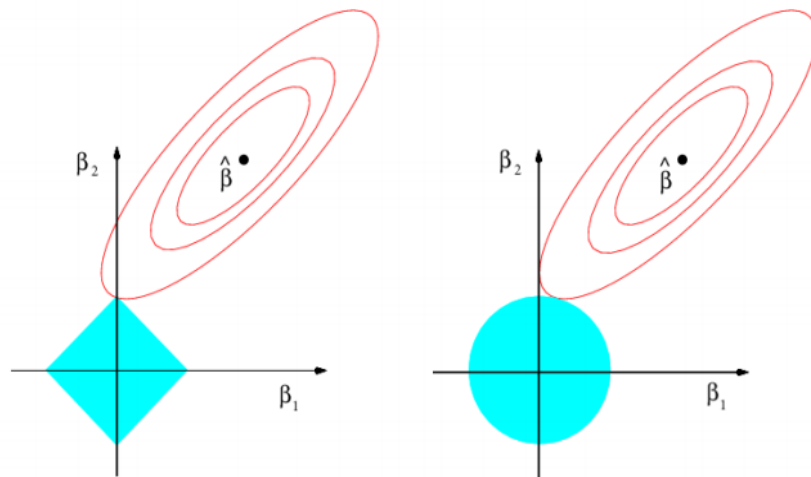
تابع جریمه قدرمطلق باعث می‌شود که جواب صریحی برای برآوردگر لاسو وجود نداشته باشد. تیبشیرانی این برآوردگر را از طریق برنامه‌ریزی درجه دوم<sup>۲۹</sup> به دست آورده بود. سپس افرون<sup>۳۰</sup> و همکاران

<sup>27</sup>Breheeny and Huang

<sup>28</sup>Tibshirani

<sup>29</sup>Quadratic programming

<sup>30</sup>Efron



شکل ۱.۲: منحنی‌های تراز و نواحی محدودیت رگرسیون لاسو (قاب سمت چپ) و رگرسیون ریج (قاب سمت راست). ناحیه محدودیت لاسو به صورت  $|\beta_1| + |\beta_2| \leq s$  و ناحیه محدودیت ریج به صورت  $\beta_1^2 + \beta_2^2 \leq s$  است. بیضی‌ها، منحنی‌های تراز تابع خطای برآورد کمترین توان‌های دوم هستند.

(۲۰۰۴) رگرسیون کمترین زاویه<sup>۳۱</sup> را معرفی کردند که برآوردهای لاسو را نتیجه می‌دهد و هزینه محاسباتی آن با برآوردهای کمترین توان‌های دوم برابر است. برای روشن‌تر شدن مفهوم انتخاب متغیر توسط تابع جریمه لاسو، حالت خاصی از مدل رگرسیونی را به صورت زیر در نظر بگیرید

$$Y = \beta_1 X_1 + \beta_2 X_2 + \varepsilon.$$

بنابراین مجموع توان‌های دوم خطا تابعی از دو پارامتر  $\beta_1$  و  $\beta_2$  است

$$Q(\beta) = \| \mathbf{Y} - \mathbf{X}\beta \|^2 = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \beta_1 X_{i1} - \beta_2 X_{i2})^2. \quad (18.2)$$

نمودار تابع  $Q(\cdot)$ ، یک رویه درجه دوم به صورت

$$Q(\beta) = a_0 + a_1 \beta_1 + a_2 \beta_2 + a_3 \beta_1^2 + a_4 \beta_2^2 + a_5 \beta_1 \beta_2 \quad (19.2)$$

است که منحنی‌های تراز آن در شکل ۱.۲ نشان داده شده است. اکنون شرط‌های مینیمم‌سازی ریج و لاسو را به صورت

$$\beta_1^2 + \beta_2^2 = s \quad \text{و} \quad |\beta_1| + |\beta_2| = s \quad (20.2)$$

در نظر بگیرید که شکل هندسی آن‌ها به صورت دایره برای رگرسیون ریج (قاب سمت راست) و لوزی برای رگرسیون لاسو (قاب سمت چپ) نشان داده شده است.

<sup>31</sup>Least angle regression

با توجه به قاب سمت چپ شکل ۱.۲، در رگرسیون لاسو حالتی می‌تواند رخ دهد که یکی از ضرایب صفر شده و دیگری همواره برابر  $s$  باشد، در صورتی که در قاب سمت راست این امکان وجود ندارد. به عبارت دقیق‌تر، در قاب سمت راست، اگر یکی از ضرایب صفر شود، آنگاه دیگری می‌تواند  $\sqrt{s}$ ، و نه  $s$ ، باشد. بنابراین برای لاسو امکان صفر برآورد شدن ضرایب وجود دارد، اما برای ريج این اتفاق هرگز رخ نمی‌دهد. در صورت صفر بودن برآوردهای هر یک از ضرایب در رگرسیون لاسو، متغیر متناظر با آن ضریب به عنوان متغیر بی‌اهمیت شناخته می‌شود.

لازم به ذکر است که در هر دو روش لاسو و ريج، با انتخاب مقدار کوچک  $s$ ، هر دو پارامتر  $\beta_1$  و  $\beta_2$  (در حالت دو بعدی) تا حد ممکن می‌توانند کوچک شوند. بنابراین هر دو روش ريج و لاسو انقباضی هستند.

اگر ماتریس طرح متعامد باشد، آنگاه می‌توان برآوردگر لاسو را به عنوان تابعی از برآوردگر کمترین توان‌های دوم معمولی  $\beta_j$  به صورت

$$\hat{\beta}_j^{Lasso} = \text{sgn}(z_j)(|z_j| - \lambda)_+, \quad j = 1, \dots, p \quad (21.2)$$

به دست آورد که در آن  $I(\cdot)$  تابع نشانگر،  $x_+ = \max\{0, x\}$  و  $z_j = \mathbf{X}_j^T \mathbf{Y}$  برآورد OLS ضریب  $\beta_j$  است که  $\mathbf{X}_j$  ستون  $j$  ام ماتریس طرح  $\mathbf{X}$  است (فن و لی، ۲۰۰۱). برآوردگر (۲۱.۲) را به صورت

$$\hat{\beta}_j^{Lasso} = S(z_j | \lambda) \quad (22.2)$$

نشان می‌دهیم که  $S(\cdot | \lambda)$  عملگر آستانه ای نرم نامیده می‌شود. رابطه (۲۱.۲) در بررسی ویژگی‌های مطلوب رگرسیون لاسو مفید است که در پایان این بخش به این موضوع خواهیم پرداخت.

رگرسیون لاسو با اینکه در بسیاری از مسایل انتخاب متغیر عملکرد خوبی از خود نشان داده، اما دارای برخی محدودیت‌ها است. به عنوان مثال، لاسو برای یک مدل رگرسیونی خطی با  $p$  متغیر پیشگو و  $n$  مشاهده، حداکثر  $n$  متغیر را انتخاب می‌کند. بنابراین اگر تعداد متغیرهای توضیحی معنی‌دار در مدل بیشتر از  $n$  باشد، برخی از آن‌ها توسط لاسو انتخاب نمی‌شوند. لاسو در بین یک مجموعه از متغیرهای به شدت همبسته، فقط یک متغیر را انتخاب می‌کند (زو و هستی، ۲۰۰۵). برای حالت معمولی  $n > p$ ، اگر همبستگی بالایی بین متغیرهای پیشگو وجود داشته باشد، کارایی پیشگویی برآوردگر لاسو به خوبی برآوردگر ريج نیست (تیبشیرانی، ۱۹۹۶). رگرسیون لاسو برآوردگری را نتیجه می‌دهد که اریبی آن زیاد است. در حالتی که ماتریس طرح متعامد است، می‌توان نشان داد که برای  $j = 1, \dots, p$  داریم

$$\begin{cases} E|\hat{\beta}_j^{Lasso} - \beta_j| = 0, & \beta_j = 0, \\ E|\hat{\beta}_j^{Lasso} - \beta_j| \approx \beta_j, & \beta_j \in [0, \lambda], \\ E|\hat{\beta}_j^{Lasso} - \beta_j| \approx \lambda, & |\beta_j| > \lambda. \end{cases} \quad (23.2)$$

بنابراین مقدار اریبی برآوردگر لاسو برای ضرایب بزرگ، نزدیک به  $\lambda$  است (برهنی و هوانگ، ۲۰۱۹). برای از بین بردن این محدودیت‌ها، توابع جریمه متعددی معرفی شدند که در ادامه به برخی از آن‌ها اشاره می‌کنیم. قبل از معرفی این توابع جریمه، ابتدا ویژگی پیشگویی را تعریف می‌کنیم.

**ویژگی پیشگویی:** فرض کنید  $\mathcal{A} = \{j : \beta_j \neq 0\}$  مجموعه اندیس ضرایب غیرصفر باشد، به طوری که  $p > |\mathcal{A}| = p_0$ . در این صورت، مدل درست تنها به زیرمجموعه‌ای از متغیرها (با اندازه  $p_0$ ) وابسته است. منظور از مدل درست مدلی است که تنها شامل متغیرهای مؤثر بر پاسخ باشد. بر اساس پیشنهاد فن و لی (۲۰۰۱)، یک روش برآورد ضرایب دارای ویژگی پیشگویی<sup>۳۲</sup> است، اگر به طور مجانبی دارای خواص زیر باشد:

۱. بتواند ضرایب غیرصفر را به درستی شناسایی کند، یعنی اگر  $\hat{\mathcal{A}} = \{j : \hat{\beta}_j \neq 0\}$ ، آنگاه

$$\lim_{n \rightarrow \infty} P(\hat{\mathcal{A}} = \mathcal{A}) = 1$$

۲. دارای نرخ برآورد بهینه باشد. به عبارت دیگر

$$\sqrt{n}(\hat{\beta}_{\mathcal{A}} - \beta_{\mathcal{A}}) \xrightarrow{D} \mathcal{N}(0, \Sigma^*),$$

که  $\Sigma^*$  ماتریس کوواریانس زیر مدل درست است.

بنابراین مطلوب است که یک برآوردگر دارای ویژگی پیشگویی باشد. زو (۲۰۰۶)<sup>۳۳</sup> نشان داد که برآوردگر لاسو دارای ویژگی پیشگویی نیست. با توجه به رابطه (۲۳.۲) ملاحظه می‌شود که مقدار اریبی برآورد لاسو توسط  $\lambda$  تعیین می‌شود. بنابراین یک روش کاهش اریبی لاسو، استفاده از جریمه وزنی به صورت  $\lambda_j = w_j \lambda$  است. اگر بتوانیم وزن‌ها را طوری انتخاب کنیم که به ضرایب بزرگتر جریمه کوچکتری تخصیص داده شود، آنگاه با حفظ ویژگی تنگی، اریبی برآوردگر لاسو کاهش می‌یابد. بنابراین زو (۲۰۰۶) لاسوی تطبیقی را معرفی کرد و نشان داد که این روش از ویژگی پیشگویی برخوردار است.

**۲. جریمه لاسوی تطبیقی:** صورت کلی تابع جریمه لاسوی تطبیقی مانند لاسو است، اما برخلاف لاسو برای ضرایب مختلف، جریمه‌های مختلفی در نظر می‌گیرد. برآوردگر لاسوی تطبیقی به صورت زیر بیان می‌شود

$$\hat{\beta}^{aLasso} = \arg \min_{\beta} \left\{ \frac{1}{2} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \lambda \sum_{j=1}^p w_j |\beta_j| \right\} \quad (24.2)$$

که  $w_j$  وزن  $\beta_j$  است تا ضرایب کوچک را بیشتر از ضرایب بزرگ جریمه کند. معمولاً در رگرسیون لاسوی تطبیقی، وزن  $w_j = \frac{1}{|\hat{\beta}_j^\circ|^\gamma}$  را انتخاب می‌کنند که  $\gamma > 0$  و  $\hat{\beta}_j^\circ$  یک برآوردگر  $\sqrt{n}$  سازگار از  $\beta_j$  مانند برآوردگر OLS، ریج یا لاسو است.

**۳. جریمه SCAD:** فن و لی (۲۰۰۱) یک تابع جریمه محدب به نام SCAD را معرفی کردند و نشان دادند که این تابع جریمه، از همه ویژگی‌های مطلوب یک تابع جریمه برخوردار است (زیربخش ۲.۴.۲ را ببینید). آن‌ها همچنین ویژگی پیشگویی را نیز برای این روش ثابت کردند. برآوردگر SCAD به صورت

<sup>32</sup>Oracle property

<sup>33</sup>Zou



$$\hat{\beta}^{SCAD} = \arg \min_{\beta} \left\{ \frac{1}{2} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\} \quad (25.2)$$

تعریف می‌شود که  $p_{\lambda}(\cdot)$  تابع جریمه SCAD است و به صورت

$$p_{\lambda}(t) = \begin{cases} \lambda t, & 0 \leq t \leq \lambda \\ \frac{\gamma a \lambda t - t^2 - \lambda^2}{\gamma(a-1)}, & \lambda < t < a\lambda \\ \frac{\lambda^2(a+1)}{\gamma}, & t \geq a\lambda \end{cases} \quad (26.2)$$

تعریف می‌شود که  $t \geq 0, p_{\lambda}(0) = 0, \lambda$  پارامتر جریمه و  $a > 2$  مقداری ثابت است. فن و لی (۲۰۰۱) مقدار  $a = 3/7$  را پیشنهاد کردند و نشان دادند که این مقدار در شرایط مختلف عملکرد خوبی را نتیجه می‌دهد.

شکل ۲.۲ نمودار تابع SCAD را به‌ازای  $\lambda = 1$  نشان می‌دهد. با توجه به شکل ۲.۲ و رابطه (۲۶.۲) ملاحظه می‌شود که در فاصله  $|\beta| \leq \lambda$  تابع جریمه SCAD بر لاسو منطبق است، و پس از آن، تا زمانی که  $|\beta| \leq a\lambda$ ، جریمه SCAD یک تابع درجه دوم است. سپس به‌ازای  $|\beta| > a\lambda$  به یک تابع ثابت تبدیل می‌شود. همچنین، با توجه به (۲۶.۲)، وقتی  $a \rightarrow \infty$ ، تابع جریمه SCAD با تابع جریمه لاسو معادل است، یعنی تابع جریمه SCAD در حالت خاص  $a = \infty$  برآوردگر لاسو را نتیجه می‌دهد. با فرض متعامد بودن ماتریس طرح، برآوردهای SCAD ضرایب را می‌توان به صورت

$$\hat{\beta}_j^{SCAD}(a, \lambda) = \begin{cases} S(z_j|\lambda), & |z_j| \leq 2\lambda \\ \frac{a-1}{a-\gamma} S(z_j|\frac{a\lambda}{a-1}), & 2\lambda < |z_j| \leq a\lambda \\ z_j, & |z_j| > a\lambda \end{cases} \quad (27.2)$$

نشان داد که  $z_j$  برآورد OLS ضریب  $\beta_j$  و  $S(\cdot|\lambda)$  عملگر آستانه‌ای نرم است (فن و لی، ۲۰۰۱).

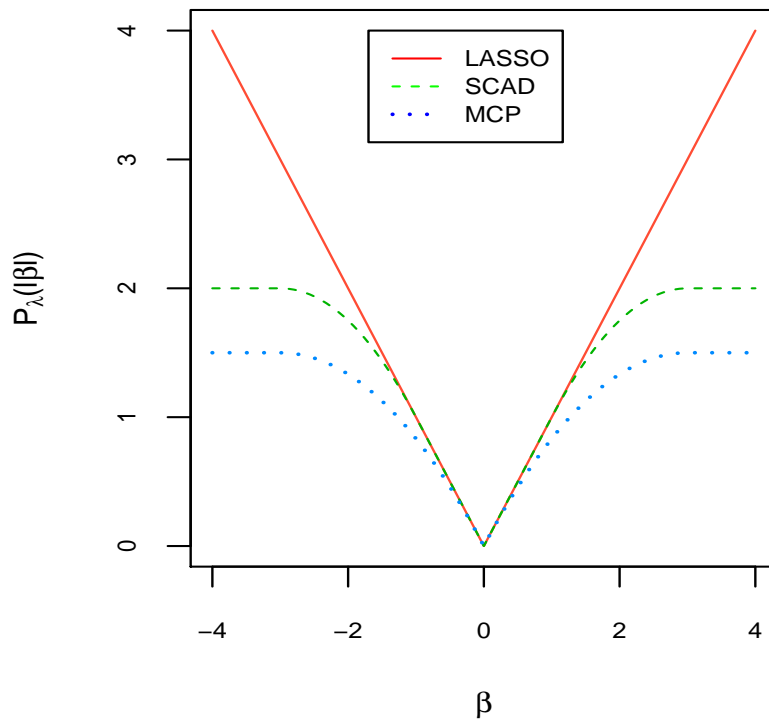
۴. جریمه MCP: ژانگ (۲۰۱۰) تابع جریمه دیگری به نام MCP را به صورت

$$p_{\lambda}(t) = \lambda \int_0^t \left( 1 - \frac{x}{\gamma\lambda} \right)_+ dx, \quad t \geq 0 \quad (28.2)$$

معرفی کرد که در آن  $\lambda$  پارامتر جریمه و  $\gamma > 0$  پارامتر تنظیم‌کننده است که میزان مقعر بودن تابع  $p_{\lambda}(\cdot)$  را کنترل می‌کند. در این حالت، برآوردگر MCP به صورت

$$\hat{\beta}^{MCP} = \arg \min_{\beta} \left\{ \frac{1}{2} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\} \quad (29.2)$$

تعریف می‌شود. تابع جریمه MCP به‌ازای  $\gamma = \infty$  تابع جریمه لاسو را نتیجه می‌دهد.



شکل ۲.۲: نمودار توابع جریمه لاسو، SCAD و MCP.

تحت فرض متعامد بودن ماتریس طرح، مشابه برآوردهای لاسو و SCAD، برآوردهای MCP ضرایب را می‌توان به صورت

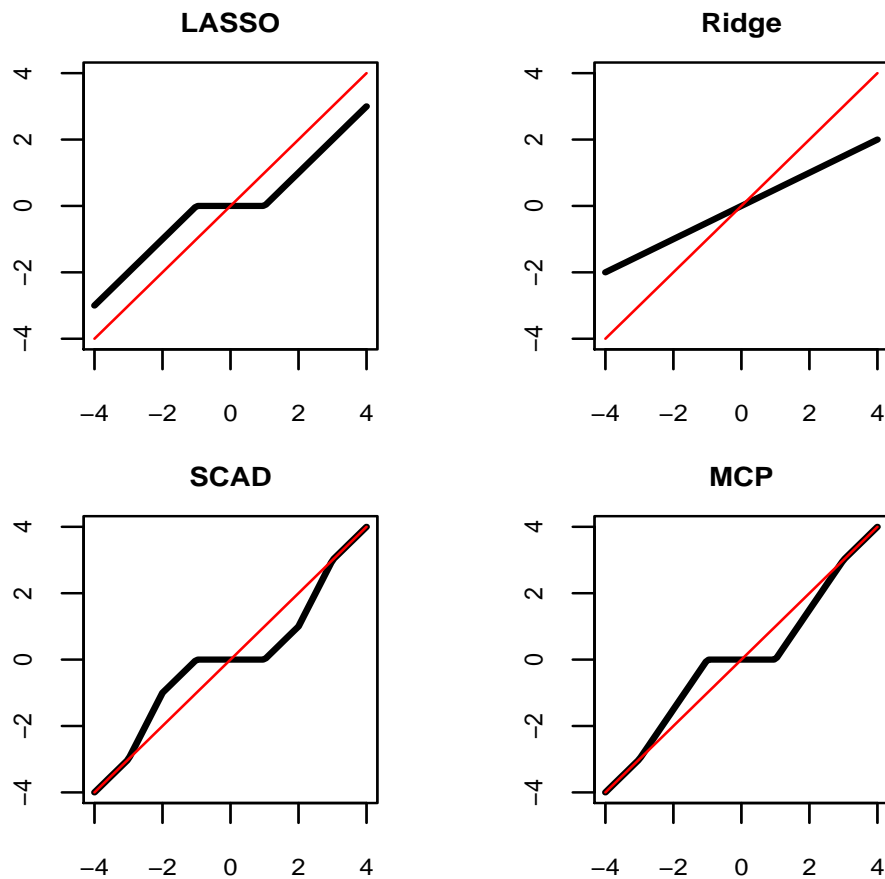
$$\hat{\beta}_j^{MCP}(\gamma, \lambda) = \begin{cases} \frac{\gamma}{\gamma-1} S(z_j|\lambda), & |z_j| \leq \gamma\lambda \\ z_j, & |z_j| > \gamma\lambda \end{cases}$$

نشان داد (برهنی و هوانگ، ۲۰۱۹). شکل ۲.۲ نمودار تابع MCP را به‌زای  $\lambda = 1$  و  $\gamma = 3$  نشان می‌دهد. تابع جریمه MCP در فاصله  $|\beta| \leq \gamma\lambda$  یک تابع درجه دوم و پس از آن یک تابع ثابت است.

### ۲.۴.۲ ویژگی‌های یک تابع جریمه مناسب

برای استفاده از رگرسیون جریمه‌شده، سوالی که پیش می‌آید این است که از چه نوع تابع جریمه باید استفاده کرد. فن و لی (۲۰۰۱) نشان دادند که یک تابع جریمه خوب باید برآوردگری با سه ویژگی مطلوب زیر را نتیجه دهد:

۱. تنکی: برآوردگر نتیجه‌شده به‌طور خودکار ضرایب برآورد شده‌ای که مقدار کوچکی دارند را برابر صفر قرار دهد تا متغیرهای مناسب انتخاب شوند. این کار پیچیدگی مدل را کاهش می‌دهد.
۲. نا اریبی: برآوردگر به دست آمده برای ضرایب رگرسیونی که مقادیر واقعی آن‌ها بزرگ است، تقریباً نااریب باشد. این ویژگی اریبی مدل را کاهش می‌دهد.
۳. پیوستگی: برآوردگر نتیجه‌شده پیوسته باشد تا بتواند ناپایداری در پیشگویی مدل را کاهش دهد.



شکل ۳.۲: نمودار برآوردهای کمترین توان‌های دوم جریمه‌شده (PLS) در مقابل برآورد کمترین توان‌های دوم معمولی (OLS) وقتی که ماتریس طرح متعامد است. محور افقی OLS و محور عمودی PLS است.

شکل ۳.۲ ارتباط برآورد OLS را با برآوردهای لاسو، ریچ، MCP و SCAD با فرض متعامد بودن ماتریس طرح  $X$  نشان می‌دهد. محور افقی نشان‌دهنده برآورد OLS و محور عمودی برآوردهای کمترین توان‌های دوم جریمه‌شده است. چون برآوردگر OLS ناریب است، انتظار داریم که برای ضرایب بزرگ، برآورد جریمه‌شده برابر یا نزدیک به برآورد OLS باشد تا شرط ناریبی برقرار شود. از طرفی برای برقراری شرط تنکی، لازم است روش کمترین توان‌های دوم جریمه‌شده، ضرایب رگرسیونی کوچک را صفر برآورد کند. همچنین لازم است برآوردگر موردنظر پیوسته باشد. با توجه به شکل، پیوستگی برآوردگر لاسو بدیهی است. همچنین برای ضرایب با مقادیر کوچک، برآورد لاسو برابر صفر است. در نتیجه شرط تنکی برقرار است. اما با بزرگ شدن مقادیر ضرایب، برآورد لاسو از برآورد OLS فاصله می‌گیرد، یعنی شرط ناریبی برقرار نیست. در نتیجه عیب روش لاسو این است که برای ضرایب با مقادیر بزرگ، برآوردهای اریب را نتیجه می‌دهد.

برآوردگر ریچ برخلاف لاسو، علی‌رغم پیوسته بودن، ضرایب کوچک را صفر برآورد نمی‌کند. همچنین با افزایش قدر مطلق ضرایب، مقدار اریبی به شدت افزایش می‌یابد. بنابراین تابع جریمه ریچ یک جریمه مناسب برای انتخاب متغیر نیست.

با توجه به نمودار برآورد SCAD در مقابل برآورد OLS، برقراری ویژگی‌های تنکی و پیوستگی واضح است. علاوه بر این، در ابتدا میزان اریبی برآوردگر جریمه شده با اریبی برآوردگر لاسو معادل است، اما با افزایش قدر مطلق ضرایب، اریبی به سمت صفر تنزل می‌یابد. لذا ویژگی ناریبی نیز برقرار است. همچنین با توجه به شکل، تابع جریمه MCP نیز برآوردگری با ویژگی‌های مطلوب را نتیجه می‌دهد. رفتار تابع جریمه MCP بسیار شبیه SCAD است، با این تفاوت که با فاصله گرفتن ضرایب از صفر، اریبی برآوردگر MCP فوراً به سمت صفر می‌رود، اما اریبی برآوردگر SCAD در یک فاصله معین برابر با اریبی لاسو است و سپس به سرعت کاهش می‌یابد. بنابراین از توابع جریمه MCP و SCAD می‌توان به عنوان توابع جریمه مناسب نام برد.

فن و لی (۲۰۰۱) همچنین نشان دادند که تحت شرایط کافی زیر برآوردگر کمترین توان‌های دوم جریمه‌شده (۱۴.۲) دارای ویژگی‌های مطلوب فوق است:

۱. تنکی:  $\min_{t \geq 0} \{t + p'_\lambda(t)\} > 0$  اگر  $\circ$  تنکی است اگر  $\circ$ .
۲. ناریبی: برآوردگر به دست آمده ناریب است اگر، برای مقادیر بزرگ  $t$ ،  $p'_\lambda(t) = 0$ .
۳. پیوستگی: برآوردگر پیوسته است اگر و فقط اگر  $\circ$   $\arg \min_{t \geq 0} \{t + p'_\lambda(t)\} = 0$ .

همان‌طور که ملاحظه می‌شود، شرایط کافی برای داشتن ویژگی‌های مطلوب بیان شده توسط فن و لی (۲۰۰۱) وابسته به مشتق تابع جریمه هست. علاوه بر این، شکل مشتق تابع جریمه در حل مسئله بهینه‌سازی جریمه‌شده تاثیر بسزایی دارد. مشتق تابع جریمه لاسو به صورت  $p'_\lambda(t) = \lambda$  و مشتق توابع جریمه SCAD و MCP، به ترتیب، به صورت

$$p'_\lambda(t) = \lambda \left\{ I(t \leq \lambda) + \frac{(a\lambda - t)_+}{(a-1)\lambda} I(t > \lambda) \right\}, \quad t \geq 0, a > 2 \quad (30.2)$$

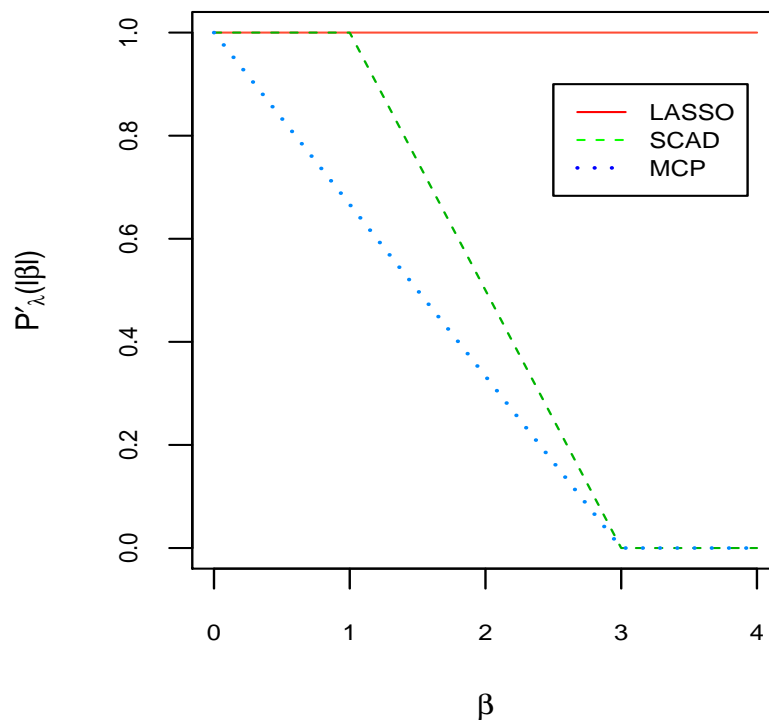
و

$$p'_\lambda(t) = \lambda \left( 1 - \frac{t}{\gamma\lambda} \right)_+, \quad t \geq 0, \gamma > 0 \quad (31.2)$$

به دست می‌آیند. بررسی شرایط کافی ویژگی‌های مطلوب توابع جریمه بسیار ساده است. به عنوان مثال، برای تابع جریمه لاسو،  $\min_{t \geq 0} \{t + p'_\lambda(t)\} = \lambda > 0$  و  $\arg \min_{t \geq 0} \{t + p'_\lambda(t)\} = 0$ ، اما  $p'_\lambda(t) = \lambda > 0$ . بنابراین تابع جریمه لاسو دارای ویژگی ناریبی نیست.

نمودار مشتق توابع جریمه لاسو، SCAD و MCP به‌ازای  $\lambda = 1$  و  $a = \gamma = 3$  در شکل ۴.۲ داده شده است. با توجه به شکل می‌بینیم که نرخ جریمه توابع SCAD و MCP در مبدأ با نرخ جریمه تابع لاسو برابر است، اما با افزایش قدرمطلق ضرایب، نرخ جریمه توابع SCAD و MCP به سمت صفر نزول می‌کند. لازم به ذکر است که نرخ جریمه تابع SCAD در ابتدا با فاصله گرفتن ضرایب از صفر معادل نرخ جریمه لاسو است و سپس با افزایش قدرمطلق ضرایب به سرعت کاهش می‌یابد، در حالی که تابع جریمه MCP به محض خارج شدن از مبدأ، نرخ جریمه آن فوراً به سمت صفر تنزل می‌یابد.

شکل ۴.۲ همچنین نشان می‌دهد که روش MCP مقعر مینیماکس است، یعنی در بین همه توابع جریمه پیوسته مشتق‌پذیر بر بازه  $(0, \infty)$  که برای هر  $t \geq \gamma\lambda$  در دو شرط  $p'_\lambda(t) = 0$  و  $p'_\lambda(0+) = \lambda$



شکل ۴.۲: نمودارهای مشتق توابع جریمه لاسو، SCAD و MCP.

صدق می‌کنند، MCP ماکزیمم تقعر زیر را مینیمم می‌کند

$$\mathcal{K} = \sup_{0 < t_1 < t_2} \frac{p'_\lambda(t_1) - p'_\lambda(t_2)}{t_2 - t_1}. \quad (32.2)$$

به خاطر داشته باشید که اغلب توابع جریمه‌ای که شرایط مطلوب تنکی، نارایی و پیوستگی را دارا هستند، غیر محدب هستند. این موضوع، مسئله بهینه‌سازی را با چالش مواجه می‌کند. مقعر مینماکس بودن تابع MCP بدین معنی است که در بین همه توابعی که دارای سه ویژگی مطلوب فوق هستند، MCP اندازه ماکزیمم تقعر (۳۲.۲) را مینیمم می‌کند. در شکل ۴.۲، مشتق‌های توابع SCAD و MCP در نقاط  $\gamma\lambda = 3$  و  $\gamma\lambda = 3$  برابر هستند، اما MCP در این ناحیه دارای تقعر  $\frac{1}{3} = \frac{1}{\gamma}$  است، در حالی که SCAD دارای ماکزیمم تقعر  $\frac{1}{\gamma-1} = \frac{1}{3}$  است. برای اطلاعات بیشتر به برهنی و هوانگ (۲۰۱۹) مراجعه شود.

### ۳.۴.۲ انتخاب پارامتر جریمه

انتخاب پارامتر جریمه در رگرسیون جریمه‌شده نقش بسیار مهمی را در کنترل پیچیدگی مدل بازی می‌کند. به عنوان مثال، به‌ازای  $\lambda = 0$  برآوردگر جریمه‌شده به برآوردگر OLS تبدیل می‌شود و همه متغیرهای توضیحی وارد مدل می‌شوند. اگر  $\lambda = \infty$ ، جزء دوم تابع (۱۴.۲) برابر با  $\infty$  می‌شود، و در نتیجه تمام ضرایب صفر شده و هیچ متغیری وارد مدل نمی‌شود. برای انتخاب بهترین مقدار بین این دو حالت فرین، از معیارهای کلاسیک مانند AIC، BIC و GCV استفاده می‌شود. به‌طور مشخص، برای

$\lambda$  داده شده، برآورد کمترین توان‌های دوم جریمه شده،  $\hat{\beta}_\lambda$ ، را به دست آورده و مقدار معیار کلاسیک انتخاب مدل را بر اساس  $\hat{\beta}_\lambda$  محاسبه می‌کنیم. این معیارهای کلاسیک به صورت زیر بیان می‌شوند

$$AIC_\lambda = \| \mathbf{Y} - \mathbf{X}\hat{\beta}_\lambda \|^2 + 2df_\lambda \hat{\sigma}^2$$

$$BIC_\lambda = \| \mathbf{Y} - \mathbf{X}\hat{\beta}_\lambda \|^2 + \log(n)df_\lambda \hat{\sigma}^2$$

$$GCV_\lambda = \frac{\frac{1}{n} \| \mathbf{Y} - \mathbf{X}\hat{\beta}_\lambda \|^2}{(1 - df_\lambda/n)^2}$$

که  $\hat{\sigma}^2 = \frac{RSS_p}{n-p}$  و  $RSS_p$  مجموع توان‌های دوم باقیمانده‌ها تحت مدل کامل است. منظور از مدل کامل مدلی است که همه متغیرهای توضیحی در مدل حضور دارند. در معیارهای فوق، درجه آزادی مدل برآورد شده  $df_\lambda$  معمولاً با استفاده از تعداد ضرایب غیر صفر به صورت

$$df_\lambda = \sum_{j=1}^p I(\hat{\beta}_{j,\lambda} \neq 0)$$

به دست می‌آید که در آن  $I(\cdot)$  تابع نشانگر است. راه دیگر محاسبه  $df_\lambda$  استفاده از ماتریس طرح به صورت

$$df_\lambda = \text{tr}(\mathbf{X}_\lambda(\mathbf{X}_\lambda^T \mathbf{X}_\lambda + n\boldsymbol{\Sigma}_\lambda)^{-1} \mathbf{X}_\lambda^T)$$

است که در آن ماتریس طرح مدل متناظر با  $\lambda$  داده شده و  $\boldsymbol{\Sigma}_\lambda$  ماتریس بلوکی به صورت

$$\boldsymbol{\Sigma}_\lambda = \text{diag}_{\hat{\beta}_{j,\lambda} \neq 0} \{P'_\lambda(|\hat{\beta}_{j,\lambda}|)/|\hat{\beta}_{j,\lambda}|\}$$

است که  $\hat{\beta}_{j,\lambda}$ ،  $j$  - امین مولفه از برآورد کمترین توان‌های دوم جریمه شده  $\hat{\beta}_\lambda$  است. در پایان، برای یک مجموعه از نقاط داده شده  $\lambda_1, \dots, \lambda_M$ ، مقداری از  $\lambda$  را به عنوان مقدار بهینه انتخاب می‌کنیم که دارای کمترین مقدار معیار کلاسیک باشد.

یک روش دیگر برای انتخاب پارامتر جریمه استفاده از روش اعتبارسنجی متقابل  $k$  گروهی<sup>۳۴</sup> است. در این روش، برای یک  $\lambda$  داده شده، مجموعه مشاهدات را به طور تصادفی به  $k$  گروه با اندازه یکسان افزایش می‌کنیم. سپس گروه اول را به عنوان داده‌های آزمون<sup>۳۵</sup> و سایر گروه‌ها را به عنوان داده‌های آموزشی<sup>۳۶</sup> در نظر می‌گیریم. پس از برازش مدل با استفاده از داده‌های آموزشی، مقدار متغیر پاسخ به ازای داده‌های آزمون پیشگویی می‌شود و خطای پیشگویی به صورت  $PE_1 = \sum_{i \in I_1} (y_i - \hat{y}_i)^2$  به دست می‌آید که  $I_1$  اندیس مشاهدات گروه اول است. این روند  $k$  بار تکرار می‌شود و در هر تکرار یکی از گروه‌ها به عنوان داده‌های آزمون و بقیه گروه‌ها به عنوان داده‌های آموزشی در نظر گرفته می‌شوند. در نتیجه  $k$  برآورد خطای آزمون  $PE_1, \dots, PE_k$  به دست آمده و در نهایت، امتیاز اعتبارسنجی متقابل به صورت

$$CV(\lambda) = \frac{1}{k} \sum_{i=1}^k PE_i$$

<sup>34</sup>k-fold cross validation

<sup>35</sup>Test dataset

<sup>36</sup>Training dataset

محاسبه می‌شود. این کار برای  $\lambda$  های مختلف تکرار شده و سپس  $\lambda$  ای انتخاب می‌شود که مقدار  $CV(\lambda)$  را مینیمم کند. در عمل معمولاً مقدار  $k$  برابر ۵ یا ۱۰ در نظر گرفته می‌شود. یکی از مشکلات روش لاسو، که قبلاً به آن اشاره شد، این است که در بین یک گروه از متغیرهای توضیحی با همبستگی بالا، تنها یک متغیر را انتخاب می‌کند. بنابراین، این روش توانایی انتخاب متغیر گروهی را ندارد. در عمل ممکن است داده‌ها ساختار گروهی داشته باشند و لازم است متغیرها به صورت گروهی انتخاب شوند. در ادامه، چند روش انتخاب متغیر گروهی را معرفی می‌کنیم.

## ۵.۲ انتخاب متغیر گروهی

مدل رگرسیون خطی با  $p$  متغیر توضیحی را در نظر بگیرید. فرض کنید متغیرهای توضیحی به طور ذاتی به  $J$  گروه مجزا افزاز شده باشند و بتوان مدل را به شکل

$$Y = \sum_{j=1}^J X_j \beta_j + \varepsilon \quad (33.2)$$

نوشت که  $Y = (y_1, \dots, y_n)^T$  بردار مشاهدات متغیر پاسخ،  $X_j$ ، با بعد  $n \times p_j$ ، ماتریس طرح  $p_j$  متغیر توضیحی در گروه  $j$ -ام،  $\beta_j = (\beta_{j1}, \dots, \beta_{jp_j})^T$  ضرایب رگرسیونی گروه  $j$ -ام و  $\varepsilon$  بردار خطا است. هنگام برآورد و انتخاب متغیرهای مهم، هر گروه از متغیرها را به عنوان یک واحد در نظر می‌گیرند و این عمل فرآیند انتخاب متغیر را تسهیل می‌کند. در نتیجه متغیرها به صورت گروهی از مدل حذف شده یا در مدل باقی می‌مانند. ساختار گروهی متغیرها می‌تواند دلایل و اهداف متفاوتی داشته باشد. نمایش یک متغیر کیفی چند سطحی با استفاده از متغیرهای ظاهری و نمایش اثر یک متغیر تصادفی پیوسته با استفاده از توابع پایه‌ای، نمونه‌هایی از متغیرهای گروه‌بندی شده هستند. گروه‌بندی همچنین می‌تواند به منظور بهره‌گیری از دانش پیشین، که از نظر علمی معنی‌دار است، مورد استفاده قرار گیرد. به عنوان مثال، در تحلیل داده‌های بیان ژنی<sup>۳۷</sup>، ژن‌های متعلق به یک مسیر بیولوژیکی یکسان می‌توانند در یک گروه قرار گیرند. در مطالعات پیوند ژنتیکی، نشانگرهای ژنتیکی از ژن یکسان می‌توانند به عنوان یک گروه در نظر گرفته شوند. در عمل مطلوب است ساختار گروهی را در تحلیل این نوع داده‌ها مورد توجه قرار دهیم. در ادامه چند تابع جریمه برای انتخاب متغیرهای گروهی معرفی می‌کنیم.

### ۱.۵.۲ جریمه لاسوی گروهی

تابع جریمه لاسوی گروهی توسط یوان و لین<sup>۳۸</sup> (۲۰۰۶) معرفی شد و زمانی به کار می‌رود که متغیرهای توضیحی از قبل گروه‌بندی شده باشند. برآوردگر لاسوی گروهی با مینیمم کردن تابع هدف جریمه شده

$$\mathcal{L}(\beta) = \frac{1}{p} \left\| Y - \sum_{j=1}^J X_j \beta_j \right\|^2 + \lambda \sum_{j=1}^J \|\beta_j\|_{R_j} \quad (34.2)$$

<sup>37</sup>Gene expression data

<sup>38</sup>Yuan and Lin

به دست می‌آید که  $\lambda$  پارامتر جریمه،  $\beta = (\beta_1^T, \dots, \beta_J^T)^T$  بردار ضرایب،  $R_j$  یک ماتریس معین مثبت با بعد  $p_j \times p_j$  و  $\|\beta_j\|_{R_j} = (\beta_j^T R_j \beta_j)^{\frac{1}{2}}$  است.

یک مسئله مهم، نحوه انتخاب ماتریس  $R_j$  است. در حالتی که  $X_j^T X_j = nI_{p_j}$ ، یوان و لین (۲۰۰۶)،  $R_j = I_{p_j}$  را پیشنهاد دادند. با این وجود، انتخاب  $R_j = I_{p_j}$  چندان مناسب نیست، زیرا مقیاس متغیرهای توضیحی ممکن است یکسان نباشد. در حالت کلی، یک انتخاب مناسب  $R_j$ ، ماتریس گرام<sup>۳۹</sup> بر اساس  $X_j$  یعنی  $X_j^T X_j / n$  است. این انتخاب  $R_j$  معادل با استاندارد کردن متغیرها در هر گروه است.

برای برازش رگرسیون لاسوی گروهی با الگوریتم ارائه شده توسط یوان و لین (۲۰۰۶)، لازم است که  $X_j$  ها متعامد باشند. بدین منظور روش متعامدسازی گرام-اشمیت می‌تواند مورد استفاده قرار گیرد. ویژگی متعامد بودن، توسعه الگوریتم‌های محاسباتی برای برازش رگرسیون جریمه شده را تسهیل می‌کند و باعث کاهش بار محاسباتی می‌شود (برهنی و هوانگ، ۲۰۱۵). بعلاوه، این ویژگی امکان استفاده از الگوریتم‌های ارائه شده برای لاسوی گروهی را برای برازش رگرسیون SCAD گروهی و MCP گروهی فراهم می‌کند. لازم به ذکر است که بدون از دست دادن کلیت می‌توانیم فرض متعامد بودن متغیرهای درون هر گروه را بپذیریم، زیرا پس از برآورد ضرایب رگرسیونی مدل استاندارد، می‌توانیم با تبدیل ضرایب رگرسیونی مدل اصلی را به دست آوریم. در ادامه به بیان جزئیات این موضوع می‌پردازیم.

با استفاده از تجزیه چولسکی، ماتریس گرام را به صورت  $R_j = U_j^T U_j$  در نظر می‌گیریم که  $U_j$  یک ماتریس بالا مثلثی با بعد  $p_j \times p_j$  است. فرض کنید  $\tilde{X}_j = X_j U_j^{-1}$  و  $\tilde{\mathbf{b}}_j = U_j \mathbf{b}_j$ . بنابراین، می‌توان تابع جریمه شده (۳۴.۲) را به صورت

$$\mathcal{L}(\beta) = \frac{1}{2} \left\| \mathbf{Y} - \sum_{j=1}^J \tilde{X}_j \tilde{\mathbf{b}}_j \right\|^2 + \lambda \sum_{j=1}^J \|\tilde{\mathbf{b}}_j\| \quad (35.2)$$

نوشت. پس از به دست آوردن جواب با مینیمم کردن تابع (۳۵.۲)، می‌توانیم جواب مسئله اصلی (۳۴.۲) را با تبدیل  $\beta_j = U_j^{-1} \tilde{\mathbf{b}}_j$  به دست آوریم. با توجه به تعریف  $U_j$  نتیجه می‌شود

$$\begin{aligned} \tilde{X}_j^T \tilde{X}_j &= U_j^{T-1} X_j^T X_j U_j^{-1} \\ &= U_j^{T-1} n R_j U_j^{-1} \\ &= n U_j^{T-1} U_j^T U_j U_j^{-1} \\ &= n I_{p_j}. \end{aligned}$$

بنابراین با این انتخاب  $R_j$  و بدون از دست دادن کلیت، فرض می‌کنیم که به ازای  $1 \leq j \leq J$ ،  $X_j^T X_j = n I_{p_j}$ . توجه کنید که در اینجا فرض نکرده‌ایم که برای  $j \neq k$ ،  $X_j$  و  $X_k$  متعامد هستند. هنگامی که  $p_j = 1$  ( $1 \leq j \leq J$ )، لاسوی گروهی به لاسوی استاندارد تبدیل می‌شود. در این حالت،  $R_j = \frac{\|X_j\|^2}{n}$  معادل با واریانس نمونه‌ای  $j$ -امین متغیر توضیحی است. بنابراین، در نظر گرفتن  $R_j$  به صورت ماتریس گرام معادل استاندارد کردن متغیرهای توضیحی قبل از تحلیل است که اغلب برای برازش رگرسیون لاسو توصیه می‌شود.

<sup>39</sup>Gram matrix



## ۲.۵.۲ لاسوی گروهی تطبیقی

همان‌طور که در بخش ۴.۲ بیان شد، برآوردگر لاسو دارای ویژگی پیشگویی نیست. لنگ<sup>۴۰</sup> و همکاران (۲۰۰۶) نشان دادند که اگر پارامتر جریمه با مینیمم کردن خطای پیشگویی انتخاب شود، لاسو دارای ویژگی سازگاری در انتخاب نیست. دلیل این مشکل این است که لاسو برای تمام ضرایب رگرسیونی میزان انقباض یکسانی را در نظر می‌گیرد. مطلوب این است که برای ضرایب بزرگتر انقباض کوچکتر و برای ضرایب نزدیک به صفر انقباض بزرگتری را در نظر بگیرد. این موضوع باعث می‌شود که لاسو مدلی بزرگتر از مدل درست را انتخاب کند، چون قادر به تشخیص ضرایب کوچک از ضرایب صفر نیست. تابع جریمه لاسوی گروهی هم رفتاری شبیه لاسو دارد و گروه‌های بی‌اهمیت زیادی را وارد مدل می‌کند. برای بهبود عملکرد لاسوی گروهی، تابع مجموع توان‌های دوم جریمه‌شده به صورت

$$\mathcal{L}(\beta) = \frac{1}{p} \left\| \mathbf{Y} - \sum_{j=1}^J \mathbf{X}_j \beta_j \right\|^2 + \lambda \sum_{j=1}^J \omega_j \|\beta_j\|_{\mathbf{R}_j} \quad (۳۶.۲)$$

تعریف می‌شود که  $\omega_j = (\|\tilde{\beta}_j\|_{\mathbf{R}_j})^{-1}$  یک مقدار اولیه برای  $\beta_j$  است. برآورد لاسوی گروهی می‌تواند به عنوان مقدار اولیه در این روش استفاده شود. همان‌طور که می‌بینید لاسوی گروهی تطبیقی برای گروه‌های مختلف، جریمه‌های متفاوتی را در نظر می‌گیرد.

## ۳.۵.۲ سایر توابع جریمه گروهی

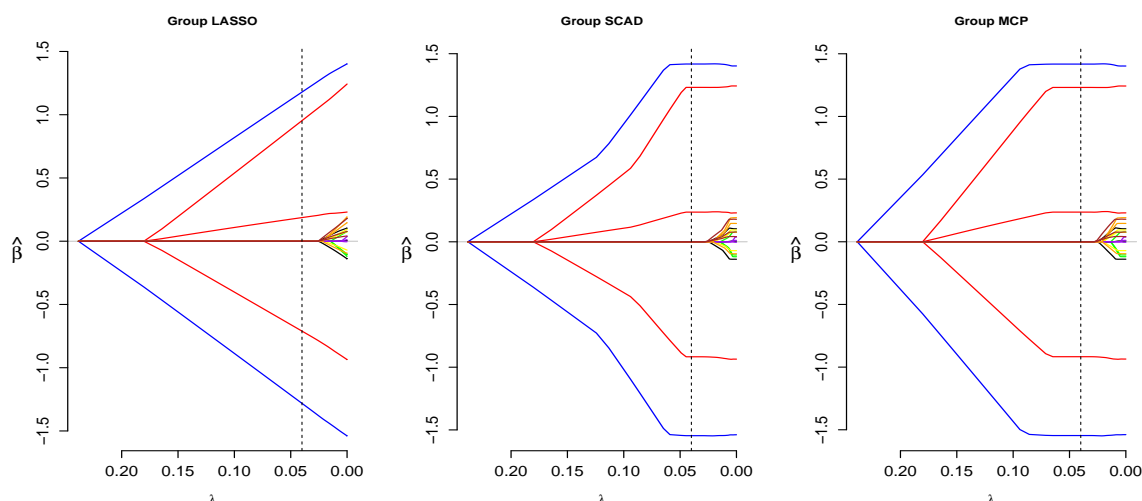
واضح است که تابع جریمه گروهی، با قرار دادن  $\|\beta_j\|_{\mathbf{R}_j}$  به جای  $|\beta_j|$  در تابع جریمه معمولی به دست می‌آید. حال اگر  $p_\lambda(t)$  یک تابع جریمه باشد، یک روش انتخاب متغیر گروهی بر اساس معیار

$$\mathcal{L}(\beta) = \frac{1}{p} \left\| \mathbf{Y} - \sum_{j=1}^J \mathbf{X}_j \beta_j \right\|^2 + \sum_{j=1}^J p_{\lambda, \gamma}(\|\beta_j\|_{\mathbf{R}_j}) \quad (۳۷.۲)$$

به دست می‌آید. اینجا  $\gamma$  یک پارامتر اضافی است که می‌تواند برای اصلاح  $p_\lambda(\cdot)$  استفاده شود. بنابراین می‌توان مشابه لاسوی گروهی، جریمه‌های SCAD گروهی، MCP گروهی و بریج گروهی را نیز تعریف کرد.

برای درک بهتر مفهوم انتخاب متغیر گروهی، نمونه‌ای با  $J = ۲۰$  گروه در نظر بگیرید به طوری که دو گروه اول دارای ضرایب غیر صفر و ضرایب مربوط به سایر گروه‌ها صفر می‌باشند. فرض کنید  $\beta_1 = (-\sqrt{2}, \sqrt{2})^T$  و  $\beta_2 = (-0.5, 1, 0.5)^T$ . همچنین فرض کنید اندازه گروه‌های صفر برابر ۳ باشد. شکل ۵.۲ برآوردهای لاسوی گروهی، SCAD گروهی و MCP گروهی را برای مقادیر مختلف  $\lambda$  نشان می‌دهد. با توجه به شکل می‌توان دید که، برای هر مقدار  $\lambda$ ، همه ضرایب متعلق به یک گروه، یا برابر صفر برآورد می‌شوند، یا همه آن ضرایب دارای مقداری غیر صفر هستند. در نتیجه، همه متغیرهای یک گروه یا در مدل قرار می‌گیرند، یا همگی از مدل حذف می‌شوند. این موضوع، مفهوم انتخاب گروهی را بیان می‌کند. علاوه بر این برای SCAD گروهی و MCP گروهی، بازه‌ای از  $\lambda$  وجود دارد به طوری که

<sup>40</sup>Leng

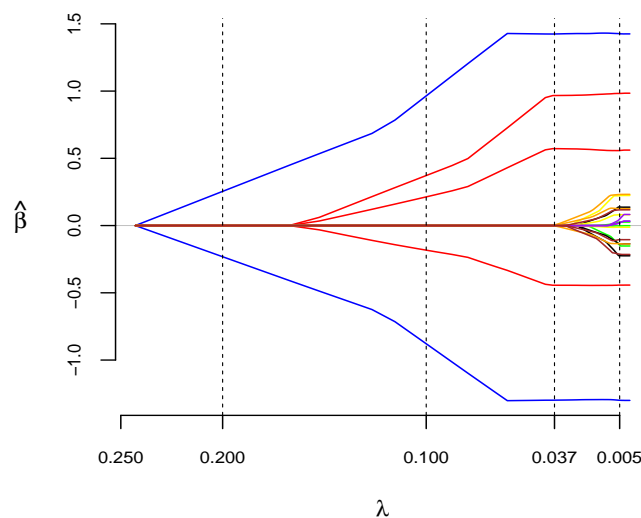


شکل ۵.۲: برآوردهای لاسوی گروهی، SCAD گروهی و MCP گروهی برای مقادیر مختلف  $\lambda$ . خط عمودی نقطه‌چین مقدار  $\lambda = 0.04$  را نشان می‌دهد. نمودارهای آبی و قرمز، به ترتیب، برآوردهای ضرایب گروه اول و دوم و بقیه نمودارها مربوط به گروه‌های صفر هستند.

در آن بازه، برآوردها به مقادیر واقعی پارامترها بسیار نزدیک هستند. به‌عنوان مثال، به ازای  $\lambda = 0.04$ ، MCP گروهی و SCAD گروهی همان برآوردهای کمترین توان‌های دوم را تولید می‌کنند و مدلی را نتیجه می‌دهند که تنها شامل متغیرهای مهم است. این مدل را مدل پیشگو<sup>۴۱</sup> گویند. اما برای لاسوی گروهی، برآوردها همواره اریب هستند، مگر به ازای  $\lambda = 0$ . در این حالت برآورد لاسو با برآورد کمترین توان‌های دوم معادل است.

برای بررسی دقیق‌تر، شکل ۶.۲ را برای تابع جریمه SCAD گروهی در نظر بگیرید. با توجه به این شکل، برای  $\lambda = 0.05$ ، تقریباً متغیرهای توضیحی همه گروه‌ها در مدل حضور دارند، سپس با افزایش  $\lambda$ ، همه ضرایب به سمت صفر منقبض می‌شوند. به‌عنوان مثال، به ازای  $\lambda = 0.10$  فقط متغیرهای گروه‌های مهم در مدل حضور دارند، اما اریبی برآوردهای به‌دست آمده زیاد است. به ازای  $\lambda = 0.20$  تنها متغیرهای گروه اول در مدل حضور دارند و متغیرهای گروه دوم، که متغیرهای مؤثر بر پاسخ هستند، از مدل حذف می‌شوند. در حالت خاص  $\lambda = 0.37$  همه ضرایب گروه‌های بی‌اهمیت صفر می‌شوند، درحالی‌که برای گروه‌های مهم، برآوردهای به‌دست آمده با SCAD گروهی با برآوردهای کمترین توان‌های دوم معادل هستند. به عبارت دیگر، به ازای  $\lambda = 0.37$  برآوردگر SCAD گروهی به خوبی برآوردگر پیشگو است. بنابراین  $\lambda = 0.37$  می‌تواند یک مقدار بهینه برای  $\lambda$  باشد. البته در عمل مقدار بهینه  $\lambda$  با استفاده از روش‌های معرفی‌شده در زیربخش ۳.۴.۲ به‌دست می‌آید.

<sup>41</sup>Oracle model



شکل ۶.۲: برآوردهای SCAD گروهی برای مقادیر مختلف  $\lambda$ .

## ۶.۲ انتخاب متغیر در داده‌های با بعد بسیار بالا

به‌طور کلی، داده‌ها را از نظر بعد می‌توان به دو گروه تقسیم کرد: ۱- بعد پایین<sup>۴۲</sup> ۲- بعد بالا<sup>۴۳</sup>. در داده‌های با بعد پایین، تعداد متغیرها،  $p$ ، بسیار کوچکتر از تعداد مشاهدات،  $n$ ، است. به‌عنوان مثال، مجموعه داده با حجم  $n = 100$  و بعد  $p = 10$  یک مجموعه داده با بعد پایین است. بعلاوه، در این داده‌ها، برای نظریه‌های جانبی فرض می‌شود که با افزایش  $n$ ، مقدار  $p$  ثابت می‌ماند. در داده‌های با بعد بالا،  $p$  از  $n$  بزرگتر است. به‌عنوان مثال،  $n = 100$  و  $p = 200$  تحقق از بعد بالا است. در این نوع داده‌ها با افزایش  $n$ ، مقدار  $p$  می‌تواند افزایش یابد. امروزه محققین در زمینه‌های مختلف علمی و تحقیقاتی مانند ژنومیک، علوم اعصاب، علوم زمین و مالی با این نوع داده‌ها مواجه هستند. اگرچه محققین در زمینه‌های مختلف اهداف متفاوتی را دنبال می‌کنند، اما یک موضوع مشترک دارند: تحقیق آن‌ها به شدت به استخراج اطلاعات مفید از داده‌های حجیم وابسته است و تعداد متغیرهای توضیحی در مقایسه با حجم نمونه می‌تواند بزرگ باشد. مسئله “ $p$  بزرگ،  $n$  کوچک” چالش‌های زیادی را پیش روی آماردانان قرار داده است و تحلیل این نوع داده‌ها مستلزم به کارگیری روش‌های نوین آماری است. یک فرض اساسی که استنباط پیرامون این نوع داده‌ها را ممکن می‌سازد، فرض تنگی<sup>۴۴</sup> است که بر اساس این فرض، فقط مجموعه کوچکی از متغیرهای توضیحی بر متغیر پاسخ تأثیرگذارند و بقیه متغیرهای توضیحی ارتباطی با متغیر پاسخ ندارند. بنابراین کاهش بعد و انتخاب متغیر نقش مهمی را در تحلیل این نوع داده‌ها ایفا می‌کند.

تاکنون روش‌های آماری مختلفی برای انتخاب متغیر در داده‌های با بعد بالا معرفی شده‌اند که از

<sup>۴۲</sup>Low dimension

<sup>۴۳</sup>High dimension

<sup>۴۴</sup>Sparsity assumption

جمله می‌توان لاسو، لاسوی تطبیقی، SCAD، الاستیک نت، گاروت نامنی<sup>۴۵</sup> و انتخاب گر دن‌زیگ<sup>۴۶</sup> را نام برد. همه این روش‌ها می‌توانند برای تحلیل داده‌هایی به کار روند که تعداد متغیرهای پیشگو،  $p$ ، بزرگتر از حجم نمونه،  $n$ ، است.

در حالت خاص، ممکن است  $p$  با نرخ نمایی رشد کند، به عبارت دیگر، رابطه  $\log(p) = O(n^a)$  برای برخی مقادیر مثبت  $a$  برقرار باشد. داده‌های بیان ژنی، ریز آرایه‌ها، داده‌های تابعی، داده‌های مالی، داده‌های پردازش تصویر، داده‌های تصویربرداری پزشکی و طبقه‌بندی تومور نمونه‌هایی از این نوع داده‌ها هستند. به‌عنوان مثال، در مطالعه رابطه بین ویژگی‌های ظاهری<sup>۴۷</sup> و ویژگی‌های ژنتیکی<sup>۴۸</sup> از میلیون‌ها SNP<sup>۴۹</sup> استفاده می‌کنند، یا در طبقه‌بندی نوعی تومور از داده‌های ریز آرایه استفاده می‌کنند که شامل هزاران عامل ژنتیکی است و هدف شناسایی ژن‌های موثر بر بیماری است. در این مثال‌ها، تعداد متغیرهای توضیحی بسیار بزرگتر از تعداد مشاهدات است. برای تمایز این نوع داده‌ها از داده‌های با بعد بالای  $p$  با نرخ رشد چندجمله‌ای، آن‌ها را داده‌های با بعد بسیار بالا<sup>۵۰</sup> می‌نامند. بنابراین، در این رساله، عبارت “بعد بالا” برای نرخ رشد چندجمله‌ای ( $p = O(n^b)$ ) و عبارت “بعد بسیار بالا” برای نرخ رشد نمایی ( $p = \exp\{O(n^a)\}$ ) استفاده می‌شود. برای درک بهتر مطلب، دو حالت  $p = n^b$  و  $p = \exp\{n^a\}$  را به ترتیب برای بعد بالا و بعد بسیار بالا در نظر گرفته و نسبت  $\frac{p}{n}$  را برای مقادیر مختلف  $(n, a, b)$  در جدول ۱.۲ گزارش کرده‌ایم. با توجه به جدول ۱.۲ ملاحظه می‌شود که در بعد بسیار بالا، مقدار  $p$  با افزایش  $n$  به سرعت رشد می‌کند و به سمت بینهایت می‌رود. برای آگاهی بیشتر درباره داده‌های با بعد بالا و چالش‌های تحلیل این نوع داده‌ها به فن و لی (۲۰۰۶)، فن و لیو<sup>۵۱</sup> (۲۰۰۸)، فن و همکاران (۲۰۰۹)، فن و همکاران (۲۰۱۴)، جان‌استون و تیتترینگتون<sup>۵۲</sup> (۲۰۰۹)، بولمن و فن دگیر<sup>۵۳</sup> (۲۰۱۱) و گیراد<sup>۵۴</sup> (۲۰۱۴) مراجعه کنید.

تحلیل داده‌های با بعد بسیار بالا با سه چالش همزمان هزینه محاسباتی، دقت آماری و پایداری الگوریتمی روبروست (فن و لیو، ۲۰۰۹). این چالش‌ها کاربرد مستقیم روش‌های انتخاب متغیر اشاره شده را برای داده‌های با بعد بسیار بالا با مشکل مواجه می‌کند. قالب کلی پیشنهاد شده برای رفع این مشکل یک طرح دو مرحله‌ای است که در مرحله اول، بعد مدل توسط یک روش غربالگری کاهش یافته و در مرحله دوم از روش‌های انتخاب متغیر مبتنی بر جریمه، که در بخش‌های ۴.۲ و ۵.۲ به آن‌ها اشاره شد، برای انتخاب متغیر از زیر مدل غربال شده استفاده می‌شود. بدیهی است یکی از این روش‌ها که هزینه محاسباتی چندانی ندارد، غربالگری بر اساس همبستگی متغیرهای توضیحی با متغیر پاسخ است. در این راستا، یک راهکار ساده می‌تواند این باشد که اگر همبستگی یک متغیر توضیحی با پاسخ کمتر

<sup>45</sup> Nonnegative garrote

<sup>46</sup> Dantzig selector

<sup>47</sup> Phenotypes

<sup>48</sup> Genotypes

<sup>49</sup> Single-nucleotide polymorphism

<sup>50</sup> Ultrahigh-dimensional data

<sup>51</sup> Lv

<sup>52</sup> Johnstone and Titterington

<sup>53</sup> Buhlmann and van de Geer

<sup>54</sup> Giraud

جدول ۱.۲: مقادیر  $\frac{p}{n}$  برای داده‌های با ابعاد بالا و بسیار بالا.

بعد بسیار بالا (مقادیر مختلف $a$ )				بعد بالا (مقادیر مختلف $b$ )				$n$
۱	$\frac{7}{8}$	$\frac{5}{8}$	$\frac{3}{8}$	۳	۲.۵	۲	۱.۵	
۲۲۰۲/۶۴	۱۸۰/۶۱	۶/۷۸	۱/۰۷	۱۰۰	۳۱/۶۲	۱۰	۳/۱۶	۱۰
۳۵۶۲۱۵۸۱۹۳۸۴	۱۰۹۵۱۲۱۴	۱۴۵/۱۷	۱/۱۹	۹۰۰	۱۶۴/۳۱	۳۰	۵/۴۷	۳۰
$۲,۶۸۸۱۱۷e + ۴۱$	$۲,۶۴۳۴۶۴e + ۲۲$	۵۲۸۴۰۸	۲/۷۶	۱۰۰۰۰	۱۰۰۰	۱۰۰	۱۰	۱۰۰
$۳,۶۱۲۹۸۷e + ۸۴$	$۳,۰۸۶۲۶۸e + ۴۲$	۴۰۶۸۳۳۷۱۹۴	۷/۳۴	۴۰۰۰۰	۲۸۲۸/۴۲	۲۰۰	۱۴/۱۴	۲۰۰
$۲,۸۰۷۱۸۴e + ۲۱۴$	$۱,۴۴۲۳۶۹e + ۹۷$	$۲,۶۲۰۸۱۱e + ۱۸$	۵۸/۴۵	۲۵۰۰۰۰	۱۱۱۸۰/۳۴	۵۰۰	۲۲/۳۶	۵۰۰

از یک حد آستانه از پیش تعیین شده باشد، آن متغیر توضیحی را از مدل حذف می‌کنیم. براساس این راهکار، فن و لیو (۲۰۰۸) نخستین روش غربالگری را با استفاده از همبستگی پیرسن برای مدل‌های خطی معرفی کردند و آن را غربالگری مستقل مطمئن (SIS) نامگذاری کردند. تاکنون، علاوه بر SIS، چندین روش غربالگری دیگر نیز برای مدل‌های آماری مختلف معرفی شده‌اند که در ادامه به توصیف دو روش SIRS و DC-SIS می‌پردازیم.

## ۱.۶.۲ غربالگری مستقل مطمئن

مدل خطی زیر را در نظر بگیرید

$$Y = X\beta + \varepsilon \quad (38.2)$$

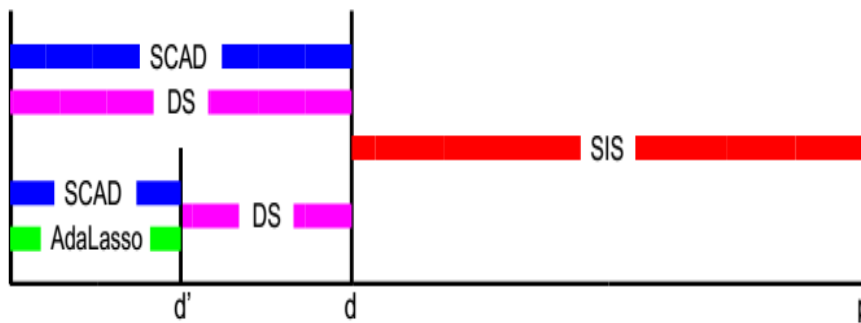
که  $Y = (y_1, \dots, y_n)^T$  متغیر پاسخ،  $X = (x_1, \dots, x_n)^T$  ماتریس طرح و  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$  خطاهای تصادفی می‌باشند. فرض کنید  $M = \{k : \beta_k \neq 0\}$  مجموعه اندیس متغیرهای مهم (مدل درست) با اندازه  $|M| = s$  و  $\hat{M}$  مجموعه اندیس متغیرهای انتخاب شده با اندازه  $|\hat{M}|$  باشد. برای رتبه‌بندی متغیرها، ابتدا بردار  $p$  بعدی  $\omega = (\omega_1, \dots, \omega_p)^T$  را به صورت

$$\omega_k = \frac{1}{n} X_k^T Y, \quad k = 1, \dots, p \quad (39.2)$$

در نظر بگیرید که متغیرهای توضیحی و پاسخ استاندارد شده‌اند، به طوری که، برای  $k = 1, \dots, p$ ،  $\frac{1}{n} \sum_{i=1}^n x_{ik}^2 = 1$ ،  $\frac{1}{n} \sum_{i=1}^n x_{ik} = 0$  و  $\frac{1}{n} \sum_{i=1}^n y_i^2 = 1$ ،  $\frac{1}{n} \sum_{i=1}^n y_i = 0$  همان ضریب همبستگی نمونه‌ای بین  $k$ -امین متغیر توضیحی و متغیر پاسخ است.

روش غربالگری مستقل مطمئن، با نام اختصاری SIS، متغیرهای توضیحی را بر اساس مقادیر  $|\omega_k|$  رتبه‌بندی می‌کند و متغیرهایی را انتخاب می‌کند که دارای همبستگی بالایی با متغیر پاسخ باشند. به طور واضح‌تر، برای هر مقدار  $\gamma \in (0, 1)$  مولفه‌های  $|\omega_k|$  را به ترتیب صعودی مرتب کرده و سپس زیر مدل  $\hat{M}$  را به صورت

$$\hat{M} = \{k : |\omega_k| \text{ در بین } [\gamma n] \text{ مقدار بزرگ } |\omega| \text{ قرار دارد.}\} \quad (40.2)$$



شکل ۷.۲: روش‌های انتخاب متغیر در ابعاد خیلی بالا

انتخاب می‌کنیم که  $\gamma$  یک مقدار ثابت از پیش تعیین شده،  $[\gamma n]$  جزء صحیح  $\gamma n$  و  $R_k = |\omega_k|$  معیار مطلوبیت حاشیه‌ای<sup>۵۵</sup> نامیده می‌شود. به وضوح، انتخاب مقدار بزرگتری برای  $\gamma$  منجر به انتخاب تعداد بیشتری از متغیرها می‌شود. در عمل، انتخاب یک مقدار مناسب برای  $\gamma$  چندان ساده به نظر نمی‌رسد. در این راستا، فن و لیو (۲۰۰۸) انتخاب  $d$  متغیر با بزرگترین معیار مطلوبیت حاشیه‌ای را به‌عنوان یک روش جایگزین پیشنهاد دادند که  $d$  به صورت  $d = [n/\log(n)]$  یا  $d = 2[n/\log(n)]$  تعیین می‌شود. فن و لیو (۲۰۰۸) تحت قضیه‌ای نشان دادند که روش SIS دارای ویژگی غربالگری مطمئن<sup>۵۶</sup> است، یعنی این روش با احتمال نزدیک به یک همه متغیرهای مهم را انتخاب می‌کند. به عبارت دیگر، به ازای برخی مقادیر  $\gamma$  داریم

$$P(M \subseteq \hat{M}) \rightarrow 1, \quad n \rightarrow \infty. \quad (۴۱.۲)$$

مطلب مهم در این خصوص، شناسایی شرایطی است که تحت آن‌ها (۴۱.۲) برقرار باشد. شرایط لازم برای برقراری این ویژگی در قضیه مذکور مطرح شده است. این روش غربالگری، بعد متغیرهای توضیحی را از مقیاس بسیار بزرگ  $p$  به مقیاس نسبتاً متوسط  $d = [\gamma n] < p$  کاهش می‌دهد. بعد از مرحله غربالگری، روش‌های انتخاب متغیر مبتنی بر جریمه بر روی مدل  $\hat{M}$  برای انتخاب مدل نهایی به کار می‌روند. شکل ۷.۲ می‌تواند در فهم مطلب موثر باشد.

اگر پس از غربالگری، در مرحله دوم از جریمه SCAD استفاده شود، این روش دو مرحله‌ای را با SIS-SCAD نشان می‌دهیم. علاوه بر این، می‌توان پس از به‌کارگیری SIS، دوبار به‌طور متوالی از روش‌های جریمه استفاده کرد. به عنوان مثال، ابتدا با استفاده از SIS بعد مدل را از  $p$  به  $d < n$  کاهش داده و سپس از انتخاب‌گر دن‌زیگ برای انتخاب متغیر روی مدل غربال شده و کاهش بعد آن به  $d' < d$  استفاده شود. در پایان، برای انتخاب مدل نهایی، روش لاسوی تطبیقی یا SCAD روی مدل کاهش یافته به کار می‌رود. در این حالت، این روش سه مرحله‌ای را با SIS-DS-SCAD یا SIS-DS-AdaLasso نشان می‌دهند.

پس از معرفی روش SIS توسط فن و لیو (۲۰۰۸)، این روش به مدل‌های آماری مختلف مانند مدل

<sup>۵۵</sup>Marginal utility measure

<sup>۵۶</sup>Sure screening property

خطی تعمیم‌یافته (فن و همکاران، ۲۰۰۹؛ فن و سانگ، ۲۰۱۰)، مدل خطرات متناسب کاکس<sup>۵۷</sup> (فن و همکاران، ۲۰۱۰)، مدل جمعی ناپارامتری (فن و همکاران، ۲۰۱۱) و مدل ضریب متغیر (فن و همکاران، ۲۰۱۴؛ لیو<sup>۵۸</sup> و همکاران، ۲۰۱۴؛ سانگ<sup>۵۹</sup> و همکاران، ۲۰۱۴؛ زیا و همکاران، ۲۰۱۶) تعمیم داده شد. روش SIS را می‌توان برای کاهش بعد فضای متغیرها در مسائل رده‌بندی هم به کار برد. بدین منظور، فرض کنید  $n_1$  و  $n_0$  تعداد مشاهداتی از  $X_k$  باشند که رده آن‌ها به ترتیب  $\circ$  و  $1$  است. همچنین فرض کنید  $\bar{X}_{k,1}$  و  $\bar{X}_{k,0}$  به ترتیب میانگین متغیر توضیحی  $X_k$  در رده های  $\circ$  و  $1$ ، و  $Se(X_k)$  انحراف معیار متغیر توضیحی  $X_k$  باشد. تحت این شرایط  $k$ -امین مولفه بردار  $\omega$  به صورت

$$\omega_k = \frac{n_1 \bar{X}_{k,1} - n_0 \bar{X}_{k,0}}{Se(X_k)}$$

محاسبه می‌شود. بقیه مراحل انتخاب متغیر مشابه قبل انجام می‌شود.

### • غربالگری مستقل مطمئن تکراری

ایده اصلی روش SIS رده‌بندی متغیرها بر اساس یک معیار مطلوبیت حاشیه‌ای است. بنابراین، سه مسئله مهم زیر ممکن است عملکرد این روش را با مشکل مواجه کنند:

۱. متغیر مهمی که به صورت حاشیه‌ای با متغیر پاسخ ناهمبسته است، توسط این روش انتخاب نمی‌شود.
۲. برخی از متغیرهای بی‌اهمیت که همبستگی بالایی با متغیرهای مهم دارند، نسبت به متغیرهای مهمی که دارای همبستگی نسبتاً ضعیفی با متغیر پاسخ هستند، در انتخاب متغیر توسط SIS، برتری دارند.
۳. این روش، همخطی بین متغیرهای توضیحی را در نظر نمی‌گیرد.

برای درک بهتر مشکلات روش SIS، دو مثال زیر را در نظر بگیرید.

#### مثال ۱.۶.۲. مدل خطی زیر را در نظر بگیرید

$$Y = \mathbf{X}^T \boldsymbol{\beta}^* + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1) \quad (42.2)$$

که  $p = 2000$  و  $\boldsymbol{\beta}^* = (3, 3, 3, 3, 3, -7/5, 0, \dots, 0)^T$  بنابراین، مجموعه اندیس متغیرهای مهم به صورت  $M = \{j : \beta_j^* \neq 0\} = \{1, 2, 3, 4, 5, 6\}$  است.

فرض کنید متغیرهای توضیحی دارای توزیع نرمال استاندارد و ضریب همبستگی بین همه متغیرهای توضیحی برابر  $0.5$  باشد. بنابراین  $Cov(X_6, Y) = 0$ . لذا،  $X_6$  و  $Y$  به‌طور حاشیه‌ای ناهمبسته هستند. در نتیجه متغیر مهم  $X_6$  توسط SIS انتخاب نمی‌شود.

<sup>57</sup>Cox's proportional hazards model

<sup>58</sup>Liu

<sup>59</sup>Song

**مثال ۲.۶.۲.** دوباره مدل خطی (۴۲.۲) را با  $\beta^* = (1, 0, \dots, 0, 1)^T$  در نظر بگیرید. در این مدل، مجموعه اندیس متغیرهای مهم شامل  $M = \{1, 2000\}$  است. فرض کنید ضریب همبستگی بین تمام متغیرهای توضیحی، به جز  $X_{2000}$ ، برابر  $0.9$  و  $X_{2000}$  مستقل از بقیه متغیرها است. در نتیجه

$$\text{Cov}(X_1, Y) = 1, \quad \text{Cov}(X_{2000}, Y) = 1, \quad \text{Cov}(X_j, Y) = 0, \quad j \neq 1, 2000.$$

در این حالت، ضریب همبستگی حاشیه‌ای همه متغیرهای بی‌اهمیت بزرگتر از ضریب همبستگی متغیر مهم  $X_{2000}$  با پاسخ است. بنابراین تمام متغیرهای بی‌اهمیت در فرآیند انتخاب توسط SIS نسبت به  $X_{2000}$  برتری دارند.

اگر ماتریس طرح داده‌ها متعامد باشد، روش SIS یک روش انتخاب متغیر قابل اعتماد بوده و نتایج آن مطلوب خواهد بود. اما در داده‌های با بعد بسیار بالا، موضوع همخطی امری اجتناب‌ناپذیر است. بنابراین با اعمال روش SIS، ممکن است سه مشکل فوق، عملکرد این روش را تحت تاثیر قرار دهند و برخی از متغیرهای مهم را از دست بدهیم. برای رفع این مشکل، فن و لیو (۲۰۰۸) روش SIS را به‌طور مکرر به کار بردند و این روش را غربالگری مستقل مطمئن تکراری (ISIS) نامیدند. الگوریتم روش ISIS به صورت زیر است:

- مرحله اول: یک زیرمجموعه  $k_1$  تایی از متغیرها، مانند  $A_1 = \{X_{i_1}, \dots, X_{i_{k_1}}\}$ ، را با استفاده از یک روش غربالگری مبتنی بر SIS، مانند SIS-SCAD یا SIS-Lasso انتخاب می‌کنیم. سپس یک مدل خطی با این مجموعه از متغیرهای توضیحی بر روی متغیر پاسخ برازش می‌دهیم.
- مرحله دوم: باقیمانده‌های مدل برازش شده در مرحله قبل را به‌عنوان متغیر پاسخ جدید در نظر گرفته و مرحله اول را برای  $p - k_1$  متغیر باقیمانده اجرا می‌کنیم تا زیرمجموعه جدید با  $k_2$  متغیر توضیحی و جدا از  $A_1$  به صورت  $A_2 = \{X_{j_1}, \dots, X_{j_{k_2}}\}$  حاصل شود. این کار تا زمانی ادامه می‌یابد که  $l$  مجموعه جدا از هم  $A_1, \dots, A_l$  حاصل شوند تا مجموعه  $A = \bigcup_{i=1}^l A_i$  با اندازه  $d$  را تشکیل دهند. در عمل  $l$  را می‌توان طوری انتخاب کرد که  $|A| < n$  باشد. در نهایت مجموعه  $A$  شامل همه متغیرهای مهم است که با روش ISIS انتخاب می‌شوند. سپس می‌توان روش‌های انتخاب متغیر مبتنی بر جریمه را برای متغیرهای مجموعه  $A$  به کار برد.

باید توجه داشته باشیم که برازش باقیمانده‌های مرحله اول بر روی متغیرهای توضیحی باقیمانده، می‌تواند به‌طور قابل توجهی اولویت انتخاب آن دسته از متغیرهای بی‌اهمیتی که همبستگی بالایی با متغیر پاسخ از طریق متغیرهای توضیحی  $\{X_{i_1}, \dots, X_{i_{k_1}}\}$  دارند را تضعیف کند، زیرا باقیمانده‌ها با این متغیرهای توضیحی ناهمبسته‌اند و این امر موجب رفع مشکل همخطی می‌شود. این عمل، همچنین باعث می‌شود متغیرهای توضیحی مهم از دست رفته در مرحله اول، احتمالاً، حفظ شوند. روش‌های غربالگری معرفی شده فوق مدل مبنا بوده و در صورت نادرست بودن مدل آماری مفروض، ممکن است منجر به حذف متغیرهای مهم شوند. در مدل‌سازی با بعد بالا، بدون داشتن اطلاعات پیشین، تعیین ساختار مدل واقعی بسیار مشکل است. بنابراین، برای اجتناب از تشخیص نادرست ساختار مدل، معرفی یک روش غربالگری آزاد-مدل<sup>۶۰</sup>، که به نوع مدل بستگی نداشته باشد، ضروری به

<sup>60</sup>Model free



نظر می‌رسد. در این راستا، ژو<sup>۶۱</sup> و همکاران (۲۰۱۱) یک روش غربالگری آزاد-مدل به نام ”غربالگری و رتبه بندی مستقل مطمئن“ (SIRS) را معرفی کردند که در ادامه به بیان جزئیات این روش می‌پردازیم.

## ۲.۶.۲ غربالگری و رتبه بندی مستقل مطمئن

فرض کنید  $Y$  متغیر پاسخ با تکیه‌گاه  $\Psi_y$  و  $\mathbf{X} = (X_1, \dots, X_p)^T$  بردار متغیرهای توضیحی باشد. تابع توزیع شرطی  $Y$  به شرط  $\mathbf{X}$  را به صورت  $F(y|\mathbf{X}) = P(Y < y|\mathbf{X})$  در نظر گرفته و فرض کنید، اندیس متغیرهای مهم به صورت زیر تعریف شود

$$\mathbb{M} = \{k : \text{تابع } F(y|\mathbf{X}) \text{ به طور تابعی وابسته به } X_k \text{ است} : k\}.$$

اگر  $k \in \mathbb{M}$ ، آنگاه  $X_k$  یک متغیر مهم و در غیر این صورت  $X_k$  یک متغیر بی‌اهمیت است. فرض کنید به‌ازای  $1 \leq k \leq p$ ،  $E(X_k) = 0$  و  $Var(X_k) = 1$ . حال  $\Omega(y) = E\{\mathbf{X}F(y|\mathbf{X})\}$  را در نظر بگیرید. با استفاده از رابطه  $F(y|\mathbf{X}) = E\{I(Y < y)|\mathbf{X}\}$  و قانون امید ریاضی مکرر داریم

$$\begin{aligned} \Omega(y) &= E[\mathbf{X}E\{I(Y < y)|\mathbf{X}\}] \\ &= E[E\{\mathbf{X}I(Y < y)|\mathbf{X}\}] \\ &= E[\mathbf{X}I(Y < y)] \\ &= Cov(\mathbf{X}, I(Y < y)) \end{aligned} \quad (۴۳.۲)$$

که  $I(\cdot)$  تابع نشانگر است. فرض کنید  $\Omega_k(y)$ ،  $k$  - امین مولفه از  $\Omega(y)$  باشد و

$$\omega_k = E\{\Omega_k^2(Y)\}, \quad k = 1, \dots, p \quad (۴۴.۲)$$

را تعریف کنید. ژو و همکاران (۲۰۱۱) از  $\omega_k$  به‌عنوان معیار مطلوبیت حاشیه‌ای برای رتبه‌بندی متغیرها استفاده کردند، یعنی برای اندازه‌گیری اهمیت متغیر  $X_k$  بر پاسخ، همبستگی بین  $X_k$  و  $I(Y < y)$  را به‌کار بردند. به‌سادگی می‌توان دید که، اگر  $X_k$  و  $Y$  مستقل باشند، آنگاه  $X_k$  و  $I(Y < y)$  نیز مستقل‌اند. بنابراین، برای هر  $y \in \Psi_y$ ،  $\Omega_k(y) = 0$  و در نتیجه  $\omega_k = 0$ . از طرف دیگر، اگر  $X_k$  و  $Y$  مرتبط باشند، آنگاه برای برخی مقادیر  $y \in \Psi_y$ ،  $\Omega_k(y) \neq 0$ ، و بنابراین  $\omega_k$  مقداری مثبت است. بنابراین (۴۴.۲) می‌تواند معیار مناسبی برای رتبه‌بندی متغیرها باشد.

برای پیاده‌سازی این روش غربالگری، لازم است معیار (۴۴.۲) بر اساس یک نمونه تصادفی برآورد شود. فرض کنید  $(\mathbf{X}_i, Y_i)_{i=1}^n$  یک نمونه تصادفی از  $(\mathbf{X}, Y)$  باشد. برای سهولت، فرض می‌کنیم متغیرهای توضیحی استاندارد شده‌اند، به‌طوری‌که،  $\frac{1}{n} \sum_{i=1}^n X_{ij} = 0$  و  $\frac{1}{n} \sum_{i=1}^n X_{ij}^2 = 1$ . برای هر  $y$  داده شده، برآوردگر گشتاوری نمونه‌ای  $\Omega_k(y)$  به صورت

$$\hat{\Omega}_k(y) = \frac{1}{n} \sum_{i=1}^n X_{ik} I(Y_i < y), \quad k = 1, \dots, p \quad (۴۵.۲)$$

محاسبه می‌شود که  $X_{ik}$ ،  $k$  - امین مولفه از  $\mathbf{X}_i$  است. در نتیجه  $\omega_k$  به صورت

$$\hat{\omega}_k = \frac{1}{n} \sum_{j=1}^n \left\{ \frac{1}{n} \sum_{i=1}^n X_{ik} I(Y_i < Y_j) \right\}^2, \quad k = 1, \dots, p \quad (46.2)$$

برآورد می‌شود. سپس با رتبه‌بندی متغیرهای توضیحی بر اساس  $\hat{\omega}_k$ ، متغیرهای توضیحی با مقدار  $\hat{\omega}_k$  بزرگتر از یک مقدار آستانه را به‌عنوان مجموعه متغیرهای مهم انتخاب می‌کنیم. برای اطلاعات بیشتر درباره انتخاب این مقدار آستانه به ژو و همکاران (۲۰۱۱) مراجعه کنید.  
ژو و همکاران (۲۰۱۱) تحت قضیه‌ای نشان دادند

$$\max_{k \in \mathbb{M}^c} \omega_k < \min_{k \in \mathbb{M}} \omega_k$$

و  $\omega_k = 0$  اگر و فقط اگر  $k \in \mathbb{M}^c$ . علاوه بر این، آن‌ها ویژگی سازگاری در رتبه‌بندی<sup>۶۲</sup> را برای روش ارائه شده ثابت کردند. این ویژگی بدین معنی است که، در احتمال، معیار  $\omega_k$  همواره یک متغیر مهم را در رتبه بالاتری از یک متغیر بی اهمیت قرار می‌دهد. این ویژگی، سازگاری در انتخاب را نتیجه می‌دهد. قضایای فوق نشان می‌دهند که  $\omega_k$  معیار حاشیه‌ای مناسبی برای غربالگری متغیرها در داده‌های با بعد بسیار بالاست.

### ۳.۶.۲ غربالگری با همبستگی فاصله‌ای

روش دیگر برای رسیدن به یک غربالگری آزاد-مدل، استفاده از یک معیار وابستگی است که به‌طور کارا بتواند روابط خطی و غیرخطی بین متغیرهای توضیحی و متغیر پاسخ را شناسایی کند. در این راستا، لی و همکاران (۲۰۱۲) یک روش غربالگری آزاد-مدل بر اساس همبستگی فاصله‌ای ارائه دادند. برخلاف ضریب همبستگی پیرسن که برای دو متغیر تصادفی تعریف می‌شود، ضریب همبستگی فاصله‌ای برای دو بردار تصادفی تعریف می‌شود که ممکن است طول این دو بردار متفاوت باشد.

**تعریف ۳.۶.۲.** کوواریانس فاصله‌ای بین بردارهای تصادفی  $\mathbf{u} \in \mathbb{R}^q$  و  $\mathbf{v} \in \mathbb{R}^r$ ، با گشتاورهای مرتبه اول متناهی، یک عدد نامنفی به صورت

$$dcov(\mathbf{u}, \mathbf{v}) = \left[ \int_{\mathbb{R}^{q+r}} \|\phi_{\mathbf{u}, \mathbf{v}}(\mathbf{t}, \mathbf{s}) - \phi_{\mathbf{u}}(\mathbf{t})\phi_{\mathbf{v}}(\mathbf{s})\|^2 \omega(\mathbf{t}, \mathbf{s}) dt ds \right]^{\frac{1}{2}} \quad (47.2)$$

است که  $\phi_{\mathbf{u}}(\mathbf{t})$  و  $\phi_{\mathbf{v}}(\mathbf{s})$ ، به ترتیب، توابع مشخصه حاشیه‌ای  $\mathbf{u}$  و  $\mathbf{v}$ ،  $\phi_{\mathbf{u}, \mathbf{v}}(\mathbf{t}, \mathbf{s})$  تابع مشخصه توام  $(\mathbf{u}, \mathbf{v})$  و  $\omega(\mathbf{t}, \mathbf{s}) = \{c_q c_r \|\mathbf{t}\|_q^{q+1} \|\mathbf{s}\|_r^{r+1}\}^{-1}$  که  $c_d = \pi^{(1+d)/2} / \Gamma\{\pi^{(1+d)/2}\}$  همچنین،  $\|\mathbf{a}\|_d$  نرم اقلیدسی بردار  $\mathbf{a} \in \mathbb{R}^d$  و  $\|\phi\|^2 = \phi \bar{\phi}$  است که  $\phi$  یک تابع مختلط و  $\bar{\phi}$  مزدوج  $\phi$  می‌باشد.

**تعریف ۴.۶.۲.** همبستگی فاصله‌ای بین دو بردار تصادفی  $\mathbf{u} \in \mathbb{R}^q$  و  $\mathbf{v} \in \mathbb{R}^r$ ، با گشتاورهای مرتبه اول متناهی به صورت

$$dcorr(\mathbf{u}, \mathbf{v}) = \begin{cases} \frac{dcov(\mathbf{u}, \mathbf{v})}{\sqrt{dcov(\mathbf{u}, \mathbf{u})dcov(\mathbf{v}, \mathbf{v})}}, & dcov(\mathbf{u}, \mathbf{u})dcov(\mathbf{v}, \mathbf{v}) > 0 \\ 0, & dcov(\mathbf{u}, \mathbf{u})dcov(\mathbf{v}, \mathbf{v}) = 0. \end{cases}$$

<sup>62</sup>Consistency in ranking

تعریف می‌شود و در شرط  $1 \leq dcorr(\mathbf{u}, \mathbf{v}) \leq 0$  صدق می‌کند. در حالت خاص، اگر  $X$  و  $Y$  دو متغیر تصادفی نرمال یک متغیره با ضریب همبستگی پیرسن  $\rho$  و ضریب همبستگی فاصله‌ای  $\mathcal{R}$  باشند، آنگاه  $|\rho| = \mathcal{R}$ . تساوی هنگامی برقرار است که  $|\rho| = 1$ .

سزکلی<sup>۶۳</sup> و همکاران (۲۰۰۷) نشان دادند

$$dcov^2(\mathbf{u}, \mathbf{v}) = S_1 + S_2 - 2S_3 \quad (48.2)$$

به طوری که  $S_1 = E\{\|\mathbf{u} - \tilde{\mathbf{u}}\|_q \|\mathbf{v} - \tilde{\mathbf{v}}\|_r\}$ ،  $S_2 = E\{\|\mathbf{u} - \tilde{\mathbf{u}}\|_q\} E\{\|\mathbf{v} - \tilde{\mathbf{v}}\|_r\}$ ،  $S_3 = E\{E(\|\mathbf{u} - \tilde{\mathbf{u}}\|_q | \mathbf{u}) E(\|\mathbf{v} - \tilde{\mathbf{v}}\|_r | \mathbf{v})\}$  و  $(\tilde{\mathbf{u}}, \tilde{\mathbf{v}})$  بردارهای تصادفی مستقل از  $(\mathbf{u}, \mathbf{v})$  و هم توزیع با آن‌ها هستند. بنابراین کوواریانس فاصله‌ای بین  $\mathbf{u}$  و  $\mathbf{v}$  را می‌توان به صورت

$$\widehat{dcov}^2(\mathbf{u}, \mathbf{v}) = \hat{S}_1 + \hat{S}_2 - 2\hat{S}_3$$

برآورد کرد که

$$\begin{aligned} \hat{S}_1 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{u}_i - \mathbf{u}_j\|_q \|\mathbf{v}_i - \mathbf{v}_j\|_r \\ \hat{S}_2 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{u}_i - \mathbf{u}_j\|_q \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{v}_i - \mathbf{v}_j\|_r \\ \hat{S}_3 &= \frac{1}{n^4} \sum_{i=1}^n \sum_{j=1}^n \sum_{l=1}^n \|\mathbf{u}_i - \mathbf{u}_l\|_q \|\mathbf{v}_j - \mathbf{v}_l\|_r. \end{aligned}$$

همچنین کوواریانس‌های فاصله‌ای نمونه‌ای  $dcov(\mathbf{u}, \mathbf{u})$  و  $dcov(\mathbf{v}, \mathbf{v})$  را نیز می‌توان به طور مشابه برآورد کرد. در نتیجه همبستگی فاصله‌ای نمونه‌ای بین  $\mathbf{u}$  و  $\mathbf{v}$  برابر است با

$$\widehat{dcorr}(\mathbf{u}, \mathbf{v}) = \frac{\widehat{dcov}(\mathbf{u}, \mathbf{v})}{\sqrt{\widehat{dcov}(\mathbf{u}, \mathbf{u})} \sqrt{\widehat{dcov}(\mathbf{v}, \mathbf{v})}}. \quad (49.2)$$

ضریب همبستگی فاصله‌ای دارای دو ویژگی مطلوب زیر است. ویژگی اول، رابطه بین ضریب همبستگی فاصله‌ای و ضریب همبستگی پیرسن است. ضریب همبستگی فاصله‌ای دو متغیر تصادفی نرمال یک متغیره تابعی اکیداً صعودی از قدر مطلق ضریب همبستگی پیرسن آن دو متغیر نرمال است (سزکلی و همکاران، ۲۰۰۷). ویژگی دوم بیان می‌کند که بردارهای تصادفی  $\mathbf{u}$  و  $\mathbf{v}$  مستقل‌اند اگر و تنها اگر  $dcorr(\mathbf{u}, \mathbf{v}) = 0$  (سزکلی و همکاران، ۲۰۰۷؛ سزکلی و ریزو<sup>۶۴</sup>، ۲۰۰۹).

به دلیل دو ویژگی فوق، لی و همکاران (۲۰۱۲) از همبستگی فاصله‌ای برای رتبه‌بندی متغیرهای توضیحی استفاده کردند و یک روش غربالگری مستقل مطمئن موسوم به DC-SIS را ارائه دادند. آن‌ها از ضریب همبستگی فاصله‌ای بین متغیرهای توضیحی و متغیر پاسخ برای رتبه‌بندی متغیرها استفاده کردند و معیار مطلوبیت حاشیه‌ای را به صورت زیر تعریف کردند

$$\omega_k = dcorr^2(X_k, Y), \quad k = 1, \dots, p. \quad (50.2)$$

<sup>63</sup>Székely

<sup>64</sup>Rizzo

ویژگی اول ضریب همبستگی فاصله‌ای نتیجه می‌دهد که در صورت نرمال بودن توزیع خطا و توزیع متغیرهای توضیحی، DC-SIS با معادل با SIS است، یعنی در مدل خطی گاوسی در نظر گرفته‌شده توسط فن و لیو (۲۰۰۸)، DC-SIS به خوبی SIS عمل می‌کند. ویژگی دوم، آزاد-مدل بودن روش غربالگری DC-SIS را نتیجه می‌دهد، یعنی در این روش نیازی به در نظر گرفتن یک مدل مفروض برای متغیرهای توضیحی و پاسخ نیست. ویژگی آزاد-مدل بودن بسیار حائز اهمیت است، چون تعیین یک مدل پارامتری معین برای داده‌های با بعد بسیار بالا یک چالش است.

برای پیاده‌سازی این روش غربالگری، کافی است معیار مطلوبیت حاشیه‌ای (۵۰.۲) بر اساس یک نمونه تصادفی برآورد شود. فرض کنید  $(X_i, Y_i)_{i=1}^n$  یک نمونه تصادفی از  $(X, Y)$  باشد. بنابراین معیار (۵۰.۲) را به صورت

$$\widehat{\omega}_k = \widehat{dcorr}^2(X_k, Y)$$

برآورد می‌کنیم. سپس، مجموعه‌ای از متغیرهای توضیحی که دارای مقادیر بزرگ  $\widehat{\omega}_k$  هستند، به‌عنوان مجموعه متغیرهای مهم انتخاب می‌شوند.

لی و همکاران (۲۰۱۲) نشان دادند که روش DC-SIS دارای ویژگی غربالگری مطمئن است. این روش دارای چندین مزیت است: الف) برای پیاده‌سازی این روش نیازی به استفاده از هیچ الگوریتم بهینه‌سازی نیست. ب) این روش را می‌توان مستقیماً برای پاسخ چندگانه یا متغیرهای توضیحی با ساختار گروهی به کار برد. ج) این روش برای هر نوع متغیر پاسخ پیوسته، گسسته یا شمارشی قابل استفاده است. بنابراین DC-SIS یک روش آزاد-مدل مناسب برای داده‌های با بعد بسیار بالا است. در سال‌های اخیر مطالعات زیادی در زمینه غربالگری متغیرها انجام شده است. از جمله این کارها می‌توان هال و میلر<sup>۶۵</sup> (۲۰۰۹)، بولمن و همکاران (۲۰۱۰)، بارت<sup>۶۶</sup> و همکاران (۲۰۱۶)، وانگ<sup>۶۷</sup> و لنگ (۲۰۱۶)، کنگ<sup>۶۸</sup> و همکاران (۲۰۱۷)، چن<sup>۶۹</sup> و همکاران (۲۰۱۸)، هوانگ و ژو (۲۰۱۶)، لیو<sup>۷۰</sup> (۲۰۱۶)، لی و همکاران (۲۰۱۸)، چن (۲۰۱۸) و سانگ و همکاران (۲۰۱۸) را نام برد.

## ۷.۲ الگوریتم محاسباتی روش کمترین توان‌های دوم جریمه‌شده

تابع کمترین توان‌های دوم جریمه‌شده زیر را در نظر بگیرید

$$\mathcal{L}(\beta) = \frac{1}{p} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \sum_{j=1}^p p_\lambda(|\beta_j|). \quad (51.2)$$

با یک تابع جریمه محذب، مانند لاسو، تابع هدف (۵۱.۲) یک تابع محذب است و لذا برآوردهای کمترین توان‌های دوم جریمه‌شده را می‌توان با بهینه‌سازی محذب به‌دست آورد. اما اغلب توابع جریمه‌ای که

<sup>65</sup>Hall and Miller

<sup>66</sup>Barut

<sup>67</sup>Wang

<sup>68</sup>Kong

<sup>69</sup>Chen

<sup>70</sup>Liu

دارای سه ویژگی مطلوب تنکی، ناریبی و پیوستگی می‌باشند، تابعی غیر محدب بوده و این موضوع باعث می‌شود که تابع هدف (۵۱.۲) محدب نباشد. محدب نبودن این تابع، مسئله مینیمم‌سازی را با مشکل مواجه می‌کند. به هر حال، می‌توان توابع جریمه غیر محدب را با برخی توابع محدب تقریب زد و مسئله بهینه‌سازی غیر محدب را با استفاده از الگوریتم‌های بهینه‌سازی محدب حل کرد. در این بخش دو تقریب زیر را برای حل مشکل فوق معرفی می‌کنیم.

### ۱.۷.۲ تقریب درجه دو موضعی

چون عبارت اول در تابع کمترین توان‌های دوم جریمه‌شده (۵۱.۲) یک تابع محدب است، لذا کافی است تابع جریمه غیر محدب  $P_\lambda(|\beta_j|)$  را با یک تابع محدب تقریب کنیم. در این راستا، فن و لی (۲۰۰۱) الگوریتم تقریب درجه دو موضعی (LQA) را برای توابع جریمه غیر محدب ارائه دادند. فرض کنید  $\beta^{(0)} = (\beta_1^{(0)}, \dots, \beta_p^{(0)})^T$  یک مقدار اولیه برای  $\beta$  باشد. اگر  $\beta_j^{(0)}$  نزدیک به صفر باشد، پس قرار دهید  $\hat{\beta}_j = 0$ ، در غیر اینصورت با استفاده از تقریب موضعی داریم

$$[p_\lambda(|\beta_j|)]' = p'_\lambda(|\beta_j|) \frac{\beta_j}{|\beta_j|} \approx p'_\lambda(|\beta_j^{(0)}|) \frac{\beta_j}{|\beta_j^{(0)}|}, \quad \beta_j \approx \beta_j^{(0)}. \quad (52.2)$$

با انتگرال گیری از طرفین نسبت به  $\beta_j$  از  $\beta_j^{(0)}$  تا  $\beta_j$  داریم

$$p_\lambda(|\beta_j|) \approx p_\lambda(|\beta_j^{(0)}|) + \frac{1}{2} \frac{p'_\lambda(|\beta_j^{(0)}|)}{|\beta_j^{(0)}|} \{\beta_j^2 - \beta_j^{(0)2}\}. \quad (53.2)$$

بنابراین

$$\begin{aligned} \mathcal{L}(\beta) &= \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p p_\lambda(|\beta_j|) \\ &\approx \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p \left\{ p_\lambda(|\beta_j^{(0)}|) + \frac{1}{2} \frac{p'_\lambda(|\beta_j^{(0)}|)}{|\beta_j^{(0)}|} (\beta_j^2 - \beta_j^{(0)2}) \right\} \\ &= \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \sum_{j=1}^p \frac{1}{2} \frac{p'_\lambda(|\beta_j^{(0)}|)}{|\beta_j^{(0)}|} \beta_j^2 + c \end{aligned} \quad (54.2)$$

که ثابت  $c$  تابعی از  $\beta^{(0)}$  است. با استفاده از تقریب فوق و حذف جملات ثابت، تابع هدف (۵۱.۲) به یک تابع درجه دو تبدیل می‌شود. در نتیجه، برآورد کمترین توان‌های دوم جریمه‌شده را می‌توان با محاسبه تکراری رگرسیون ریبج به صورت

$$\beta^{(1)} = \{\mathbf{X}^T \mathbf{X} + \Sigma_\lambda(\beta^{(0)})\}^{-1} \mathbf{X}^T \mathbf{Y} \quad (55.2)$$

به‌دست آورد که

$$\Sigma_\lambda(\beta^{(0)}) = \text{diag} \left\{ \frac{p'_\lambda(|\beta_1^{(0)}|)}{|\beta_1^{(0)}|}, \dots, \frac{p'_\lambda(|\beta_p^{(0)}|)}{|\beta_p^{(0)}|} \right\}.$$

با توجه به این واقعیت که رگرسیون ریج نمی‌تواند متغیرهای معنی‌دار را انتخاب کند، ضرایب رگرسیونی برآورد شده را که بسیار نزدیک به صفر می‌باشند، برابر صفر قرار می‌دهیم تا متغیرهای بی‌اهمیت از مدل حذف شوند.

یکی از مزایای تقریب LQA تعمیم‌پذیری آن به هر تابع زیان هموار دیگری (غیر از تابع زیان کمترین توان‌های دوم) است. با در نظر گرفتن تابع زیان  $l(\beta)$  تابع هدف جریمه‌شده به صورت

$$l(\beta) + \sum_{j=1}^p P_{\lambda}(|\beta_j|) \quad (56.2)$$

تعریف می‌شود. در این حالت، ابتدا تابع زیان  $l(\beta)$  را با یک تابع درجه دو به صورت

$$l(\beta) \approx l(\beta^{(0)}) + \nabla l(\beta^{(0)})^T (\beta - \beta^{(0)}) + \frac{1}{2} (\beta - \beta^{(0)})^T \nabla^2 l(\beta^{(0)}) (\beta - \beta^{(0)}) \quad (57.2)$$

تقریب می‌کنیم که

$$\nabla l(\beta^{(0)}) = \frac{\partial l(\beta^{(0)})}{\partial \beta}, \quad \nabla^2 l(\beta^{(0)}) = \frac{\partial^2 l(\beta^{(0)})}{\partial \beta \partial \beta^T}.$$

با حذف جمله ثابت، تابع هدف (51.2) را می‌توان به صورت

$$\nabla l(\beta^{(0)})^T (\beta - \beta^{(0)}) + \frac{1}{2} (\beta - \beta^{(0)})^T \nabla^2 l(\beta^{(0)}) (\beta - \beta^{(0)}) + \frac{1}{2} \beta^T \Sigma_{\lambda}(\beta^{(0)}) \beta \quad (58.2)$$

نوشت و سپس با استفاده از الگوریتم نیوتن-رافسون مینیم تابع (51.2) را به دست آورد. به‌طور دقیق‌تر، با اجرای مراحل زیر، الگوریتم LQA را می‌توان برای تقریب تابع (56.2) به کار برد:

- مرحله اول: مقدار اولیه  $\beta^{(0)}$  را برای  $\beta$  انتخاب کنید. این مقدار اولیه می‌تواند برآورد OLS باشد.

- مرحله دوم: تابع هدف (56.2) را بر اساس (58.2) تقریب کنید.

- مرحله سوم: الگوریتم نیوتن-رافسون را برای محاسبه مقدار مینیمم‌کننده تابع درجه دو به دست آمده در مرحله دوم به کار برید تا مقدار  $\beta_j^{(0)}$  به‌روزرسانی شود. سپس اگر  $|\beta_j^{(0)}| < \epsilon$  باشد، متغیر  $X_j$  را از مدل حذف کنید.

- مرحله چهارم: مراحل دوم و سوم را تا رسیدن به همگرایی تکرار کنید.

فن و لی (۲۰۰۱) نشان دادند که در صورت انتخاب یک مقدار مناسب برای  $\beta^{(0)}$ ، الگوریتم LQA پس از تعداد کمی تکرار به همگرایی می‌رسد. علاوه بر این، با یک مقدار اولیه مناسب، کارایی الگوریتم LQA با یک تکرار معادل با تکرار کامل الگوریتم تا رسیدن به همگرایی است. اما ایراد وارد بر این الگوریتم این است که اگر یک متغیر در هر تکرار حذف شود، در تکرار بعدی وارد مدل نمی‌شود و در نتیجه از مدل نهایی حذف می‌شود. بنابراین، این روش مشابه روش انتخاب پیشرو عمل می‌کند.

## ۲.۷.۲ تقریب خطی موضعی

برای غلبه بر مشکل بیان شده در تقریب LQA، زو<sup>۷۱</sup> و لی (۲۰۰۸) الگوریتم تقریب خطی موضعی (LLA) را برای تقریب توابع جریمه غیر محدب پیشنهاد دادند. با استفاده از این تقریب برای تابع جریمه

$$p_{\lambda}(|\beta_j|) \approx p_{\lambda}(|\beta_j^{(0)}|) + p'_{\lambda}(|\beta_j^{(0)}|) \left\{ |\beta_j| - |\beta_j^{(0)}| \right\}, \quad \beta_j \approx \beta_j^{(0)} \quad (۵۹.۲)$$

می‌توان تابع هدف (۵۱.۲) را به صورت

$$\begin{aligned} \mathcal{L}(\beta) &= \frac{1}{\gamma} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \sum_{j=1}^p p_{\lambda}(|\beta_j|) \\ &\approx \frac{1}{\gamma} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \sum_{j=1}^p \left\{ p_{\lambda}(|\beta_j^{(0)}|) + p'_{\lambda}(|\beta_j^{(0)}|) (|\beta_j| - |\beta_j^{(0)}|) \right\} \\ &= \frac{1}{\gamma} \| \mathbf{Y} - \mathbf{X}\beta \|^2 + \sum_{j=1}^p p'_{\lambda}(|\beta_j^{(0)}|) |\beta_j| + c \end{aligned} \quad (۶۰.۲)$$

تقریب کرد که ثابت  $c$  تابعی از  $\beta^{(0)}$  است. بنابراین تقریب LLA تابع هدف جریمه‌شده غیر محدب را به تابع هدف لاسوی وزنی تبدیل می‌کند. این تابع هدف تبدیل‌یافته را می‌توان با استفاده از برنامه‌ریزی درجه دو یا الگوریتم کمترین زاویه حل کرد.

شکل ۸.۲ تقریب LQA و LLA را برای توابع جریمه  $L_{0.5}$  و SCAD نشان می‌دهد. در این شکل برای هر دو تابع جریمه  $\lambda = 2$  و دو مقدار اولیه مختلف در نظر گرفته شده‌است.

تاکنون، علاوه بر دو تقریب فوق، الگوریتم‌های دیگری نیز برای حل مسئله کمترین توان‌های دوم جریمه‌شده معرفی شده‌اند. الگوریتم پیش‌برنده<sup>۷۲</sup> (فو<sup>۷۳</sup>، ۱۹۹۸)، الگوریتم مختصات نزولی<sup>۷۴</sup> برای لاسوی گروهی (وو و لانگ<sup>۷۵</sup>، ۲۰۰۸)، الگوریتم ماکزیمم‌سازی شرطی تکراری<sup>۷۶</sup> (ژانگ و لی، ۲۰۰۹)، الگوریتم مختصات نزولی برای توابع جریمه غیرمحدب (برهنی و هوانگ، ۲۰۱۱) و الگوریتم مختصات نزولی گروهی (برهنی و هوانگ، ۲۰۱۵) نمونه‌هایی از این الگوریتم‌ها می‌باشند.

<sup>71</sup>Zou

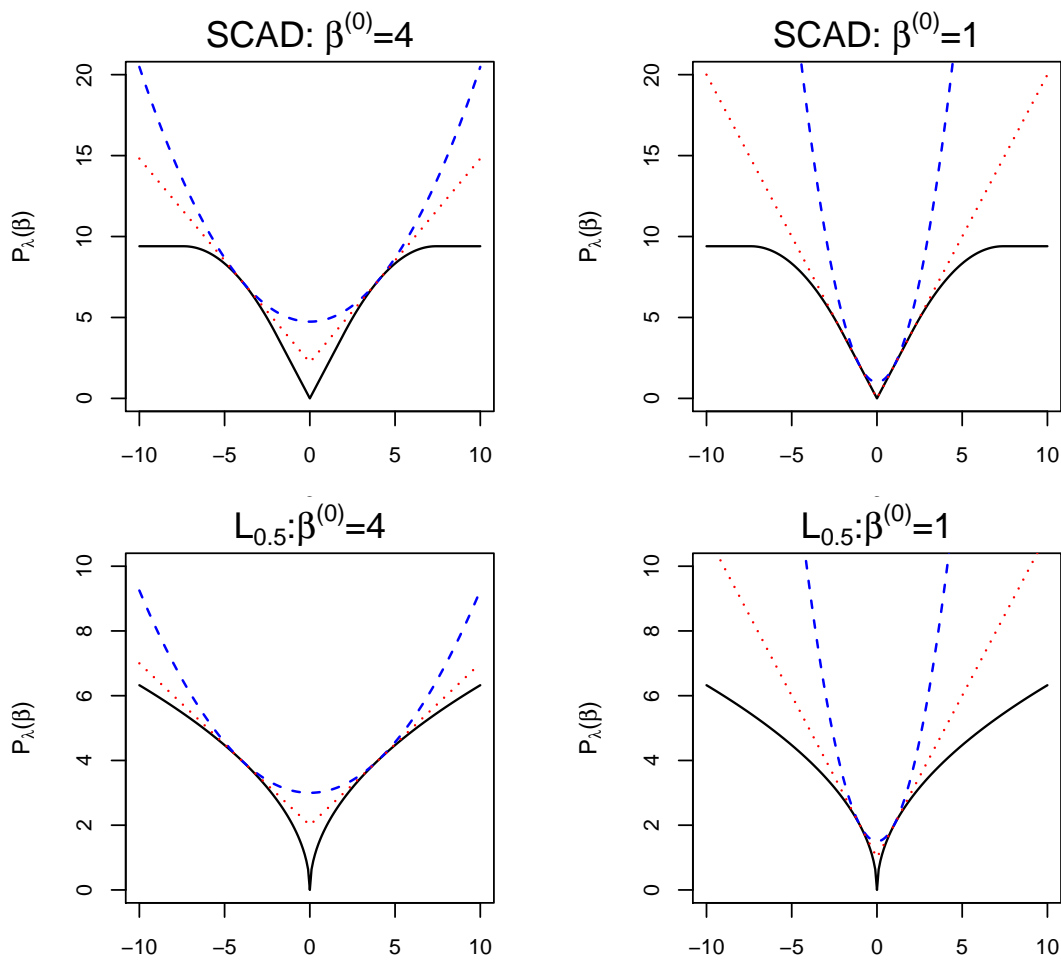
<sup>72</sup>Shooting algorithm

<sup>73</sup>Fu

<sup>74</sup>Coordinate descent algorithm

<sup>75</sup>Wu and Lang

<sup>76</sup>Iterative conditional maximization algorithm



شکل ۸.۲: تقریب درجه دو موضعی و تقریب خطی موضعی برای توابع جریمه  $L_{0.5}$  و SCAD. نمودار آبی تقریب درجه دو موضعی، نمودار قرمز تقریب خطی موضعی و نمودار مشکی تابع جریمه را نشان می‌دهد.



## فصل ۳

# انتخاب متغیر و تشخیص ساختار در مدل جمعی خطی- جزئی

در این فصل، یک روش دو مرحله‌ای برای انتخاب متغیر و تشخیص مؤلفه‌های خطی و غیرخطی در مدل جمعی با بعد بسیار بالا معرفی می‌شود. در مرحله اول، با استفاده از یک روش غربالگری آزاد-مدل، بعد فضای متغیرهای توضیحی را کاهش داده و سپس در مرحله دوم، از روشی مبتنی بر دو تابع جریمه برای انتخاب همزمان مؤلفه‌های غیرصفر و خطی استفاده می‌شود. از این فصل، مقاله‌های کاظمی و همکاران (۱۳۹۷)، کاظمی و همکاران (۲۰۱۸) و کاظمی و همکاران (۲۰۱۹a) استخراج شده است.

### ۱.۳ مقدمه

مدل جمعی ناپارامتری را به صورت

$$Y_i = \sum_{j=1}^p f_j(X_{ij}) + \varepsilon_i, \quad 1 \leq i \leq n \quad (1.3)$$

در نظر بگیرید که در آن  $Y_i$  متغیر پاسخ،  $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^T$  بردار متغیرهای توضیحی و  $f_j$ ها توابع یک متغیره هموار<sup>۱</sup> نامعلوم هستند. همچنین  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$  بردار خطای تصادفی با میانگین تابع هموار تابعی است که بر روی یک بازه مشخص، تا مرتبه معینی مشتق پذیر باشد.

صفر و  $E(\varepsilon\varepsilon^T) = \sigma^2 \mathbf{I}_n$  است. از آنجا که در عمل ممکن است تعداد زیادی از متغیرهای توضیحی دارای اثر خطی یا حتی بدون اثر روی متغیر پاسخ باشند و بقیه متغیرها به صورت غیرخطی وارد مدل شوند، اوسومر و راپرت<sup>۲</sup> (۱۹۹۹) مدل جمعی خطی-جزئی را به صورت

$$Y_i = \sum_{j \in S_1} \beta_j X_{ij} + \sum_{j \in S_2} f_j(X_{ij}) + \varepsilon_i, \quad 1 \leq i \leq n \quad (2.3)$$

معرفی کردند، که در آن  $S_1$  و  $S_2$  دو زیر مجموعه مجزا از  $\{1, \dots, p\}$  است. این مدل دو ویژگی مثبت تفسیرپذیری مدل های خطی و انعطاف پذیری مدل های ناپارامتری را با هم ترکیب می کند. از نظر آماری، مؤلفه های پارامتری دارای همگرایی سریع تر از مؤلفه های ناپارامتری هستند. در نتیجه غیرخطی در نظر گرفتن مؤلفه هایی که ذاتاً خطی هستند، کارایی برآورد را کاهش می دهد. بنابراین استفاده از مدل های جمعی خطی-جزئی در برخی از مسائل، نسبت به مدل های جمعی ناپارامتری، مناسب تر است. تاکنون بسیاری از محققین به مسئله برآورد و انتخاب متغیر در مدل جمعی خطی-جزئی پرداخته اند که از جمله می توان به لیو<sup>۳</sup> و همکاران (۲۰۱۱)، لیان (۲۰۱۲a)، گو<sup>۴</sup> و همکاران (۲۰۱۳)، دو<sup>۵</sup> و همکاران (۲۰۱۵)، روزبه (۲۰۱۵)، روزبه (۲۰۱۶)، آرشی و روزبه (۲۰۱۶)، لیو<sup>۶</sup> و همکاران (۲۰۱۷) و اکدنیز<sup>۷</sup> و همکاران (۲۰۱۸) اشاره کرد.

مهمترین پذیره در استفاده از مدل (۲.۳) این است که اجزای خطی و غیرخطی مدل، از پیش تعیین شده باشند. اما در عمل چنین اطلاعات پیشین به ندرت در دسترس است، به ویژه وقتی که تعداد متغیرهای توضیحی زیاد است، یا به عبارتی مدلی با بعد بالا داشته باشیم. بنابراین، علاوه بر تشخیص مؤلفه های غیرصفر، لازم است مؤلفه های خطی و غیرخطی نیز شناسایی شوند.

یک روش ممکن برای تشخیص ساختار خطی-جزئی در مدل جمعی، استفاده از آزمون فرض است، یعنی ابتدا مدل را به صورت جمعی ناپارامتری (۱.۳) در نظر گرفته و سپس با انجام آزمون مؤلفه به مؤلفه اثرات خطی و غیرخطی را شناسایی می کنیم (چن<sup>۸</sup> و همکاران، ۲۰۱۱). روش دیگر برای تشخیص اثرات ثابت و متغیر، به کارگیری دستورات عمل زیر است: ابتدا همه متغیرهای توضیحی پیوسته را در قسمت ناپارامتری و همه متغیرهای توضیحی گسسته را در قسمت پارامتری قرار دهید. اگر نتایج برآورد نشان دهد که اثر برخی از متغیرهای پیوسته را می توان با توابع پارامتری مشخص، مانند مدل خطی، توصیف کرد، با قرار دادن آن متغیرهای پیوسته در قسمت پارامتری، مدل جدید را برآزش دهید (لیو<sup>۹</sup> و همکاران، ۲۰۱۱). به کارگیری روش های مذکور در داده های با بعد بالا بسیار طاقت فرسا است. روش متداول برای انتخاب متغیر و تشخیص ساختار در مدل های جمعی با بعد بالا استفاده از روش های جریمه است.

<sup>2</sup>Opsomer and Ruppert

<sup>3</sup>Liu

<sup>4</sup>Guo

<sup>5</sup>Du

<sup>6</sup>Lv

<sup>7</sup>Akdeniz

<sup>8</sup>Chen

<sup>9</sup>Liu

ژانگ و همکاران (۲۰۱۱) از دو تابع جریمه برای شناسایی همزمان مؤلفه‌های خطی و غیرخطی در مدل جمعی استفاده کردند، اما نتوانستند ویژگی سازگاری در انتخاب را برای روش ارائه شده ثابت کنند. سپس هوانگ و همکاران (۲۰۱۲a) روشی موسوم به تعقیب رگرسیون نیمه پارامتری را برای تشخیص مؤلفه‌های خطی و غیرخطی با استفاده از تابع جریمه MCP ارائه دادند و ویژگی سازگاری را برای روش ارائه شده ثابت کردند، اما به موضوع انتخاب متغیر نپرداختند. لیان (۲۰۱۲b) روشی برای تشخیص مؤلفه‌های خطی با استفاده از تابع جریمه SCAD معرفی کرد، اما این روش انتخاب متغیر را انجام نمی‌دهد. لیان (۲۰۱۲c) با به کار بردن دو تابع جریمه SCAD توانست مؤلفه‌های خطی و غیرصفر را به طور همزمان تشخیص دهد. سپس لیان و همکاران (۲۰۱۲) دو تابع جریمه SCAD را برای رگرسیون چندکی به کار بردند و یک روش نیرومند برای انتخاب متغیر و تشخیص ساختار ارائه دادند. همچنین لیان و همکاران (۲۰۱۵) با استفاده از دو تابع جریمه LASSO گروهی تطبیقی به شناسایی مؤلفه‌های خطی و غیرصفر پرداختند.

مطالعات انجام شده بر اساس توابع جریمه که در بالا به آن‌ها اشاره شد، برای انتخاب متغیر در داده‌های با بعد بسیار بالا، به دلایل هزینه محاسباتی، دقت آماری و ناپایداری الگوریتمی، مناسب نیستند و لذا شناسایی متغیرهای مهم در این نوع داده‌ها و تشخیص نوع تاثیر آن‌ها بر متغیر پاسخ به لحاظ خطی یا غیرخطی بودن، همواره مورد توجه محققین است.

به‌طور کلی، مسئله انتخاب متغیر و تشخیص ساختار در بعد بالا، بر اساس رهیافت توابع جریمه، یک روش دو مرحله‌ای است که در مرحله اول، ابتدا بعد مدل توسط یک روش غربالگری مستقل کاراً کاهش یافته و در مرحله دوم از روشی مبتنی بر دو تابع جریمه برای شناسایی مؤلفه‌های غیرصفر و خطی در زیر مدل غربال شده استفاده می‌شود.

در این فصل ضمن رعایت الگوی متعارف دو مرحله‌ای ذکر شده، ابتدا ایده خود را برای کاهش بعد مدل، در قالب یک روش SIS بر اساس همبستگی فاصله‌ای بین متغیرهای توضیحی و تابع توزیع حاشیه‌ای متغیر پاسخ ارائه می‌دهیم. این روش، آزاد – مدل بوده و نشان می‌دهیم وقتی که متغیر پاسخ دم سنگین، چوله یا دارای مقادیر فرین است، عملکرد خوبی دارد. سپس دو تابع جریمه SCAD گروهی را برای تشخیص همزمان متغیرهای مهم و مؤلفه‌های خطی در زیر مدل غربال شده در مرحله قبل به کار می‌بریم. سرانجام عملکرد این روش را با مطالعات شبیه‌سازی و مثال واقعی مورد بررسی قرار می‌دهیم.

## ۲.۳ غربالگری مستقل مطمئن نیرومند بر اساس همبستگی فاصله‌ای

در این فصل، اگرچه هدف اصلی ما تمرکز بر مدل (۲.۳) است، اما از آنجایی که در مسائل واقعی مؤلفه‌های خطی و غیرخطی مشخص نیستند، ابتدا فرض می‌شود که مدل به صورت جمعی ناپارامتری (۱.۳) است، سپس روشی ارائه خواهد شد تا مؤلفه‌های خطی در مدل (۱.۳) تعیین شوند و صورت کلی رابطه (۲.۳) بدست آید. برای این منظور، ابتدا با استفاده از یک روش غربالگری نیرومند، بعد فضای متغیرها به بعد کوچکتر از  $n$  کاهش می‌یابد. این روش غربالگری بر اساس همبستگی فاصله‌ای بین متغیرهای توضیحی و تابع توزیع حاشیه‌ای متغیر پاسخ تعریف می‌شود. استفاده از تابع توزیع متغیر پاسخ به جای خود متغیر

پاسخ باعث می‌شود که روش ارائه شده نیرومند باشد. دلیل استفاده از روش غربالگری این است که کاربرد مستقیم روش‌های جریمه برای انتخاب متغیر و برآورد در مدل‌های با بعد بسیار بالا کارایی لازم را ندارند (فن و همکاران، ۲۰۰۹). بنابراین لازم است ابتدا با یک روش غربالگری بعد مدل را کاهش داده، سپس از توابع جریمه برای تشخیص ساختار مدل استفاده شود.

سزکلی و همکاران (۲۰۰۷) ضریب همبستگی فاصله‌ای را به عنوان یک معیار وابستگی بین دو بردار تصادفی معرفی کردند. به دلیل ویژگی‌های مطلوب ضریب همبستگی فاصله‌ای، که در زیر بخش ۳.۶.۲ به آن‌ها اشاره شد، لی و همکاران (۲۰۱۲) یک روش غربالگری مستقل مطمئن موسوم به DC-SIS را برای رتبه‌بندی متغیرهای توضیحی با استفاده از همبستگی فاصله‌ای آن‌ها با متغیر پاسخ ارائه دادند. در زیر بخش ۳.۶.۲ به جزئیات این روش غربالگری پرداخته شده است.

فرض کنید  $\mathbf{X} = (X_1, \dots, X_p)^T$  بردار متغیرهای توضیحی و  $\Psi y$  تکیه‌گاه متغیر پاسخ  $Y$  باشد. مجموعه اندیس متغیرهای مهم را به صورت

$$\mathbb{M} = \{j \text{ برای برخی مقادیر } y \in \Psi y, \text{ تابع } F(y|\mathbf{X}) \text{ به طور تابعی وابسته به } X_j \text{ است} : j\}$$

در نظر بگیرید که  $F(y|\mathbf{X}) = P(Y < y|\mathbf{X})$  تابع توزیع شرطی  $Y$  به شرط  $\mathbf{X}$  است. تعریف کنید  $s_n = |\mathbb{M}|$  و فرض کنید  $s_n < p$  که مقدار  $s_n$  ممکن است با  $n$  تغییر کند. در این فصل برای اندازه‌گیری همبستگی بین متغیرهای توضیحی و متغیر پاسخ، از  $F(Y)$  به جای  $Y$  در روش لی و همکاران (۲۰۱۲) استفاده می‌شود، یعنی معیار مطلوبیت حاشیه‌ای برای رتبه‌بندی متغیرها به صورت

$$\omega_j = \text{dcorr}(X_j, F(Y)) = \frac{\text{dcov}(X_j, F(Y))}{\sqrt{\text{dcov}(X_j, X_j)} \sqrt{\text{dcov}(F(Y), F(Y))}}, \quad 1 \leq j \leq p \quad (3.3)$$

تعریف می‌شود، که در آن  $F(y)$  تابع توزیع حاشیه‌ای  $Y$  و  $\text{dcov}(U, V)$  کوواریانس فاصله‌ای بین  $U$  و  $V$  است که نحوه محاسبه آن در زیر بخش ۳.۶.۲ داده شده است. این روش غربالگری نسبت به روش‌های موجود دارای دو مزیت است:

الف) بنا به ویژگی همبستگی فاصله‌ای،  $X_j$  و  $F(Y)$  مستقل اند، اگر و تنها اگر  $\text{dcorr}(X_j, F(Y)) = 0$ . بنابراین، این روش آزاد-مدل بوده و برای غربالگری متغیرها نیازی به مشخص کردن ساختار مدل نیست.

ب) چون  $F(Y)$  برای هر نوع متغیر پاسخ یک تابع کراندار است، می‌توان انتظار داشت که این روش برای متغیر پاسخ دم سنگین یا دارای مقادیر فرین، به دلیل جایگزینی  $Y$  با  $F(Y)$ ، نسبت به DC-SIS عملکرد بهتری داشته باشد.

برای پیاده‌سازی این روش غربالگری، کافی است معیار مطلوبیت حاشیه‌ای (۳.۳) بر اساس یک نمونه تصادفی برآورد شود. فرض کنید  $(\mathbf{X}_i, Y_i)_{i=1}^n$  یک نمونه تصادفی از مدل جمعی ناپارامتری (۱.۳) باشد. با استفاده از رابطه (۴.۸.۲)، نتیجه می‌شود

$$\widehat{\text{dcov}}^2(X_j, F(Y)) = \hat{S}_{j,1} + \hat{S}_{j,2} - 2\hat{S}_{j,3}$$

که در آن

$$\begin{aligned}\hat{S}_{j,1} &= \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| |F_n(Y_i) - F_n(Y_k)| \\ \hat{S}_{j,2} &= \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^n |F_n(Y_i) - F_n(Y_k)| \\ \hat{S}_{j,3} &= \frac{1}{n^3} \sum_{i=1}^n \sum_{k=1}^n \sum_{l=1}^n |X_{ij} - X_{lj}| |F_n(Y_k) - F_n(Y_l)|\end{aligned}$$

و  $F_n(y) = \frac{1}{n} \sum_{i=1}^n I(Y_i \leq y)$  تابع توزیع تجربی  $Y$  است. به طور مشابه، می‌توان کوواریانس‌های فاصله‌ای نمونه‌ای  $dcov(X_j, X_j)$  و  $dcov(F(Y), F(Y))$  را نیز برآورد کرد. در نتیجه

$$\hat{\omega}_j = \widehat{dcorr}(X_j, F(Y)) = \frac{\widehat{dcov}(X_j, F(Y))}{\sqrt{\widehat{dcov}(X_j, X_j)} \sqrt{\widehat{dcov}(F(Y), F(Y))}}.$$

برای کاهش بعد فضای متغیرها، مجموعه‌ای از متغیرهای توضیحی را که دارای مقادیر بزرگ  $\hat{\omega}_j$  هستند، به عنوان مجموعه متغیرهای مهم انتخاب می‌شوند. به عبارت دیگر

$$\hat{M} = \{1 \leq j \leq p : \hat{\omega}_j \geq cn^{-\kappa}\} \quad (۴.۳)$$

که در آن  $c > 0$  و  $0 \leq \kappa < 1/2$  و مقادیر ثابت از پیش تعیین شده و در شرط  $C_2$  از زیر بخش ۱.۲.۳ صدق می‌کنند. این روش غربالگری که بعد فضای متغیرها را از مقیاس بسیار بزرگ  $p$  به یک فضای بسیار کوچکتر با اندازه  $|\hat{M}|$  کاهش می‌دهد، غربالگری مستقل مطمئن نیرومند بر اساس همبستگی فاصله‌ای (RDC-SIS) نامیده می‌شود. استفاده از روش غربالگری RDC-SIS مستلزم تعیین یک مقدار آستانه است که در عمل، تعیین مقدار آن معمولاً مشکل است. در راستای یافتن راهکاری مناسب برای این مسئله، یک روش جایگزین، انتخاب  $d$  ( $d < n$ ) متغیر با بیشترین مطلوبیت حاشیه‌ای یا بیشترین همبستگی است. انتخاب مقدار  $d$  نقش مهمی را در مرحله غربالگری ایفا می‌کند. فن و ليو (۲۰۰۸) ضریبی از  $[n/\log(n)]$ ، مانند  $d_1 = [n/\log(n)]$ ،  $d_2 = 2[n/\log(n)]$  یا  $d_3 = 3[n/\log(n)]$  را به عنوان یک مقدار مناسب پیشنهاد دادند. البته یک انتخاب محافظه کارانه می‌تواند  $d = n - 1$  باشد.

### ۱.۲.۳ بررسی خواص نظری روش RDC-SIS

در این بخش به اثبات سازگاری برآوردگر  $\hat{\omega}_j$  و ویژگی غربالگری مطمئن روش RDC-SIS می‌پردازیم. بدین منظور، فرض می‌کنیم دو شرط زیر برقرار هستند.

C۱. ثابت‌های مثبت  $C$  و  $t$  موجودند به طوری که

$$\max_{1 \leq j \leq p} E\{\exp(t|X_j|)\} \leq C < \infty, \quad 0 < t \leq t_0.$$

C۲. معیار مطلوبیت حاشیه‌ای (۳.۳) به ازای برخی مقادیر ثابت  $c > 0$  و  $0 \leq \kappa < 1/2$  در نامساوی زیر صدق می‌کند

$$\min_{j \in \hat{M}} \omega_j \geq 2cn^{-\kappa}.$$

**قضیه ۱.۲.۳** (ویژگی سازگاری). تحت شرط C۱، برای هر  $\kappa - \frac{1}{p} < \gamma \leq \frac{1}{p}$ ، ثابت‌های مثبت  $c_۱$  و  $c_۲$  موجودند به طوری که

$$P(\max_{1 \leq j \leq p} |\hat{\omega}_j - \omega_j| \geq cn^{-\kappa}) \leq O\left(p \left[ \exp\{-c_۱ n^{1-2(\kappa+\gamma)}\} + n \exp(-c_۲ n^\gamma)\right]\right).$$

این قضیه نشان می‌دهد که  $\hat{\omega}_j$  یک برآوردگر سازگار برای  $\omega_j$  است، بدین معنی که با افزایش  $n$ ، برآورد بدست آمده بسیار نزدیک به پارامتر است، یا احتمال اینکه اختلاف بین برآورد و پارامتر بزرگ باشد، برابر با صفر است.

**قضیه ۲.۲.۳** (ویژگی غربالگری مطمئن). تحت شرط‌های C۱ و C۲، داریم

$$P(\mathbb{M} \subseteq \hat{\mathbb{M}}) \geq 1 - O\left(s_n \left[ \exp\{-c_۱ n^{1-2(\kappa+\gamma)}\} + n \exp(-c_۲ n^\gamma)\right]\right)$$

که  $s_n = |\mathbb{M}|$ . بنابراین نتیجه می‌شود

$$P(\mathbb{M} \subseteq \hat{\mathbb{M}}) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

روش‌های غربالگری مستقل با مطلوبیت حاشیه‌ای، تنها از اطلاعات حاشیه‌ای متغیرها به جای مدل کامل استفاده می‌کنند. این موضوع ممکن است عملکرد این روش‌ها را با مشکل مواجه کند. لذا روش RDC-SIS نیز مشابه سایر روش‌های غربالگری مستقل، در صورت عدم برقراری شرایط لازم برای داشتن ویژگی غربالگری مطمئن، ممکن است منجر به حذف متغیرهای مهم از مدل شود. برای رفع این چالش، یک روش غربالگری تکراری را در ادامه معرفی می‌کنیم.

### ۲.۲.۳ غربالگری مستقل مطمئن تکراری

همانطور که در زیربخش ۱.۶.۲ بیان شد، در برخی از موارد، به کاربردن SIS ممکن است منجر به حذف برخی از متغیرهای مهم شود. فن و لیو (۲۰۰۸) برای رفع مشکل و افزایش کارایی انتخاب متغیر، روش SIS را به طور مکرر برای مدل‌های خطی به کار گرفتند و این روش را ISIS نامگذاری کردند. تاثیر روش‌های SIS و ISIS به فرض خطی بودن مدل وابسته است، لذا ایده بهبود SIS توسط ISIS را نمی‌توان مستقیماً به روش غربالگری آزاد-مدل RDC-SIS تعمیم داد، مگر اینکه یک مدل مفروض برای  $X$  و  $Y$  در نظر گرفته شود. برای ارائه یک روش تکراری به منظور افزایش کارایی RDC-SIS از یک مثال استفاده می‌کنیم. مدل

$$Y = ۵X_۱ + ۵X_۲ + ۵X_۳ - ۱۵\sqrt{\rho}X_۴ + \varepsilon \quad (۵.۳)$$

را در نظر بگیرید که خطای  $\varepsilon$  و هر کدام از متغیرهای توضیحی آن از توزیع  $\mathcal{N}(0, 1)$  تولید می‌شوند. در اینجا فرض می‌شود که به استثنای  $X_۴$ ، ضریب همبستگی هر متغیر با سایر متغیرها یکسان و برابر با  $\rho \neq 0$  است و  $X_۴$  دارای همبستگی  $\sqrt{\rho}$  با سایر متغیرها است. این مثال از فن و لیو (۲۰۰۸) گرفته شده است. در این مثال، متغیر  $X_۴$  به طور توأم مهم اما به طور حاشیه‌ای با  $Y$  ناهمبسته است.

بنابراین روش غربالگری RDC-SIS در شناسایی متغیر مهم  $X_4$  ناتوان است. یک روش ممکن برای از بین بردن ناهمبستگی حاشیه‌ای بین  $X_4$  و  $Y$  و تقویت تأثیر حاشیه‌ای  $X_4$  روی  $Y$ ، حذف همبستگی بین  $X_4$  و  $(X_1, X_2, X_3)$  است. روش رایج برای از بین بردن این همبستگی، مدل کردن  $X_j$  روی  $(X_1, X_2, X_3)$  به صورت خطی برای  $4 \leq j \leq p$  است. باقیمانده‌های بدست آمده از این رگرسیون‌های خطی با  $(X_1, X_2, X_3)$  ناهمبسته هستند. سپس روش RDC-SIS برای این باقیمانده‌ها (به جای متغیر  $X_j$ ) و  $Y$  به کار برده می‌شود. باقیمانده‌های متناظر با  $X_4$  با  $Y$  ناهمبسته نیستند. بنابراین  $X_4$  به عنوان متغیر مهم انتخاب می‌شود.

با توجه به بحث فوق، یک روش تکراری کلی برای RDC-SIS ارائه می‌شود که شامل سه گام زیر است. فرض کنید  $Y = (Y_1, \dots, Y_n)^T$  و  $X$  ماتریس طرح با بعد  $n \times p$  باشد:

• گام ۱: ابتدا روش RDC-SIS را برای  $Y$  و  $X$  به کار برید. فرض کنید در این مرحله  $d_1$  متغیر توضیحی به صورت  $X_{M_1} = \{X_j : j \in M_1\}$  انتخاب می‌شوند که  $M_1$  مجموعه اندیس متغیرهای انتخاب شده با اندازه  $d$  و  $d_1 < d$  یک مقدار از پیش تعیین شده است. در اینجا برای سهولت از  $d = \lceil 2[n/\log(n)] \rceil$  استفاده می‌شود.

• گام ۲: فرض کنید  $X_1$  ماتریس طرح متناظر با متغیرهای مجموعه  $M_1$  و  $X_1^c$  ماتریس طرح متناظر با متغیرهای مجموعه  $M_1^c$  است. ماتریس  $X_{new} = \{I_n - X_1(X_1^T X_1)^{-1} X_1^T\} X_1^c$  را محاسبه کنید و سپس روش RDC-SIS را برای  $Y$  و تمام ستونهای ماتریس  $X_{new}$  به کار برید. فرض کنید در این مرحله  $d_2$  متغیر توضیحی انتخاب می‌شوند و مجموعه اندیس متغیرهای انتخاب شده را با  $M_2$  نشان می‌دهیم. مجموعه  $M_1$  را با  $M_2$  به روز رسانی کنید.

• گام ۳: گام ۲ را  $k - 1$  بار تکرار کنید تا تعداد متغیرهای توضیحی انتخاب شده از  $d$  تجاوز کند، یعنی  $d_1 + \dots + d_k \geq d$ . سرانجام، مجموعه متغیرهای انتخاب شده  $M_1 \cup \dots \cup M_k$  است.

از آنجا که متغیرهای حذف شده در مرحله پیشین، مجدداً در مرحله کنونی برای ورود به مدل بررسی می‌شوند، این الگوریتم قادر است احتمال حذف متغیرهای مهم را کاهش دهد. روش غربالگری بالا از ایده تصویرسازی متعامد استفاده می‌کند که در ژو و همکاران (۲۰۱۱) برای غربالگری آزاد-مدل استفاده شده است. در اینجا مقدار  $d_1$  در گام ۱ توسط کاربر تعیین می‌شود. در عمل،  $d_1$  به عنوان یک پارامتر تنظیم کننده در نظر گرفته می‌شود و مقدار بهینه آن با مینیمم کردن میانگین توان دوم خطای پیش بینی تعیین می‌گردد.

انتخاب یک مقدار بزرگ‌تر برای  $d$ ، احتمال انتخاب متغیرهای بی‌اهمیت را برای ورود به مدل افزایش می‌دهد. در ادامه، یک روش مبتنی بر جریمه برای حذف متغیرهای اضافی از مدل و تشخیص مؤلفه‌های خطی و غیرخطی ارائه می‌شود.

### ۳.۳ انتخاب متغیر و تشخیص ساختار

فرض کنید  $d$  متغیر در مرحله غربالگری برای ورود به مدل انتخاب شده‌اند. مدل جمعی ناپارامتری

$$Y_i = \sum_{j=1}^d f_j(X_{ij}) + \epsilon_i, \quad 1 \leq i \leq n$$

را در نظر بگیرید. برای تقریب توابع هموار نامعلوم از توابع پایه‌ای  $B$  - اسپلاین به صورت

$$f_j(x) \approx \sum_{k=1}^K \beta_{jk} B_{jk}(x) \quad j = 1, \dots, d$$

استفاده می‌شود که  $K$  تعداد توابع پایه‌ای استفاده شده است. برای اطلاعات بیشتر درباره توابع  $B$  - اسپلاین به دی بور<sup>۱۰</sup> (۲۰۰۱) و پیوست پ مراجعه کنید.

یک روش رایج برای برآورد و انتخاب متغیر به‌طور همزمان، استفاده از رگرسیون جریمه شده است. در اینجا برای انتخاب متغیر و تشخیص ساختار مدل، از دو تابع جریمه به‌طور همزمان استفاده می‌شود، یعنی بردار ضرایب  $\beta = (\beta_1^T, \dots, \beta_d^T)^T$ ، که  $\beta_j = (\beta_{j1}, \dots, \beta_{jK})^T$ ، با حل مسئله بهینه‌سازی

$$\hat{\beta} = \arg \min_{\beta} \left\{ \frac{1}{2} \sum_{i=1}^n (Y_i - \sum_{j=1}^d \sum_{k=1}^K \beta_{jk} B_{jk}(X_{ij}))^2 + n \sum_{j=1}^d p_{\lambda_1}(\|\beta_j\|_{A_j}) + n \sum_{j=1}^d p_{\lambda_2}(\|\beta_j\|_{D_j}) \right\} \quad (۶.۳)$$

برآورد می‌شود که در آن  $p_{\lambda_1}(\cdot)$  و  $p_{\lambda_2}(\cdot)$  توابع جریمه SCAD با پارامترهای تنظیم کننده  $\lambda_1$  و  $\lambda_2$  به ترتیب برای شناسایی مؤلفه‌های صفر و خطی به کار می‌روند. علاوه بر این،  $A_j$  و  $D_j$  دو ماتریس  $K \times K$  هستند.  $\|\beta_j\|_{D_j} = (\beta_j^T D_j \beta_j)^{\frac{1}{2}}$  و  $\|\beta_j\|_{A_j} = (\beta_j^T A_j \beta_j)^{\frac{1}{2}}$ ، انتخاب ماتریس‌های  $A_j$  و  $D_j$  بسیار حائز اهمیت است. این ماتریس‌ها باید طوری انتخاب شوند که  $\|\beta_j\|_{A_j} = 0$ ، اگر فقط  $\sum_k \beta_{jk} B_{jk}(x) \equiv 0$  و به‌طور مشابه  $\|\beta_j\|_{D_j} = 0$ ، اگر و فقط اگر  $\sum_k \beta_{jk} B_{jk}(x)$  یک تابع خطی باشد، یعنی دو تابع جریمه بتوانند برای شناسایی مؤلفه‌های صفر و خطی به کار روند. در اینجا  $A_j$  و  $D_j$  به صورت  $A_j = \{\int_0^1 B_{jk}(x) B_{jk'}(x) dx\}_{k,k'=1}^K$  و  $D_j = \{\int_0^1 B_{jk}''(x) B_{jk'}''(x) dx\}_{k,k'=1}^K$  انتخاب می‌شوند. انتخاب  $D_j$  به صورت فوق از این واقعیت نتیجه می‌شود که مشتق دوم یک تابع خطی برابر صفر است. فرض کنید  $Y = (Y_1, \dots, Y_n)^T$  و  $Z = (Z_1, \dots, Z_d)$  که در آن

$$Z_j = \begin{pmatrix} B_{j1}(X_{1j}) & B_{j2}(X_{1j}) & \dots & B_{jK}(X_{1j}) \\ \vdots & \vdots & & \vdots \\ B_{j1}(X_{nj}) & B_{j2}(X_{nj}) & \dots & B_{jK}(X_{nj}) \end{pmatrix}.$$

<sup>10</sup>De Boor



پس رابطه (۶.۳) را می‌توان به صورت ماتریسی

$$\hat{\beta} = \arg \min_{\beta} \left\{ \frac{1}{\nu} \|\mathbf{Y} - \mathbf{Z}\beta\|^2 + n \sum_{j=1}^d p_{\lambda_1}(\|\beta_j\|_{\mathbf{A}_j}) + n \sum_{j=1}^d p_{\lambda_2}(\|\beta_j\|_{\mathbf{D}_j}) \right\}. \quad (7.3)$$

نوشت. به دلیل سهولت، تابع هدف سمت راست رابطه (۷.۳) با  $Q(\beta)$  نشان داده می‌شود. برای حل مسئله بهینه‌سازی (۷.۳) از تقریب موضعی درجه دو (LQA) معرفی شده در زیر بخش ۷.۲.۱ استفاده می‌کنیم. با استفاده از بسط تیلور و داشتن یک مقدار اولیه  $\beta_j^{(0)}$ ، توابع جریمه را می‌توان به صورت

$$p_{\lambda_1}(\|\beta_j\|_{\mathbf{A}_j}) \approx p_{\lambda_1}(\|\beta_j^{(0)}\|_{\mathbf{A}_j}) + \frac{1}{\nu} \frac{p'_{\lambda_1}(\|\beta_j^{(0)}\|_{\mathbf{A}_j})}{\|\beta_j^{(0)}\|_{\mathbf{A}_j}} \{ \|\beta_j\|_{\mathbf{A}_j} - \|\beta_j^{(0)}\|_{\mathbf{A}_j} \}$$

و

$$p_{\lambda_2}(\|\beta_j\|_{\mathbf{D}_j}) \approx p_{\lambda_2}(\|\beta_j^{(0)}\|_{\mathbf{D}_j}) + \frac{1}{\nu} \frac{p'_{\lambda_2}(\|\beta_j^{(0)}\|_{\mathbf{D}_j})}{\|\beta_j^{(0)}\|_{\mathbf{D}_j}} \{ \|\beta_j\|_{\mathbf{D}_j} - \|\beta_j^{(0)}\|_{\mathbf{D}_j} \}$$

تقریب زد. با به کار بردن تقریب‌های فوق و حذف جملات ثابت می‌توان رابطه (۷.۳) را به صورت

$$Q(\beta) = \frac{1}{n} \|\mathbf{Y} - \mathbf{Z}\beta\|^2 + \frac{1}{\nu} \beta^T (\Omega_1 + \Omega_2) \beta \quad (8.3)$$

خلاصه کرد، که در آن  $\Omega_1$  و  $\Omega_2$  دو ماتریس بلوکی با بعد  $dK \times dK$  به صورت

$$\Omega_1 = \text{diag} \left( \frac{p'_{\lambda_1}(\|\beta_1^{(0)}\|_{\mathbf{A}_1})}{\|\beta_1^{(0)}\|_{\mathbf{A}_1}} \mathbf{A}_1, \dots, \frac{p'_{\lambda_1}(\|\beta_d^{(0)}\|_{\mathbf{A}_d})}{\|\beta_d^{(0)}\|_{\mathbf{A}_d}} \mathbf{A}_d \right)$$

و

$$\Omega_2 = \text{diag} \left( \frac{p'_{\lambda_2}(\|\beta_1^{(0)}\|_{\mathbf{D}_1})}{\|\beta_1^{(0)}\|_{\mathbf{D}_1}} \mathbf{D}_1, \dots, \frac{p'_{\lambda_2}(\|\beta_d^{(0)}\|_{\mathbf{D}_d})}{\|\beta_d^{(0)}\|_{\mathbf{D}_d}} \mathbf{D}_d \right)$$

هستند. توجه شود که رابطه (۸.۳) یک تابع درجه دو بر حسب  $\beta$  بوده و در نتیجه دارای جواب صریح

$$\hat{\beta} = (\mathbf{Z}^T \mathbf{Z} + n(\Omega_1 + \Omega_2))^{-1} \mathbf{Z}^T \mathbf{Y} \quad (9.3)$$

است. برآورد بدست آمده  $\hat{\beta}$ ، به عنوان برآورد اولیه در تکرار بعدی استفاده می‌شود. این الگوریتم تا رسیدن به همگرایی تکرار می‌شود. بر اساس پیشنهاد فن و لی (۲۰۰۱)، در هر تکرار، به محض اینکه برخی از مؤلفه‌های  $\|\beta_j\|_{\mathbf{A}_j}$  یا  $\|\beta_j\|_{\mathbf{D}_j}$  کوچکتر از  $10^{-6}$  شود، مولفه  $f_j(x)$  به ترتیب به عنوان تابع صفر یا تابع خطی تشخیص داده می‌شود.

همانطور که در بخش ۷.۲ اشاره شد، تقریب LQA یک الگوریتم پیش‌رونده است، یعنی اگر متغیری در یک تکرار از مدل حذف شود، در تکرار بعدی وارد مدل نمی‌شود و در نتیجه آن متغیر توضیحی از مدل نهایی حذف می‌شود. بنابراین، با بکار بردن تقریب LQA، ممکن است برخی متغیرهای توضیحی مهم را از دست بدهیم. برای رفع این مشکل، هانترو و لی<sup>۱۱</sup> (۲۰۰۵) اعمال برخی تغییرات جزئی را در

<sup>۱۱</sup>Hunter and Li

این الگوریتم پیشنهاد دادند. به طور دقیق تر، برای  $\varepsilon_1 > 0$  و  $\varepsilon_2 > 0$ ، رابطه (۹.۳) را، با اندکی تغییر، می توان بصورت زیر نوشت

$$\beta = (\mathbf{Z}^T \mathbf{Z} + n(\Omega_{\varepsilon_1} + \Omega_{\varepsilon_2}))^{-1} \mathbf{Z}^T \mathbf{Y} \quad (10.3)$$

که

$$\Omega_{\varepsilon_1} = \text{diag} \left( \frac{p'_{\lambda_1}(\|\beta_1^{(0)}\|_{\mathbf{A}_1})}{\|\beta_1^{(0)}\|_{\mathbf{A}_1} + \varepsilon_1} \mathbf{A}_1, \dots, \frac{p'_{\lambda_1}(\|\beta_d^{(0)}\|_{\mathbf{A}_d})}{\|\beta_d^{(0)}\|_{\mathbf{A}_d} + \varepsilon_1} \mathbf{A}_d \right)$$

و

$$\Omega_{\varepsilon_2} = \text{diag} \left( \frac{p'_{\lambda_2}(\|\beta_1^{(0)}\|_{\mathbf{D}_1})}{\|\beta_1^{(0)}\|_{\mathbf{D}_1} + \varepsilon_2} \mathbf{D}_1, \dots, \frac{p'_{\lambda_2}(\|\beta_d^{(0)}\|_{\mathbf{D}_d})}{\|\beta_d^{(0)}\|_{\mathbf{D}_d} + \varepsilon_2} \mathbf{D}_d \right).$$

بر اساس پیشنهاد هانتز و لی (۲۰۰۵)، در پیاده سازی این الگوریتم،  $\varepsilon_1$  و  $\varepsilon_2$  را به صورت زیر در نظر می گیریم

$$\varepsilon_1 = \frac{\tau}{2n\lambda_1} \min\{\|\beta_j^{(0)}\|_{\mathbf{A}_j}, j = 1, \dots, d\}, \quad \varepsilon_2 = \frac{\tau}{2n\lambda_2} \min\{\|\beta_j^{(0)}\|_{\mathbf{D}_j}, j = 1, \dots, d\}$$

که  $\tau = 10^{-8}$ . قرار دهید  $k = 0$  و با استفاده از الگوریتم زیر جواب نهایی را بدست آورید:

- گام ۱:  $\beta = \beta^{(k)}$  را به عنوان مقدار اولیه در نظر بگیرید.
- گام ۲: با استفاده از رابطه (۱۰.۳)، مقدار  $\beta$  را به روز رسانی کرده تا  $\beta^{(k+1)}$  حاصل شود.
- گام ۳: قرار دهید  $k = k + 1$  و گام های ۱ و ۲ را تا رسیدن به همگرایی تکرار کنید.

در ابتدا، برآورد کمترین توان های دوم به عنوان برآورد اولیه در نظر گرفته می شود. در گام ۳ تکرار تا زمانی ادامه می یابد که

$$\max_{1 \leq j \leq p} \left| \|\beta_j^{(k)}\|_{\mathbf{A}_j} - \|\beta_j^{(k+1)}\|_{\mathbf{A}_j} \right| + \max_{1 \leq j \leq p} \left| \|\beta_j^{(k)}\|_{\mathbf{D}_j} - \|\beta_j^{(k+1)}\|_{\mathbf{D}_j} \right| < 10^{-6}.$$

پس از همگرایی الگوریتم فوق، برآورد ضرایب دقیقاً برابر با صفر نیست. لذا لازم است از یک مقدار آستانه استفاده شود. در مطالعات شبیه سازی و تحلیل داده های واقعی، اگر  $\|\beta_j^{(k)}\|_{\mathbf{A}_j} < 10^{-6}$ ، قرار دهید  $\hat{\beta}_j = 0$  و به طور مشابه، اگر  $\|\beta_j\|_{\mathbf{D}_j} < 10^{-6}$ ، مولفه به عنوان تابع خطی انتخاب می شود. نتایج شبیه سازی نشان می دهند که این انتخاب مناسب است.

## ۴.۳ مطالعه شبیه سازی

در این بخش، عملکرد روش دو مرحله ای مذکور با انجام شبیه سازی بررسی می شود. ابتدا، به عنوان توضیحی در خصوص مثال های عددی، در مثال ۱.۴.۳ حساسیت روش RDC-SIS نسبت به پارامتر  $d$  تحلیل می شود. سپس در مثال ۲.۴.۳، عملکرد RDC-SIS با DC-SIS (لی و همکاران، ۲۰۱۲)، SIRS

(ژو و همکاران، ۲۰۱۱)، SIS (فن و لیو، ۲۰۰۸) و NIS (فن و همکاران، ۲۰۱۱) مقایسه می‌کنیم. در مثال‌های ۳.۴.۳ و ۴.۴.۳، عملکرد روش تکراری RDC-ISIS را با روش‌های مذکور و همچنین روش تکراری ISIS (فن و لیو، ۲۰۰۸) مقایسه می‌شود. عملکرد این روش‌ها را با سه معیار  $S, P_j, M$  ارزیابی می‌کنیم که  $M$  حداقل اندازه مدل برای در برگرفتن تمام متغیرهای مهم،  $P_j$  احتمال تجربی انتخاب متغیر مهم  $X_j$  و  $S$  احتمال انتخاب تمام متغیرهای مهم برای اندازه مدل داده شده است. به منظور استنباط بهتر، در مثال ۲.۴.۳ چندک‌های ۵٪، ۲۵٪، ۵۰٪، ۷۵٪ و ۹۵٪ معیار  $M$  در ۵۰۰ تکرار نیز ارائه شده است. توجه شود که معیار  $M$  نیازی به مشخص کردن مقدار آستانه ندارد. اگر مقادیر  $M$  مربوط به یک روش نزدیک به تعداد متغیرهای مهم باشند، آن روش از عملکرد مطلوبی برخوردار است. همچنین در یک روش غربالگری مناسب مقادیر  $P_j$  و  $S$  باید نزدیک به یک باشند.

برای پیاده‌سازی RDC-ISIS، مطالعات تجربی نشان می‌دهند که تعداد کمی از تکرارها کافی است و می‌تواند هزینه محاسبات را کاهش دهد. با تکرار بیشتر این الگوریتم ممکن است احتمال حذف متغیرهای مهم کاهش یابد، اما هزینه محاسبات افزایش می‌یابد. در این رساله، برای شبیه‌سازی با انتخاب  $d_1 = 5$  و  $d_2 = d - 5$  الگوریتم فقط یکبار تکرار می‌شود.

در مثال ۵.۴.۳ کارایی روش دو مرحله‌ای در شناسایی مؤلفه‌های غیرصفر و همچنین تشخیص مؤلفه‌های خطی و غیرخطی در یک مدل جمعی خطی-جزئی با استفاده از دو تابع جریمه بررسی می‌شود. برای انتخاب بهینه پارامترهای  $\lambda_1$  و  $\lambda_2$  از معیار BIC به صورت

$$\log\left(\frac{1}{n}\|\mathbf{Y} - \mathbf{Z}\hat{\beta}_\lambda\|^2\right) + d_1 \frac{\log(n/K)}{n/K} + d_2 \frac{\log n}{n}$$

استفاده می‌شود، که در آن مقدار مینیمم کننده (۷.۳) برای  $\lambda = (\lambda_1, \lambda_2)$  داده شده است. همچنین  $d_1$  تعداد مؤلفه‌های ناپارامتری و  $d_2$  تعداد مؤلفه‌های پارامتری به ازای مقدار داده شده  $\lambda$  است.

**مثال ۱.۴.۳.** برای تحلیل حساسیت عملکرد RDC-SIS نسبت به مقادیر مختلف  $d$ ، داده‌ها از مدل

$$Y = 5g_1(X_1) + 3g_2(X_2) + 4g_3(X_3) + 6g_4(X_4) + \sqrt{174}\varepsilon$$

تولید شده‌اند، که در آن

$$g_1(x) = x \quad g_2(x) = (2x - 1)^2 \quad g_3(x) = \sin(2\pi x)/(2 - \sin(2\pi x))$$

$$g_4(x) = 0.1 \sin(2\pi x) + 0.2 \cos(2\pi x) + 0.3 \sin(2\pi x)^2 + 0.4 \cos(2\pi x)^3 + 0.5 \sin(2\pi x)^3$$

و متغیرهای توضیحی دارای توزیع حاشیه‌ای نرمال استاندارد و همبستگی  $Cov(X_i, X_j) = 0.8^{|i-j|}$  هستند. سه مقدار مختلف  $d_1 = \lfloor n/\log(n) \rfloor$ ،  $d_2 = 2\lfloor n/\log(n) \rfloor$ ،  $d_3 = 3\lfloor n/\log(n) \rfloor$  را برای  $d$  در نظر گرفته و شبیه‌سازی برای مقادیر مختلف  $(n, p)$  انجام شده است. همچنین برای خطا دو توزیع نرمال استاندارد و تی-استودنت با یک درجه آزادی در نظر گرفته شده است. نتایج شبیه‌سازی پس از ۵۰۰ بار تکرار در جدول ۱.۳ گزارش شده است. همان‌طور که ملاحظه می‌شود، به‌ازای هر  $(n, p)$ ، با

جدول ۱.۳: احتمال تجربی  $P_j$  و احتمال  $S$  در مثال ۱.۴.۳

$\varepsilon \sim t(1)$					$\varepsilon \sim \mathcal{N}(0, 1)$					$d$	$n$	$p$
$S$	$P_4$	$P_3$	$P_2$	$P_1$	$S$	$P_4$	$P_3$	$P_2$	$P_1$			
۰.۴۱	۰.۴۵	۰.۹۲	۱.۰۰	۰.۸۴	۰.۵۳	۰.۶۰	۰.۹۶	۱.۰۰	۰.۹۱	$d_1$		
۰.۵۸	۰.۶۲	۰.۹۶	۱.۰۰	۰.۹۴	۰.۷۳	۰.۷۵	۰.۹۸	۱.۰۰	۰.۹۸	$d_2$	۱۰۰	
۰.۶۷	۰.۷۰	۰.۹۷	۱.۰۰	۰.۹۶	۰.۸۱	۰.۸۱	۰.۹۹	۱.۰۰	۱.۰۰	$d_3$		
۰.۹۳	۰.۹۳	۱.۰۰	۱.۰۰	۱.۰۰	۰.۹۹	۰.۹۹	۱.۰۰	۱.۰۰	۱.۰۰	$d_1$		
۰.۹۷	۰.۹۷	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_2$	۲۰۰	۱۰۰۰
۰.۹۸	۰.۹۸	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_3$		
۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_1$		
۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_2$	۴۰۰	
۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_3$		
۰.۴۲	۰.۵۵	۰.۹۰	۱.۰۰	۰.۷۷	۰.۵۵	۰.۶۱	۰.۹۶	۱.۰۰	۰.۹۲	$d_1$		
۰.۵۸	۰.۶۶	۰.۹۶	۱.۰۰	۰.۸۹	۰.۷۱	۰.۷۲	۰.۹۹	۱.۰۰	۰.۹۶	$d_2$	۱۰۰	
۰.۶۵	۰.۷۰	۰.۹۸	۱.۰۰	۰.۹۲	۰.۷۸	۰.۷۹	۰.۹۹	۱.۰۰	۰.۹۹	$d_3$		
۰.۹۲	۰.۹۷	۱.۰۰	۱.۰۰	۰.۹۹	۰.۹۷	۰.۹۷	۱.۰۰	۱.۰۰	۱.۰۰	$d_1$		
۰.۹۷	۰.۹۷	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_2$	۲۰۰	۲۰۰۰
۰.۹۸	۰.۹۸	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_3$		
۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_1$		
۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_2$	۴۰۰	
۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	۱.۰۰	$d_3$		

افزایش مقدار  $d$  عملکرد RDC-SIS بهبود می‌یابد. برای  $n = 100$ ، عملکرد RDC-SIS به ازای  $d_1$  و  $d_2$  چندان مطلوب نیست و امکان حذف شدن متغیر مهم  $X_4$  وجود دارد. لذا برای حجم نمونه کوچک مقدار  $d_3$  پیشنهاد می‌شود. اما برای  $n = 200, 400$  و هر دو نوع توزیع خطا، هر یک از مقادیر  $d_1, d_2$  و  $d_3$  را می‌توان استفاده کرد.

در مثال‌های ۲.۴.۳ تا ۴.۴.۳، حجم نمونه  $n = 200$ ، تعداد متغیرهای توضیحی را  $p = 1000$  و  $d = 2 \lceil n / \log(n) \rceil$  در نظر گرفته و شبیه‌سازی ۵۰۰ بار تکرار شده است.

### مثال ۲.۴.۳. دو مدل به صورت

$$Y = X_1 + 2X_2 + 3X_3 + 4X_4 + \varepsilon \quad \text{مدل ۲:} \quad Y = c\beta^T X + \sigma\varepsilon \quad \text{مدل ۱:}$$

را در نظر بگیرید، که در آن  $\beta = (1, 0.8, 0.6, 0.4, 0.2, 0, \dots, 0)^T$  و  $\sigma^2 = 6.83$ . در مدل ۱، به منظور کنترل نسبت سیگنال به نوفه<sup>۱۲</sup>، مقادیر  $c = 1, 2$  در نظر گرفته شده است. در هر دو مدل، بردار

<sup>12</sup>Signal-to-noise ratio

جدول ۲.۳: چندک‌های  $M$ ، احتمال تجربی  $P_j$  و احتمال  $S$  در مدل ۱

$S$	$P_j$					$M$					روش	خطا $c$	
	$P_5$	$P_4$	$P_3$	$P_2$	$P_1$	%۹۵	%۷۵	%۵۰	%۲۵	%۵			
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	RDC-SIS	۱	$\mathcal{N}(0, 1)$
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	DC-SIS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	SIS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	SIRS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	RDC-SIS	۲	
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	DC-SIS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	SIS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	SIRS		
۰/۹۹	۰/۹۹	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱۹	۵	۵	۵	۵	RDC-SIS	۱	$t(3)$
۰/۶۷	۰/۷۱	۰/۷۷	۰/۸۶	۰/۸۴	۰/۸۲	۵۸۴	۱۰۹	۱۹	۷	۵	DC-SIS		
۰/۱۰	۰/۱۶	۰/۲۰	۰/۲۱	۰/۲۱	۰/۲۰	۹۶۶	۹۱۶	۸۰۶	۴۶۷	۳۵	SIS		
۰/۹۸	۰/۹۸	۰/۹۹	۰/۹۹	۱/۰۰	۱/۰۰	۲۵	۶	۵	۵	۵	SIRS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	RDC-SIS	۲	
۰/۸۷	۰/۹۰	۰/۹۴	۰/۹۴	۰/۹۴	۰/۹۵	۱۵۶	۶	۵	۵	۵	DC-SIS		
۰/۲۳	۰/۳۰	۰/۴۲	۰/۴۴	۰/۴۳	۰/۴۴	۹۸۱	۸۶۵	۴۹۴	۹۰	۵	SIS		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۵	۵	۵	۵	SIRS		

متغیرهای توضیحی از توزیع نرمال چند متغیره با میانگین صفر و ماتریس کوواریانس  $\Sigma = (\sigma_{ij})_{p \times p}$  تولید می‌شود که  $\sigma_{ij} = 0.5^{|i-j|}$  در مدل ۱، برای خطا دو توزیع نرمال استاندارد و توزیع تی-استودنت با سه درجه آزادی در نظر گرفته شده است. در مدل ۲، علاوه بر دو توزیع فوق، توزیع نرمال چوله با پارامترهای  $\alpha = 2, \sigma = 1, \mu = 0$  نیز در نظر گرفته شده و نتایج شبیه‌سازی پس از ۵۰۰ تکرار در جداول ۲.۳ و ۳.۳ ارائه شده‌اند.

با نگاه کلی به مقادیر جدول ۲.۳، می‌توان کارایی روش غربالگری پیشنهادی را دریافت. با توجه به جدول ۲.۳، هنگامی که خطا دارای توزیع نرمال استاندارد است، برای هر دو حالت  $c = 1$  و  $c = 2$  هر چهار روش در شناسایی متغیرهای مهم  $X_1, X_2, X_3, X_4, X_5$  بسیار خوب عمل می‌کنند و همواره هر ۵ متغیر مهم را به درستی شناسایی می‌کنند، زیرا برای همه روش‌ها، مقادیر  $P_j$  ( $1 \leq j \leq 5$ ) و  $S$  برابر با یک است. اما برای خطای غیر نرمال نتایج متفاوت است. برای توزیع خطای تی-استودنت، عملکرد روش‌های DC-SIS و SIS بسیار ضعیف است، در حالی که روش RDC-SIS با احتمال تجربی تقریباً ۱۰۰٪ متغیرهای مهم را به درستی تشخیص می‌دهد. به‌عنوان مثال، برای  $c = 1$ ، مقدار  $P_1 = 1$  بدین معنی است که RDC-SIS در ۱۰۰٪ دفعات، متغیر  $X_1$  را به‌عنوان متغیر مهم انتخاب می‌کند و همچنین مقدار  $S = 99\%$  بدین معنی است که روش مذکور در ۹۹٪ اوقات ۵ متغیر مهم را به درستی تشخیص می‌دهد. روش SIRS نیز در شناسایی متغیرهای مهم خوب عمل می‌کند و عملکرد دو روش

جدول ۳.۳: چندک‌های  $M$ ، احتمال تجربی  $P_j$  و احتمال  $S$  در مدل ۲

$S$	$P_j$				$M$					روش	خطا
	$P_4$	$P_3$	$P_2$	$P_1$	%۹۵	%۷۵	%۵۰	%۲۵	%۵		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۴	۴	۴	۴	۴	RDC-SIS	$\mathcal{N}(0, 1)$
۰/۹۹	۱/۰۰	۱/۰۰	۱/۰۰	۰/۹۹	۹	۶	۵	۵	۴	DC-SIS	
۰/۹۳	۱/۰۰	۱/۰۰	۱/۰۰	۰/۹۳	۱۰۲	۱۱	۶	۵	۵	NIS	
۰/۸۵	۰/۸۵	۱/۰۰	۱/۰۰	۱/۰۰	۲۹۹	۱۹	۵	۴	۴	SIRS	
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۵	۴	۴	۴	۴	RDC-SIS	$t(3)$
۰/۹۴	۰/۹۸	۰/۹۸	۰/۹۶	۰/۹۴	۶۴	۷	۵	۵	۴	DC-SIS	
۰/۵۵	۰/۸۹	۰/۸۵	۰/۸۱	۰/۶۵	۸۶۴	۸۲	۱۱	۶	۵	NIS	
۰/۶۹	۰/۷۸	۱/۰۰	۱/۰۰	۰/۹۹	۴۶۸	۳۴	۷	۴	۴	SIRS	
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۴	۴	۴	۴	۴	RDC-SIS	SN
۰/۹۹	۱/۰۰	۱/۰۰	۱/۰۰	۰/۹۹	۸	۶	۵	۵	۴	DC-SIS	
۰/۹۲	۱/۰۰	۱/۰۰	۱/۰۰	۰/۹۲	۱۰۴	۱۱	۶	۵	۵	NIS	
۰/۸۷	۰/۸۷	۱/۰۰	۱/۰۰	۱/۰۰	۳۴۵	۲۳	۵	۴	۴	SIRS	

RDC-SIS و SIRS تقریباً یکسان و بسیار رضایت‌بخش است. با توجه به مقادیر  $P_j$  و  $S$  در جدول ۲.۳، برای خطای تی-استودنت، روش SIS شانس بسیار کمی برای انتخاب متغیرهای مهم دارد. این روش در حالت  $c = 1$  با احتمال تجربی  $10\%$  و در حالت  $c = 2$  با احتمال  $23\%$  همه متغیرهای مهم را انتخاب می‌کند.

جدول ۳.۳ که حاوی نتایج مدل ۲ است، نشان می‌دهد که برای هر سه نوع توزیع خطا، عملکرد روش RDC-SIS عالی است و نسبت به روش‌های دیگر بهتر عمل می‌کند. برای  $1 \leq j \leq 4$ ، مقادیر  $P_j$  و  $S$  مربوط به روش RDC-SIS برابر یک است. بنابراین، برای هر سه نوع توزیع خطا، روش RDC-SIS هر چهار متغیر مهم را برای ورود به مدل انتخاب می‌کند. همچنین در این مدل، عملکرد DC-SIS بسیار مشابه روش RDC-SIS است. این دو روش نسبت به NIS و SIRS برتری دارند.

همانطور که اشاره شد، در بسیاری از مدل‌ها، روش‌های غربالگری حاشیه‌ای در شناسایی متغیرهای مهم با شکست مواجه می‌شوند. در مثال ۳.۴.۳ متغیر مهم  $X_4$  به‌طور حاشیه‌ای با متغیر پاسخ ناهمبسته است و لذا توسط روش‌های غربالگری حاشیه‌ای انتخاب نمی‌شود. بنابراین از روش تکراری RDC-SIS برای بازگرداندن متغیرهای مهم به مدل استفاده می‌شود.

مثال ۳.۴.۳. در این مثال، عملکرد روش RDC-SIS را با روش‌های SIS، ISIS، SIRS، DC-SIS و RDC-SIS در مدل خطی (۵.۳)، یعنی

$$Y = 5X_1 + 5X_2 + 5X_3 - 15\sqrt{\rho}X_4 + \varepsilon \quad \text{مدل ۱:}$$

مقایسه می‌شود. علاوه بر این، یک مدل خطی دیگر، با نسبت سیگنال به نوفه ضعیف‌تر، به صورت

$$Y = 2/5X_1 + 2/5X_2 + 2/5X_3 - 7/5\sqrt{\rho}X_4 + \varepsilon \quad \text{مدل ۲:}$$

جدول ۴.۳: احتمال تجربی  $P_j$  و احتمال  $S$  در مثال ۳

مدل ۲					مدل ۱					روش	$\rho$
$S$	$P_4$	$P_2$	$P_2$	$P_1$	$S$	$P_4$	$P_2$	$P_2$	$P_1$		
۰/۰۱	۰/۰۱	۱/۰۰	۱/۰۰	۱/۰۰	۰/۰۰	۰/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	SIS	۰/۲
۰/۹۵	۰/۹۵	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	ISIS	
۰/۰۳	۰/۰۳	۱/۰۰	۱/۰۰	۱/۰۰	۰/۰۳	۰/۰۳	۱/۰۰	۱/۰۰	۱/۰۰	SIRS	
۰/۰۲	۰/۰۲	۱/۰۰	۱/۰۰	۱/۰۰	۰/۰۰	۰/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	DC-SIS	
۰/۰۴	۰/۰۴	۱/۰۰	۱/۰۰	۱/۰۰	۰/۰۳	۰/۰۳	۱/۰۰	۱/۰۰	۱/۰۰	RDC-SIS	
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	RDC-ISIS	
۰/۰۰	۰/۰۰	۰/۹۹	۰/۹۸	۱/۰۰	۰/۰۱	۰/۰۱	۱/۰۰	۱/۰۰	۱/۰۰	SIS	۰/۵
۰/۹۴	۰/۹۴	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	ISIS	
۰/۰۱	۰/۰۱	۱/۰۰	۰/۹۹	۱/۰۰	۰/۰۰	۰/۰۰	۱/۰۰	۰/۹۹	۱/۰۰	SIRS	
۰/۰۱	۰/۰۱	۰/۹۹	۰/۹۸	۱/۰۰	۰/۰۳	۰/۰۳	۰/۹۹	۱/۰۰	۱/۰۰	DC-SIS	
۰/۰۱	۰/۰۱	۱/۰۰	۰/۹۹	۱/۰۰	۰/۰۲	۰/۰۲	۱/۰۰	۰/۹۹	۱/۰۰	RDC-SIS	
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	RDC-ISIS	
۰/۰۰	۰/۰۰	۰/۸۶	۰/۸۷	۰/۸۵	۰/۰۰	۰/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	SIS	۰/۸
۰/۶۰	۰/۶۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	ISIS	
۰/۰۰	۰/۰۰	۰/۹۳	۰/۹۴	۰/۹۳	۰/۰۰	۰/۰۰	۰/۹۷	۰/۹۵	۰/۹۵	SIRS	
۰/۰۰	۰/۰۰	۰/۸۴	۰/۸۳	۰/۸۴	۰/۰۰	۰/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	DC-SIS	
۰/۰۰	۰/۰۰	۰/۹۴	۰/۹۵	۰/۹۳	۰/۰۰	۰/۰۰	۰/۹۵	۰/۹۴	۰/۹۵	RDC-SIS	
۰/۹۱	۰/۹۳	۱/۰۰	۱/۰۰	۰/۹۹	۰/۹۵	۰/۹۵	۱/۰۰	۱/۰۰	۱/۰۰	RDC-ISIS	

در نظر گرفته شده است. در این مثال، با فرض  $\varepsilon \sim \mathcal{N}(0, 1)$  و  $\rho = 0.2, 0.5, 0.8$ ، شبیه‌سازی ۵۰۰ بار تکرار شده است. در این مدل‌ها، متغیر  $X_4$  به‌طور توأم مهم، اما به‌طور حاشیه‌ای با  $Y$  ناهمبسته است. بنابراین روش‌های غربالگری حاشیه‌ای  $SIS$ ،  $SIRS$ ،  $DC-SIS$  و  $RDC-SIS$  به سختی می‌توانند متغیر مهم  $X_4$  را شناسایی کنند. اما همان‌طور که در جدول ۴.۳ ملاحظه می‌شود روش  $RDC-ISIS$  در هر دو مدل خطی و به ازای مقادیر مختلف  $\rho$  در انتخاب متغیر  $X_4$  بسیار تواناست. ما همچنین  $RDC-ISIS$  را با  $ISIS$ ، نسخه تکراری  $SIS$ ، مقایسه کرده‌ایم. نتایج نشان می‌دهد که در برخی موارد روش  $RDC-ISIS$  در انتخاب  $X_4$  حتی از  $ISIS$  نیز بهتر عمل می‌کند. به‌عنوان مثال، در مدل ۲ به ازای  $\rho = 0.8$  روش  $RDC-ISIS$  با احتمال تجربی ۹۱٪ همه متغیرهای مهم را به درستی انتخاب می‌کند، در حالیکه  $ISIS$  تنها با احتمال ۶۰٪ عمل انتخاب متغیر را به درستی انجام می‌دهد. به خاطر داشته باشید که  $RDC-ISIS$  از اطلاعات ساختاری مدل رگرسیونی استفاده نمی‌کند، در حالیکه  $ISIS$  برای غربالگری در مدل‌های خطی معرفی شده و از اطلاعات درست مدل خطی استفاده می‌کند. روش  $RDC-ISIS$  یک روش غربالگری آزاد مدل است و می‌تواند برای مدل‌های رگرسیونی مختلف جهت تشخیص روابط خطی و غیر خطی به کار رود.

در مثال زیر برخی از متغیرهای بی‌اهمیت، به دلیل همبستگی بالا با متغیر مهم  $X_1$ ، نسبت به برخی از متغیرهای مهم دیگر در فرآیند انتخاب برتری دارند. بنابراین لازم است از روش‌های تکراری برای انتخاب متغیرهای مهم استفاده شود. در این مثال عملکرد روش RDC-ISIS با روش‌های دیگر مقایسه می‌شود.

**مثال ۴.۴.۳.** در این مثال، مدل رگرسیونی جمعی ناپارامتری به صورت

$$Y = 2f_1(X_1) + \sqrt{6}f_2(X_{101}) + 3f_3(X_{201}) - 0.6f_4(X_{202}) + \varepsilon$$

است، که در آن توابع پایه‌ای به صورت

$$f_1(x) = \exp(2x/3) \quad f_2(x) = \frac{\sin(3\pi x/4 + 3/2)}{2 - \sin(3\pi x/4 + 3/2)}$$

$$f_3(x) = \begin{cases} x + 4 & x < -2 \\ |x| & |x| \leq 2 \\ 4 - x & x > 2 \end{cases} \quad f_4(x) = \log(x^2)$$

تعریف شده‌اند. با توجه به مدل، متغیرهای  $(X_1, X_{101}, X_{201}, X_{202})$  متغیرهای مؤثر بر پاسخ هستند. متغیرهای توضیحی از توزیع نرمال چند متغیره با میانگین صفر و همبستگی  $Cov(X_i, X_j) = \rho^{|i-j|}$  تولید شده‌اند. سپس متغیر  $X_k$  با متغیر جدید  $\zeta_k$   $X_k = 0.8X_1 + \zeta_k$  جایگزین می‌شود، که در آن  $\zeta_k \sim \mathcal{N}(0, 1)$  و  $k = 2, \dots, 100$ . بنابراین ۹۹ متغیر بی‌اهمیت  $(X_2, \dots, X_{100})$  دارای همبستگی بالایی با متغیر مهم  $X_1$  هستند. همچنین با در نظر گرفتن  $\rho = 0.5, 0.8$  دو نوع ماتریس کوواریانس برای متغیرهای توضیحی تعریف شده است. به منظور بررسی نیرومندی هر روش، با در نظر گرفتن دو نوع توزیع خطا، نتایج شبیه‌سازی در جدول ۵.۳ گزارش شده‌اند.

همان‌طور که ملاحظه می‌شود روش‌های غربالگری مستقل حاشیه‌ای SIS، NIS، SIRS و DC- در شناسایی متغیرهای مهم  $X_{101}, X_{201}, X_{202}$  ناتوان هستند، چون متغیرهای بی‌اهمیت  $(X_2, \dots, X_{100})$  دارای همبستگی بالایی با متغیر مهم  $X_1$  هستند و در فرآیند انتخاب متغیرها توسط این سه روش نسبت به سه متغیر مهم  $X_{101}, X_{201}, X_{202}$  اولویت دارند. بنابراین یک روش غربالگری تکراری مناسب برای حذف ۹۹ متغیر بی‌اهمیت و بازگرداندن سه متغیر مهم  $X_{101}, X_{201}, X_{202}$  به مدل نیاز است. با توجه به جدول ۵.۳ ملاحظه می‌شود که روش ISIS نمی‌تواند روابط غیرخطی بین متغیرهای توضیحی و متغیر پاسخ را شناسایی کند، در حالی که روش RDC-ISIS تمام متغیرهای مهم را در مدل جمعی ناپارامتری به درستی انتخاب می‌کند. بنابراین روش RDC-ISIS به دلیل آزاد-مدل بودن و نیرومند بودن یک روش قابل قبول برای مدل‌های غیرخطی با بعد خیلی بالا است.

مثال‌های بیان‌شده کارایی روش RDC-ISIS را برای غربالگری در داده‌های با بعد بسیار بالا نشان می‌دهند. در مثال زیر، پس از کاهش بعد اولیه با RDC-ISIS، از دو تابع جریمه SCAD به‌طور همزمان برای برازش مدل جمعی خطی-جزئی استفاده می‌شود.



جدول ۵.۳: احتمال تجربی  $P_j$  و احتمال  $S$  در مثال ۴

$\varepsilon \sim t(2)$					$\varepsilon \sim \mathcal{N}(0, 1)$					روش	$\rho$
$S$	$P_{2.2}$	$P_{2.1}$	$P_{1.1}$	$P_1$	$S$	$P_{2.2}$	$P_{2.1}$	$P_{1.1}$	$P_1$		
۰/۰۰	۰/۰۱	۰/۰۰	۰/۰۱	۰/۹۹	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۱/۰۰	SIS	۰/۵
۰/۰۰	۰/۰۵	۰/۰۵	۰/۰۶	۰/۹۸	۰/۰۰	۰/۰۵	۰/۰۶	۰/۰۷	۰/۹۸	ISIS	
۰/۳۲	۰/۶۱	۰/۷۱	۰/۷۵	۰/۹۸	۰/۵۰	۰/۷۴	۰/۸۳	۰/۸۲	۱/۰۰	NIS	
۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۱	۱/۰۰	SIRS	
۰/۰۳	۰/۲۷	۰/۳۴	۰/۳۱	۱/۰۰	۰/۰۲	۰/۲۴	۰/۲۹	۰/۲۵	۱/۰۰	DC-SIS	
۰/۱۱	۰/۳۵	۰/۳۶	۰/۳۰	۱/۰۰	۰/۱۱	۰/۳۳	۰/۴۰	۰/۳۴	۱/۰۰	RDC-SIS	
۰/۸۲	۰/۸۸	۰/۹۹	۰/۹۴	۱/۰۰	۰/۹۶	۰/۹۹	۰/۹۹	۰/۹۸	۱/۰۰	RDC-ISIS	
۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۰/۹۹	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۱/۰۰	SIS	۰/۸
۰/۰۰	۰/۰۵	۰/۰۵	۰/۰۹	۰/۹۷	۰/۰۰	۰/۰۴	۰/۰۴	۰/۰۷	۰/۹۳	ISIS	
۰/۵۹	۰/۸۸	۰/۹۳	۰/۷۲	۰/۹۹	۰/۸۱	۰/۹۶	۰/۹۹	۰/۸۵	۱/۰۰	NIS	
۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۱/۰۰	۰/۰۰	۰/۰۱	۰/۰۱	۰/۰۰	۱/۰۰	SIRS	
۰/۱۴	۰/۶۲	۰/۷۲	۰/۳۰	۱/۰۰	۰/۱۳	۰/۶۱	۰/۷۸	۰/۲۸	۱/۰۰	DC-SIS	
۰/۳۲	۰/۷۲	۰/۷۹	۰/۳۶	۱/۰۰	۰/۳۱	۰/۷۲	۰/۷۸	۰/۳۵	۱/۰۰	RDC-SIS	
۰/۸۳	۰/۹۵	۰/۹۸	۰/۹۱	۰/۹۸	۰/۹۰	۰/۹۳	۰/۹۸	۱/۰۰	۱/۰۰	RDC-ISIS	

مثال ۵.۴.۳. داده‌ها از مدل

$$Y = \sum_{j=1}^p f_j(X_j) + \varepsilon$$

تولید شده‌اند، که در آن  $f_3(x) = 3 \sin(2\pi x) / (2 - \sin(2\pi x))$ ،  $f_2(x) = 6x(1-x)$ ،  $f_1(x) = 2x$ ،  $f_4(x) = x$ ،  $f_5(x) = -x$  و به ازای  $j > 5$ ،  $f_j(x) = 0$ . بنابراین تعداد مؤلفه‌های ناپارامتری در این مدل ۲ و تعداد مؤلفه‌های پارامتری برابر ۳ است. متغیرهای توضیحی ابتدا همانند مثال ۲.۴.۳ تولید شده‌اند. سپس با استفاده از تابع توزیع تجمعی نرمال استاندارد به توزیع حاشیه‌ای یک‌نواخت روی بازه  $[0, 1]$  تبدیل شده‌اند. برای نشان دادن کارایی روش ارائه شده، برای خطا دو توزیع نرمال استاندارد و توزیع تی-استودنت با ۲ درجه آزادی که برای تولید خطای دم سنگین استفاده می‌شود، در نظر گرفته شده و شبیه‌سازی با  $n = 70, 200, 400$  و  $p = 1000, 2000$  به تعداد ۲۰۰ بار تکرار شده است. نتایج انتخاب متغیر و تشخیص ساختار مدل در جدول ۶.۳ خلاصه شده‌اند، که در آن  $N$  نشان دهنده متوسط تعداد مؤلفه‌های غیرصفر انتخاب شده،  $NN$  متوسط تعداد مؤلفه‌های ناپارامتری انتخاب شده،  $NNT$  متوسط تعداد مؤلفه‌های ناپارامتری انتخاب شده که واقعاً ناپارامتری هستند،  $NL$  متوسط تعداد مؤلفه‌های خطی انتخاب شده و  $NLT$  متوسط تعداد مؤلفه‌های خطی است که درست انتخاب شده‌اند.

برای بررسی دقت برآورد، ریشه دوم میانگین توان‌های دوم خطا (RMSE) برای پنج مؤلفه اول در

جدول ۶.۳: نتایج تشخیص ساختار مدل در مثال ۵

NLT(۳)*	NL(۳)*	NNT(۲)*	NN(۲)*	N(۵)*	خطا	n	p
۲/۶۹(۰/۶۳)	۲/۷۴(۰/۸۶)	۲(۰)	۲/۶۸(۱/۰۴)	۵/۴۲	$\mathcal{N}(۰, ۱)$	۷۰	
۲/۶۲(۰/۸۵)	۳/۶۵(۱/۱۵)	۱/۹۵(۰/۱۹)	۲/۹۶(۱/۱۰)	۶/۶۱	t(۲)		
۲/۷۴(۰/۳۲)	۲/۹۳(۰/۸۸)	۲(۰)	۲/۳۱(۰/۹۶)	۵/۲۴	$\mathcal{N}(۰, ۱)$	۲۰۰	۱۰۰۰
۲/۷۱(۰/۷۱)	۳/۳۷(۱/۰۳)	۱/۹۷(۰/۱۲)	۲/۵۲(۰/۸۸)	۵/۸۹	t(۲)		
۲/۹۶(۰/۲۶)	۳/۱۲(۰/۵۹)	۲(۰)	۲/۲۰(۰/۶۴)	۵/۳۲	$\mathcal{N}(۰, ۱)$	۴۰۰	
۲/۷۴(۰/۷۵)	۳/۲۸(۰/۶۶)	۱/۹۹(۰/۰۸)	۲/۴۹(۰/۵۴)	۵/۷۷	t(۲)		
۲/۸۵(۱/۰۵)	۳/۲۳(۱/۰۱)	۲(۰)	۲/۷۱(۱/۰۹)	۵/۹۴	$\mathcal{N}(۰, ۱)$	۷۰	
۲/۷۲(۰/۹۸)	۳/۳۸(۱/۴۸)	۱/۹۰(۰/۵۸)	۳/۲۶(۱/۳۱)	۶/۶۴	t(۲)		
۲/۸۴(۰/۶۴)	۳/۱۲(۰/۹۴)	۲(۰)	۲/۴۶(۰/۸۴)	۵/۵۸	$\mathcal{N}(۰, ۱)$	۲۰۰	۲۰۰۰
۲/۷۶(۰/۸۲)	۳/۲۴(۱/۰۱)	۱/۹۸(۰/۱۷)	۲/۷۷(۰/۷۹)	۶/۰۱	t(۲)		
۲/۹۲(۰/۶۴)	۳/۰۴(۰/۹۴)	۲(۰)	۲/۳۳(۰/۸۱)	۵/۳۷	$\mathcal{N}(۰, ۱)$	۴۰۰	
۲/۸۱(۰/۶۳)	۲/۹۷(۰/۸۵)	۱/۹۹(۰/۰۶)	۲/۵۴(۰/۷۴)	۵/۵۱	t(۲)		

\* اعداد داخل پرانتز، مقدار درست آن معیار در مدل واقعی است.

جدول ۷.۳ گزارش شده است و برای مؤلفه  $j$ ام بصورت

$$\text{RMSE}(f_j) = \sqrt{\frac{1}{500} \sum_{i=1}^{500} (\hat{f}_j(t_i) - f_j(t_i))^2}, \quad j = 1, \dots, 5$$

محاسبه می‌شود که  $(t_1, \dots, t_{500})$  نقاط شبکه‌ای، شامل نقطه با فواصل یکسان، روی بازه  $[0, 1]$  هستند و تابع  $f_j(\cdot)$  در این نقاط ارزیابی می‌شود. همچنین RMSE تابع رگرسیون  $f = \sum_{j=1}^p f_j$  نیز ارائه شده است. اعداد داخل پرانتز در جدول ۷.۳ خطای استاندارد متناظر هستند.

نتایج شبیه‌سازی نشان می‌دهد که روش دو مرحله‌ای ارائه شده در برآورد و تشخیص مؤلفه‌های غیرصفر و همچنین شناسایی مؤلفه‌های خطی و غیرخطی تواناست. برای خطای نرمال، این روش در هر ۲۰۰ تکرار هر دو مؤلفه ناپارامتری مدل را به درستی به‌عنوان مؤلفه غیرصفر ناپارامتری تشخیص می‌دهد، اما برای خطای غیرنرمال در برخی از تکرارها بعضی از مؤلفه‌های ناپارامتری به‌عنوان مؤلفه صفر یا خطی تشخیص داده می‌شوند. همچنین در برخی از تکرارها، مؤلفه‌های خطی به‌عنوان مؤلفه‌های صفر یا غیرخطی انتخاب می‌شوند. با توجه به جدول، هنگامی که  $n = 70$  و خطا دارای توزیع تی-استودنت است، این روش عملکرد چندان مطلوبی ندارد. در این حالت، مدل انتخاب شده معمولاً بزرگ‌تر از مدل واقعی است. به ویژه برای  $p = 2000$  این مشکل بسیار واضح است. اما با افزایش حجم نمونه عملکرد بهبود می‌یابد. علاوه بر این، دقت برآورد روش ارائه‌شده، تحت تاثیر توزیع خطاها قرار می‌گیرد. با توجه به جدول ۷.۳، در بیشتر حالت‌ها، مقدار RMSE متناظر با خطای نرمال کوچکتر از RMSE متناظر با خطای تی-استودنت است. این بدین دلیل است که در مرحله دوم برای انتخاب متغیر و

جدول ۷.۳: مجذور میانگین توان‌های دوم خطا برای پنج مؤلفه اول و تابع رگرسیون.

$f$	$f_5$	$f_4$	$f_3$	$f_2$	$f_1$	خطا	$n$
$p = 1000$							
۰.۳۷۷۳(۰/۱۸۶)	۰.۱۳۵۴(۰/۱۰۳)	۰.۱۴۵۶(۰/۱۰۲)	۰.۱۲۴۴(۰/۱۴۲)	۰.۲۱۵۴(۰/۱۳۵)	۰.۲۴۷۸(۰/۱۱۸)	$\mathcal{N}(0, 1)$	۷۰
۰.۴۷۲۰(۰/۳۲۰)	۰.۱۵۸۰(۰/۱۱۲)	۰.۱۵۹۷(۰/۱۱۳)	۰.۱۸۶۴(۰/۲۰۰۸)	۰.۲۵۶۵(۰/۱۵۶)	۰.۳۴۰۱(۰/۲۶۹)	$t(2)$	
۰.۳۸۶۷(۰/۰۷۴)	۰.۱۴۲۴(۰/۰۴۶)	۰.۱۴۳۷(۰/۰۴۵)	۰.۰۴۷۰(۰/۰۲۷)	۰.۲۰۳۴(۰/۰۸۱)	۰.۲۱۴۲(۰/۰۲۵)	$\mathcal{N}(0, 1)$	۲۰۰
۰.۴۰۲۴(۰/۰۶۵)	۰.۱۵۵۹(۰/۰۶۴)	۰.۱۶۲۸(۰/۰۶۱)	۰.۰۵۱۰(۰/۰۲۸)	۰.۲۸۰۰(۰/۰۶۴)	۰.۲۳۳۰(۰/۰۳۰)	$t(2)$	
۰.۲۴۶۱(۰/۰۷۹)	۰.۰۶۱۹(۰/۰۵۳)	۰.۰۶۱۲(۰/۰۵۵)	۰.۰۳۱۲(۰/۰۲۶)	۰.۱۶۷۱(۰/۰۸۰)	۰.۱۴۶۴(۰/۰۳۲)	$\mathcal{N}(0, 1)$	۴۰۰
۰.۲۷۴۸(۰/۰۹۱)	۰.۰۸۱۰(۰/۰۷۱)	۰.۱۰۰۱(۰/۰۷۵)	۰.۰۳۹۹(۰/۰۵۷)	۰.۱۷۵۰(۰/۰۶۷)	۰.۱۸۰۴(۰/۰۶۴)	$t(2)$	
$p = 2000$							
۰.۳۷۷۵(۰/۲۵۶)	۰.۱۳۲۳(۰/۱۰۵)	۰.۱۴۴۵(۰/۱۰۸)	۰.۰۹۴۹(۰/۱۰۷)	۰.۲۲۵۰(۰/۱۵۳)	۰.۲۷۱۸(۰/۲۱۵)	$\mathcal{N}(0, 1)$	۷۰
۰.۴۷۲۴(۰/۲۷۸)	۰.۱۴۴۹(۰/۱۰۸)	۰.۱۵۰۶(۰/۱۰۵)	۰.۰۹۵۲(۰/۱۳۴)	۰.۲۱۸۷(۰/۱۴۰)	۰.۳۰۵۷(۰/۲۴۶)	$t(2)$	
۰.۳۷۶۰(۰/۰۶۰)	۰.۱۳۳۶(۰/۰۳۷)	۰.۱۳۵۹(۰/۰۴۰)	۰.۰۴۷۶(۰/۰۲۴)	۰.۲۸۲۸(۰/۰۷۱)	۰.۲۱۵۹(۰/۰۲۶)	$\mathcal{N}(0, 1)$	۲۰۰
۰.۴۳۳۱(۰/۰۸۴)	۰.۱۴۸۹(۰/۰۶۴)	۰.۱۸۴۷(۰/۰۶۵)	۰.۰۶۳۴(۰/۰۴۳)	۰.۲۳۲۱(۰/۰۸۳)	۰.۲۴۵۳(۰/۰۴۲)	$t(2)$	
۰.۲۳۸۹(۰/۰۶۷)	۰.۰۶۵۴(۰/۰۵۲)	۰.۰۵۴۸(۰/۰۴۳)	۰.۰۲۸۳(۰/۰۲۴)	۰.۱۵۹۹(۰/۰۶۵)	۰.۱۵۲۸(۰/۰۴۵)	$\mathcal{N}(0, 1)$	۴۰۰
۰.۲۷۶۳(۰/۰۸۴)	۰.۰۸۲۷(۰/۰۵۵)	۰.۰۶۹۸(۰/۰۵۸)	۰.۰۳۶۱(۰/۰۲۶)	۰.۱۸۴۹(۰/۰۸۰)	۰.۱۶۸۵(۰/۰۴۹)	$t(2)$	

تشخیص ساختار از تابع زیان توان دوم استفاده شده است. این تابع زیان نسبت به داده‌های پرت و توزیع دم‌سنگین خطاها بسیار حساس است.

### ۵.۳ تحلیل داده‌های مسمومیت کبد

در این بخش، با استفاده از روش دو مرحله‌ای پیشنهادی، به تحلیل یک مجموعه داده واقعی می‌پردازیم. این مجموعه داده در بسته "mixOmics" نرم افزار R قابل دسترسی است و نخستین بار توسط بشل<sup>۱۳</sup> و همکاران (۲۰۰۷) مورد تحلیل قرار گرفته است. این داده‌ها مربوط به بیماری مسمومیت کبد بوده و شامل سطوح بیان ۳۱۱۶ ژن و ۹ اندازه‌گیری بالینی برای ۶۴ موش است. بنابراین، ۳۱۱۶ متغیر توضیحی و ۹ متغیر پاسخ کمی وجود دارند. متغیرهای پاسخ به اختصار عبارتند از: BUN، نیتروژن اوره؛ TP، پروتئین کل؛ ALB، آلبومین؛ ALT، آلانین آمینوترانسفراز؛ SDH، سوربیتول دهیدروژناز؛ AST، آسپارات آمینوترانسفراز؛ ALP، آلکالین فسفاتاز؛ TBA، کل اسیدهای صفرای؛ و CHOL، کلسترول. در این بخش ALT به‌عنوان متغیر پاسخ است و هدف پیدا کردن مؤثرترین ژن‌ها در پیش‌بینی ALT و نوع تاثیر این ژن‌هاست.

جدول ۸.۳ نشان دهنده ژن‌هایی است که توسط روش‌های غربالگری مختلف در رتبه‌های اول تا چهارم قرار می‌گیرند. به‌عنوان مثال، RDC-SIS و SIRS رتبه اول را به ژن  $A_{42\_P496622}$  اختصاص می‌دهند. در مقابل، DC-SIS ژن  $A_{43\_P14131}$  و NIS ژن  $A_{42\_P825290}$  را به‌عنوان مهمترین ژن انتخاب می‌کند.

<sup>13</sup>Bushel

جدول ۸.۳: چهار ژن بسیار مهم انتخاب شده توسط روش‌های مختلف

ژن‌های انتخاب شده				روش
$A_{42\_P469551}$	$A_{43\_P20438}$	$A_{43\_P14163}$	$A_{42\_P496622}$	RDC-SIS
$A_{42\_P469551}$	$A_{43\_P12724}$	$A_{43\_P20438}$	$A_{42\_P496622}$	SIRS
$A_{43\_P111724}$	$A_{42\_P840776}$	$A_{42\_P620915}$	$A_{43\_P14131}$	DC-SIS
$A_{43\_P14131}$	$A_{43\_P111724}$	$A_{43\_P10606}$	$A_{42\_P825290}$	NIS

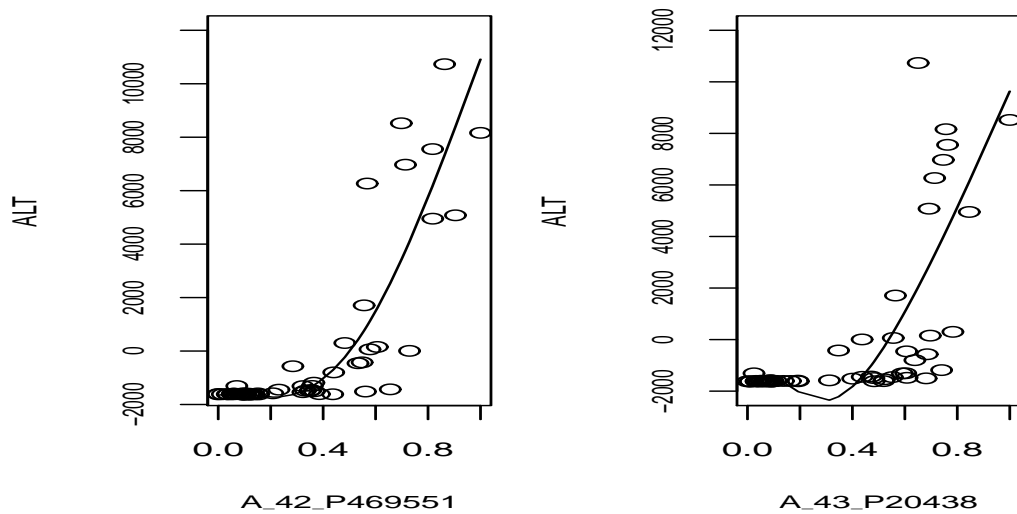
جدول ۹.۳: ژن‌های مهم انتخاب شده و نوع تأثیر آن‌ها

ژن				نوع تأثیر
$A_{43\_P14163}$	$A_{43\_P14864}$	$A_{42\_P694105}$	$A_{43\_P12724}$	غیرخطی
	$A_{43\_P20438}$	$A_{42\_P469551}$		
$A_{42\_P677628}$	$A_{42\_P496622}$	$A_{43\_P20962}$	$A_{42\_P619288}$	خطی

برای به‌کارگیری روش دو مرحله‌ای پیشنهادی در تحلیل داده‌های مسمومیت کبد، ابتدا نسخه تکراری RDC-SIS، یعنی RDC-ISIS با  $d = 28$  و  $d_1 = d/2$  را به کار برده و بعد داده‌ها را به  $d = 28$  کاهش داده‌ایم. دلیل استفاده از RDC-ISIS این است که اگر متغیر مهمی توسط RDC-SIS شناسایی نشود، آن متغیر مهم مجدداً برای ورود به مدل کاندید می‌شود و اهمیت آن مورد بررسی قرار می‌گیرد. علاوه بر این، نمودارهای پراکنش و منحنی‌های برازش اسپلاین مکعبی متغیرهای  $A_{42\_P469551}$  و  $A_{43\_P20438}$  را در مقابل متغیر پاسخ Ro1 در شکل ۱.۳ نشان داده‌ایم. این نمودار وجود یک الگوی غیرخطی بین متغیرهای توضیحی و متغیر پاسخ را نشان می‌دهد. پس از کاهش بعد اولیه، از دو تابع جریمه SCAD برای انتخاب متغیر و تشخیص ساختار مدل استفاده شده است. بعد از به‌کارگیری این روش دو مرحله‌ای، ۶ ژن دارای اثر غیرخطی و ۴ ژن دارای اثر خطی تشخیص داده شدند. ژن‌های انتخاب شده و نوع تأثیر آن‌ها در جدول ۹.۳ گزارش شده است. برای ارزیابی روش دو مرحله‌ای ارائه شده، از اعتبارسنجی متقابل استفاده می‌کنیم. ابتدا یک مشاهده را به‌عنوان داده آزمون و سایر مشاهدات را به‌عنوان داده‌های آموزشی در نظر بگیرید. پس از برازش مدل با استفاده از داده‌های آموزشی، مقدار متغیر پاسخ به ازای داده آزمون پیش‌بینی می‌شود. این روند برای تمام مشاهدات تکرار شده و نمودار مقادیر مشاهده شده متغیر پاسخ در مقابل مقادیر پیش‌بینی شده در نمودار ۲.۳ رسم شده است. با توجه به نمودار، روش معرفی شده در پیش‌بینی مقادیر ALT عملکرد خوبی دارد.

## بحث و نتیجه‌گیری

در این فصل، یک روش دو مرحله‌ای برای انتخاب متغیر و تشخیص ساختار در مدل‌های جمعی با بعد بسیار بالا ارائه شد. در مرحله اول، از یک روش غربالگری مستقل مطمئن برای کاهش بعد مدل استفاده



شکل ۱.۳: نمودارهای پراکنش متغیرپاسخ در مقابل دو ژن مهم و منحنی‌های اسپلاین مکعبی.

شد که در آن متغیرهای توضیحی براساس میزان همبستگی فاصله‌ای آنها با تابع توزیع حاشیه‌ای متغیر پاسخ رتبه‌بندی می‌شوند. کارایی این روش غربالگری با مطالعه شبیه‌سازی و تحلیل یک مجموعه داده واقعی مورد ارزیابی قرار گرفت که نتایج حاکی از عملکرد مطلوب روش ارائه شده است. در داده‌های با بعد بسیار بالا، معمولاً برخی از متغیرهای توضیحی دارای اثر خطی و برخی دیگر دارای اثر غیر خطی هستند. بنابراین تشخیص بخش‌های پارامتری و ناپارامتری بسیار حائز اهمیت است. لذا در مرحله دوم، برای انتخاب مدل نهایی و تشخیص مؤلفه‌های خطی و غیرخطی از دو تابع جریمه به‌طور همزمان استفاده شد. در اینجا نیز مطالعات شبیه‌سازی نشان دادند که این روش دو مرحله‌ای در برازش مدل جمعی خطی-جزیی کاراست.

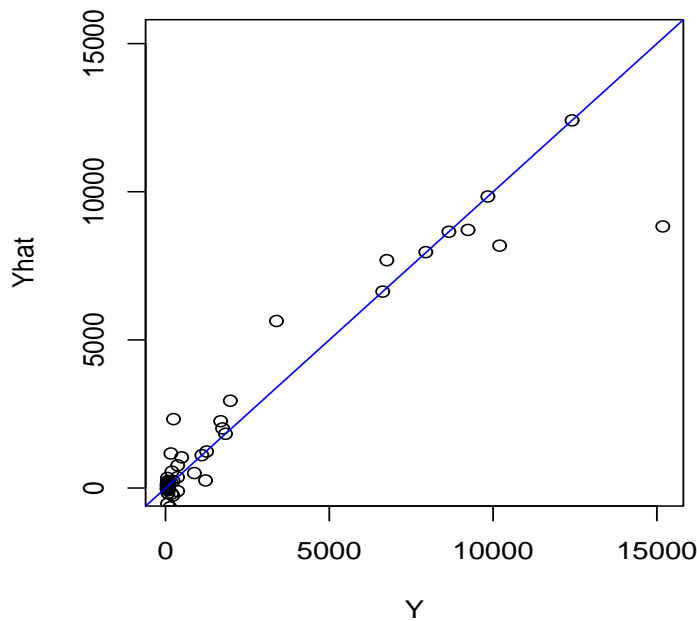
اگر چه در این فصل، مدل جمعی مد نظر بوده است، اما روش ارائه شده برای غربالگری و تشخیص ساختار را می‌توان به مدل‌های دیگر مانند مدل ضریب متغیر و مدل جمعی تعمیم یافته تعمیم داد. همچنین لازم به ذکر است که در این فصل، متغیرهای توضیحی و پاسخ صرفاً از نوع کمی در نظر گرفته شده‌اند. موضوع غربالگری در حالتی که این متغیرها از نوع کیفی و چند سطحی است، می‌تواند موضوعی برای تحقیقات آینده در نظر گرفته شود.

به‌عنوان روش دیگر برای غربالگری در مدل جمعی ناپارامتری، می‌توان از کوواریانس توابع ناپارامتری و تابع توزیع متغیر پاسخ استفاده کرد. بطور واضح‌تر، فرض کنید  $m_j(X_j) = E(Y|X_j)$  تصویر  $Y$  بر  $X_j$  و  $G(y) = P(Y \leq y)$  تابع توزیع غیر شرطی  $Y$  باشد. سپس از معیار مطلوبیت حاشیه‌ای

$$\omega_j = [Cov(m_j(X_j), G(Y))]^2, \quad 1 \leq j \leq p \quad (11.3)$$

برای رتبه‌بندی متغیرهای توضیحی استفاده می‌کنیم. برای پیاده‌سازی این روش غربالگری کافی است توابع ناپارامتری  $m_j(X_j)$  را با توابع  $B$ -اسپلاین و تابع توزیع  $G(y)$  را با تابع توزیع تجربی برآورد کنیم.

در خصوص تقلیل بعد بالای فضای متغیرهای توضیحی به بعد مرتبه  $d$  توسط روش غربالگری به‌کار



شکل ۲.۳: نمودار مقادیر مشاهده شده متغیر پاسخ در مقابل مقادیر پیش بینی شده

رفته در این فصل، می‌توان گفت که انتخاب  $d$ ، همانند انتخاب پارامتر جریمه در روش‌های انقباضی، دارای اهمیت بسزائی است. تاکنون اغلب روش‌های معرفی شده برای تشخیص مقدار  $d$  وابسته به ساختار مدل هستند، بنابراین معرفی روشی جدید برای انتخاب  $d$  در روش‌های غربالگری آزاد-مدل بسیار حائز اهمیت است.

نوع تابع جریمه در تشخیص ساختار مدل، مسئله دیگری است که بایستی به آن توجه شود. اگر چه ما از تابع جریمه SCAD استفاده کرده‌ایم، اما می‌توان از ترکیب روش غربالگری با سایر توابع جریمه نظیر Lasso، MCP یا aLasso نیز استفاده کرد.

## فصل ۴

# انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی - جزئی

در این فصل، با استفاده از رگرسیون نما، یک روش جدید برای انتخاب متغیر و تشخیص اثرات ثابت و متغیر در مدل ضریب متغیر ارائه می‌شود. این روش نسبت به خطاهای غیر نرمال و داده‌های پرت در پاسخ نیرومند است. از این فصل، مقاله کاظمی و همکاران (۲۰۱۹b) استخراج شده است.

### ۱.۴ مقدمه

مدل نیمه پارامتری ضریب متغیر خطی - جزئی، تعمیمی از مدل ضریب متغیر و مدل خطی - جزئی است که در آن برخی از ضرایب رگرسیونی، ثابت و برخی دیگر تابعی از یک متغیر شاخص، مانند سن یا زمان، هستند. این مدل علاوه بر بهره‌گیری از انعطاف‌پذیری مدل‌های ناپارامتری و تفسیرپذیری مدل‌های پارامتری، اثرات متقابل بین متغیر شاخص و متغیرهای توضیحی را نیز به صورت ناپارامتری در نظر می‌گیرد. تاکنون بسیاری از محققین به مسئله برآورد در مدل‌های ضریب متغیر خطی - جزئی پرداخته‌اند. به عنوان مثال، می‌توان به احمد<sup>۱</sup> و همکاران (۲۰۰۵)، فن و هوانگ (۲۰۰۵)، ژو و لیانگ<sup>۲</sup>

<sup>۱</sup>Ahmad

<sup>۲</sup>Zhou and Liang

(۲۰۰۹)، وانگ و همکاران (۲۰۰۹)، سان و لین<sup>۳</sup> (۲۰۱۳) و شن<sup>۴</sup> و همکاران (۲۰۱۴) اشاره کرد. در همه مطالعات مذکور، فرض بر این است که زیرمجموعه متغیرهایی که دارای اثر ثابت یا اثر متغیر بر پاسخ هستند، از پیش تعیین شده‌اند. این یک فرض زیربنایی است که در ساخت برآوردگرها و بررسی ویژگی‌های مجانبی آنها اهمیت بسزایی دارد. با این وجود، در عمل، تشخیص اینکه متغیری دارای اثر ثابت یا متغیر بر پاسخ است، بسیار حائز اهمیت است. علاوه بر این، اغلب تعداد زیادی از متغیرهای توضیحی در مدل حضور دارند. با توجه به ضرورت حذف متغیرهای بی‌اهمیت از مدل، لازم است که در مدل‌های نیمه‌پارامتری ضریب متغیر، انتخاب متغیر و تشخیص ساختار به‌طور همزمان انجام شود. تاکنون تعداد معدودی از محققین به موضوع انتخاب متغیر در مدل ضریب متغیر خطی- جزئی پرداخته‌اند. از جمله این کارها می‌توان به وانگ و همکاران (۲۰۰۸)، وانگ و زیا (۲۰۰۹)، ژائو و زوئی<sup>۵</sup> (۲۰۰۹)، ژائو و زوئی (۲۰۱۰)، کی<sup>۶</sup> و همکاران (۲۰۱۱)، ژائو و همکاران (۲۰۱۳)، ژانگ و همکاران (۲۰۱۳) اشاره کرد. اما تاکنون انتخاب متغیر به‌طور همزمان با تشخیص ساختار مورد بررسی قرار نگرفته است. لذا ارائه یک روش جدید که بتواند انتخاب متغیرهای مهم و تشخیص ضرایب ثابت و متغیر را بطور همزمان انجام دهد، مورد توجه محققین است که در این فصل به آن می‌پردازیم. زیا و همکاران (۲۰۰۴)، لنگ (۲۰۰۹)، نه و کیلگوم<sup>۷</sup> (۲۰۱۲)، هو<sup>۸</sup> و زیا (۲۰۱۲) روش‌های مختلفی را برای تشخیص اثرات ثابت در این مدل معرفی کردند، اما این روش‌ها متغیرهای مهم را انتخاب نمی‌کنند. همچنین، تانگ<sup>۹</sup> و همکاران (۲۰۱۲) با استفاده از رگرسیون‌های چندکی و میانگین، یک روش واحد برای مدل‌های ضریب متغیر خطی- جزئی معرفی کردند که علاوه بر انتخاب متغیرهای مهم، اثرات ثابت یا متغیر را نیز شناسایی می‌کند، اما این روش مستلزم استفاده از یک الگوریتم دو مرحله‌ای تکراری است که در عمل چندان کارا نیست، زیرا علاوه بر هزینه محاسباتی زیاد، هیچ تضمینی برای همگرایی این الگوریتم به جواب درست وجود ندارد. همچنین روش معرفی شده بر اساس رگرسیون میانگین نسبت به مشاهدات پرت یا توزیع دم‌سنگین خطاها بسیار حساس است. علاوه بر این رگرسیون چندکی نیز در صورت نرمال بودن توزیع خطاها و عدم وجود مشاهدات پرت، ممکن است از کارایی لازم برخوردار نباشد. بنابراین ارائه یک روش نیرومند برای انتخاب متغیر و تشخیص ساختار در مدل‌های ضریب متغیر بسیار حائز اهمیت است.

اخیراً، یائو<sup>۱۰</sup> و همکاران (۲۰۱۲) و یائو و لی (۲۰۱۴)، براساس نمای شرطی متغیر پاسخ  $Y$ ، یک روش نیرومند به نام "رگرسیون نما"<sup>۱۱</sup> را معرفی کردند. این روش در تحلیل داده‌های چوله کاربرد زیادی دارد. در این فصل، یک روش نیرومند برای انتخاب متغیر و تشخیص ساختار در مدل‌های ضریب متغیر خطی- جزئی براساس رگرسیون نما ارائه می‌شود. ویژگی‌های مجانبی برآوردگرها مورد بررسی قرار

<sup>3</sup>Sun and Lin

<sup>4</sup>Shen

<sup>5</sup>Zhao and Zue

<sup>6</sup>Kai

<sup>7</sup>Noh and Keilegom

<sup>8</sup>Hu

<sup>9</sup>Tang

<sup>10</sup>Yao

<sup>11</sup>Modal regression



گرفته و عملکرد این روش با مطالعات شبیه‌سازی و تحلیل داده‌های واقعی ارزیابی می‌شود.

## ۱.۱.۴ رگرسیون نما

رگرسیون نما، نمای متغیر پاسخ  $Y$  را، به شرط بردار متغیرهای توضیحی  $\mathbf{X} = (X_1, \dots, X_p)^T$ ، به صورت یک تابع خطی از  $\mathbf{X}$  مدل‌سازی می‌کند. این مدل رگرسیونی، متفاوت از رگرسیون خطی معمولی است، که فرض می‌شود میانگین شرطی  $Y$  تابعی خطی از  $\mathbf{X}$  است. رگرسیون نما، به جای میانگین شرطی، از "محتمل‌ترین" مقادیر شرطی برای سنجش مرکز توزیع متغیر پاسخ استفاده می‌کند. بنابراین، در مقایسه با رگرسیون میانگین، رگرسیون نما دارای چندین مزیت است (یائو و همکاران، ۲۰۱۴). برخی از مزیت‌ها عبارتند از:

۱. به دلیل نیرومندی نما، رگرسیون نما نسبت به مشاهدات پرت و توزیع دم‌سنگین خطا نیرومند است.
۲. هنگامی که چولگی توزیع خطا بسیار زیاد است، رگرسیون نما پیشگویی نقطه‌ای معنی‌دارتری از رگرسیون میانگین نتیجه می‌دهد.
۳. با توجه به تمرکز رگرسیون نما بر ناحیه با بزرگترین چگالی شرطی، انتظار می‌رود که یک فاصله کوتاه حول برآورد رگرسیون نما نسبت به یک فاصله با همان طول حول برآورد رگرسیون میانگین، دارای احتمال پوشش بزرگتری باشد.

بنابراین رگرسیون نما یک ابزار مفید برای تحلیل داده‌های چوله است. برای برآورد رگرسیون نما یک تابع هدف هسته مبنا به کار می‌رود. برای ماکزیمم کردن این تابع هدف از الگوریتم MEM استفاده می‌شود. فرض کنید چگالی شرطی  $f(y|\mathbf{x})$  به شرط  $\mathbf{X}$  باشد و نمای تابع  $f(y|\mathbf{x})$ ، که به صورت  $\text{Mode}(Y|\mathbf{X}) = \arg \max_y (f(y|\mathbf{x}))$  تعریف می‌شود، یکتاست. در رگرسیون نما، فرض می‌شود  $\text{Mode}(Y|\mathbf{X})$  تابعی خطی از  $\mathbf{X}$  است، یعنی

$$\text{Mode}(Y|\mathbf{X}) = \mathbf{X}^T \boldsymbol{\beta} \quad (1.4)$$

که در آن  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  بردار ضرایب رگرسیونی است. فرض کنید  $\varepsilon = Y - \mathbf{X}^T \boldsymbol{\beta}$  خطای تصادفی و تابع  $g(\varepsilon|\mathbf{X})$  چگالی شرطی  $\varepsilon$  به شرط  $\mathbf{X}$  باشد. اگر  $g(\varepsilon|\mathbf{X})$  حول صفر متقارن باشد، بردار ضرایب  $\boldsymbol{\beta}$  در (۱.۴) همان ضرایب رگرسیون خطی میانگین است. اما اگر  $g(\varepsilon|\mathbf{X})$  چوله باشد، این دو بردار ضرایب متفاوت خواهند بود.

برای درک بهتر تفاوت بین تابع رگرسیون نما و تابع رگرسیون میانگین، مثال زیر را در نظر بگیرید.

مثال ۱.۱.۴. مدل زیر را در نظر بگیرید

$$Y = m(\mathbf{X}) + \sigma(\mathbf{X})\varepsilon \quad (2.4)$$

که  $\varepsilon$  دارای تابع چگالی چوله  $h(\cdot)$  با میانگین صفر و مد ۱ است. اکنون دو حالت زیر را در نظر بگیرید.

۱. اگر  $m(\mathbf{X}) = \mathbf{X}^T \boldsymbol{\beta}$  و  $\sigma(\mathbf{X}) = \alpha$ ، آنگاه

$$E(Y|\mathbf{X}) = \mathbf{X}^T \boldsymbol{\beta}, \quad \text{Mode}(Y|\mathbf{X}) = \mathbf{X}^T (\alpha + \boldsymbol{\beta}).$$

بنابراین، در هر دو رگرسیون میانگین و نما،  $Y$  بطور خطی به  $\mathbf{X}$  وابسته است، اگرچه دو رگرسیون دارای ضرایب رگرسیونی متفاوت هستند.

۲. اگر  $m(\mathbf{X}) = \alpha$  و  $\sigma(\mathbf{X}) = \boldsymbol{\beta}$ ، پس

$$E(Y|\mathbf{X}) = \alpha, \quad \text{Mode}(Y|\mathbf{X}) = \mathbf{x}^T \boldsymbol{\beta}.$$

بنابراین، از منظر رگرسیون میانگین،  $Y$  وابسته به  $\mathbf{X}$  نیست، اما براساس رگرسیون نما،  $Y$  دارای رابطه خطی با  $\mathbf{x}$  است. از این مثال می‌توان نتیجه گرفت که روش‌های انتخاب متغیر براساس رگرسیون نما، در مقایسه با رگرسیون میانگین، ممکن است در شناسایی برخی متغیرهای توضیحی مهم مفیدتر باشند.

فرض کنید  $\{(\mathbf{X}_i, Y_i), i = 1, \dots, n\}$  یک نمونه تصادفی از  $(\mathbf{X}, Y)$  باشد. برای برآزش رگرسیون نما، بردار ضرایب  $\boldsymbol{\beta}$  در (۱.۴) با ماکزیمم کردن

$$Q_h(\boldsymbol{\beta}) = \sum_{i=1}^n \phi_h(Y_i - \mathbf{X}_i^T \boldsymbol{\beta}) \quad (۳.۴)$$

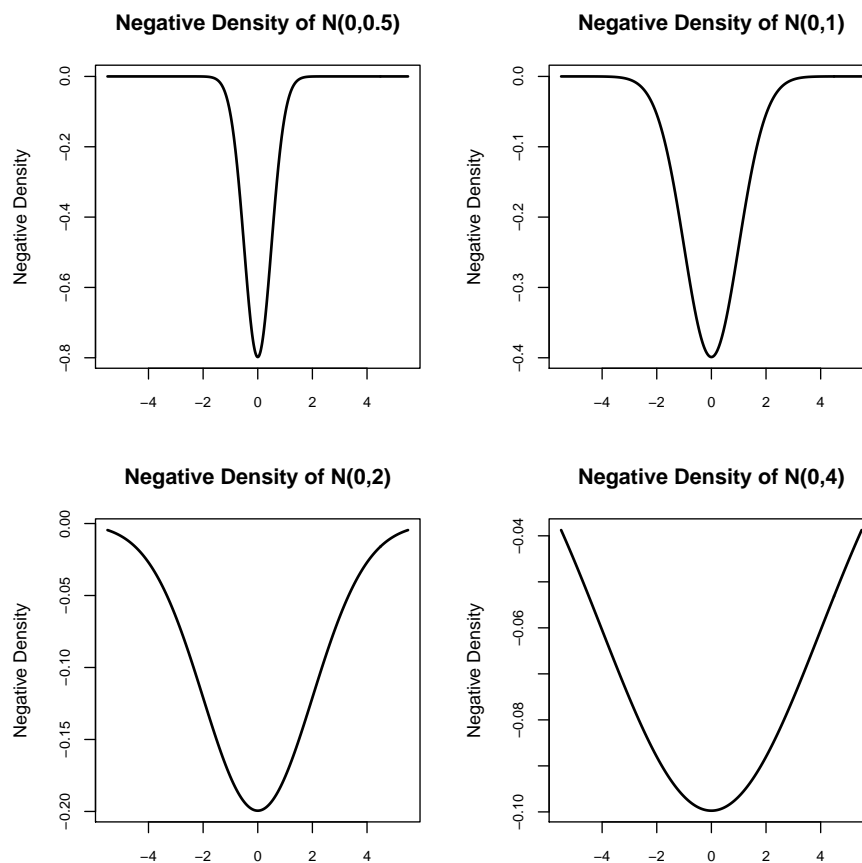
برآورد می‌شود که در آن  $\phi_h(t) = \frac{1}{h} \phi(t/h)$ ،  $\phi(\cdot)$  یک تابع چگالی هسته متقارن حول صفر و  $h$  یک پارامتر هموارساز است که پهنای باند نامیده می‌شود.

به منظور درک علت استفاده از (۳.۴) برای برآورد رگرسیون نما، مدل تنها شامل عرض از مبدأ را در نظر بگیرید. در این مدل، هیچ متغیر توضیحی در آن حضور ندارد، به عبارتی  $\boldsymbol{\beta}_0 = \beta_0$ . در این صورت تابع هدف عبارتست از

$$Q_h(\beta_0) = \sum_{i=1}^n \phi_h(Y_i - \beta_0). \quad (۴.۴)$$

$Q_h(\beta_0)$ ، به عنوان یک تابع از  $\beta_0$ ، برآورد هسته تابع چگالی  $Y$  است. بنابراین ماکزیمم‌کننده (۴.۴)، نمای تابع چگالی هسته براساس  $(Y_n, \dots, Y_1)$  است. هنگامی که  $n \rightarrow \infty$  و  $h \rightarrow 0$ ، نمای تابع چگالی هسته به نمای توزیع  $Y$  همگرا می‌شود. این برآوردگر نما توسط پارزن<sup>۱۲</sup> (۱۹۶۲) معرفی شد. هنگامی که متغیرهای توضیحی در مدل حضور دارند، برای هر  $\boldsymbol{\beta}$  ثابت، در (۳.۴) مقدار تابع چگالی هسته براساس باقیمانده‌های  $\varepsilon_i = Y_i - \mathbf{X}_i^T \boldsymbol{\beta}$  در  $\varepsilon = 0$  است. ماکزیمم کردن (۳.۴) نسبت به  $\boldsymbol{\beta}$  خط  $\mathbf{X}^T \hat{\boldsymbol{\beta}}$  را نتیجه می‌دهد به طوری که تابع چگالی هسته باقیمانده‌های  $\varepsilon_i$  دارای بیشترین مقدار در  $\varepsilon = 0$  است. در حالت خاص، اگر هسته یکنواخت  $\phi_h(t) = \frac{1}{2h} I(|t| \leq h)$  استفاده شود، ماکزیمم کردن (۳.۴) خط  $\mathbf{X}^T \hat{\boldsymbol{\beta}}$  را نتیجه می‌دهد به طوری که باند  $\mathbf{X}^T \hat{\boldsymbol{\beta}} \pm h$  بیشترین تعداد از مقادیر پاسخ  $Y_i$  را شامل می‌شود.

<sup>12</sup>Parzen



شکل ۱.۴: نمودار  $-\phi_h(\cdot)$  به ازای پارامترهای مختلف.

یائو و لی (۲۰۱۴) ثابت کردند اگر وقتی که  $n \rightarrow \infty$  داشته باشیم  $h \rightarrow 0$ ، آنگاه برای هر تابع چگالی خطا (بدون فرض متقارن بودن توزیع خطا)،  $\hat{\beta}$  به دست آمده با ماکزیمم کردن تابع هدف (۳.۴)، یک برآوردگر سازگار از بردار ضرایب  $\beta$  در (۱.۴) است.

در حالت کلی، نمای نمونه‌ای، به عنوان برآوردگری از نمای جامعه، نسبت به داده‌های پرت حساس نیست (یائو و همکاران، ۲۰۱۲). علاوه بر این، اگر  $-\phi_h(\cdot)$  را به عنوان یک تابع زیان در نظر بگیریم، رگرسیون نما یک  $M$ -برآوردگر<sup>۱۳</sup> را نتیجه می‌دهد که برآوردگری نیرومند است. پهنای باند  $h$  میزان نیرومندی برآوردگر را تعیین می‌کند. فرض کنید  $\phi(\cdot)$  چگالی نرمال استاندارد باشد. شکل ۱.۴ چگونگی نیرومندی برآوردگر رگرسیون نما را نشان می‌دهد. توجه کنید که  $h^2$  متناظر با واریانس چگالی نرمال است. از شکل ۱.۴ می‌توان دید که  $-\phi_h(\cdot)$  برای  $h$  کوچک، مانند  $h = 0.5$ ، به صورت یک تابع زیان پرت-مقاوم<sup>۱۴</sup> به نظر می‌رسد، در حالی که شکل این تابع برای  $h$  بزرگ، مانند  $h = 4$ ، شبیه به تابع زیان توان دوم است. در عمل، مقدار بهینه  $h$  با یک روش مناسب مبتنی بر داده تعیین می‌شود که در ادامه به آن خواهیم پرداخت.

<sup>13</sup>M-estimator

<sup>14</sup>Outlier-resistant

## ۲.۴ روش نیرومند انتخاب متغیر و تشخیص ساختار

فرض کنید مجموعه داده‌های مشاهده شده به صورت

$$D_n = \{(Y_i, X_i, U_i), i = 1, \dots, n\}$$

باشد که  $Y_i$  پاسخ  $i$  امین مشاهده،  $X_i \in \mathbb{R}^p$  بردار متغیرهای توضیحی  $i$  امین مشاهده و  $U_i$  متغیر شاخص که، بدون از دست دادن کلیت، فرض می‌شود دارای توزیع یکنواخت  $[0, 1]$  است. لذا مدل ضریب متغیر خطی-جزئی با ساختار درست به صورت زیر بیان می‌شود

$$Y_i = \sum_{j \in I_V} \theta_j(U_i) X_{ij} + \sum_{j \in I_C} \beta_j X_{ij} + \sum_{j \in I_Z} \alpha_j(U_i) X_{ij} + \varepsilon_i \quad (5.4)$$

که  $\{\theta_j(\cdot), j \in I_V\}$  توابع هموار نامعلوم،  $\{\beta_j, j \in I_C\}$  ضرایب رگرسیونی نامعلوم و  $\alpha_j(\cdot)$  نشان‌دهنده ضرایب صفر است. همچنین مجموعه‌های  $I_V, I_C$  و  $I_Z$  مجموعه‌های اندیس، به ترتیب، برای  $X$  های با اثرات متغیر، اثرات ثابت غیر صفر و اثرات صفر هستند. این مجموعه‌ها دو به دو مجزا و در شرط  $I_V \cup I_C \cup I_Z = \{1, \dots, p\}$  صدق می‌کنند.

نخستین گام برای استفاده از مدل (۵.۴)، تعیین مجموعه‌های  $I_V, I_C$  و  $I_Z$  می‌باشد. شناسایی متغیرهای توضیحی با اثر ثابت و اثر متغیر در عمل بسیار حائز اهمیت است و لذا در این فصل، هدف ما تشخیص این مجموعه‌ها و همچنین برآورد ضرایب متغیر هموار  $\{\theta_j(U), j \in I_V\}$  و ضرایب ثابت  $\{\beta_j, j \in I_C\}$  بطور کارا و نیرومند است.

در ابتدای تحلیل، چون ساختار خطی-جزئی نامعلوم است، مدل (۵.۴) را به صورت مدل ضریب متغیر زیر در نظر می‌گیریم

$$Y_i = \sum_{j=1}^p \alpha_j(U_i) X_{ij} + \varepsilon_i, \quad i = 1, \dots, n. \quad (6.4)$$

حال اگر، برای هر  $U \in [0, 1]$ ، داشته باشیم  $\alpha_j(U) \equiv 0$ ، آنگاه  $X_j$  یک متغیر بی اهمیت یا نامرتبط است. اگر  $\alpha_j(U)$  یک مقدار ثابت غیر صفر باشد،  $X_j$  دارای اثر ثابت و در غیر این صورت،  $\alpha_j(U)$  یک تابع متغیر و  $X_j$  دارای اثر متغیر است. بنابراین، می‌خواهیم تعیین کنیم که کدام یک از ضرایب  $\alpha_j(\cdot)$  توابعی متغیر، صفر یا ثابت غیر صفر هستند.

بدین منظور، ابتدا  $\alpha_j(U)$  را به صورت زیر به دو جزء ثابت و متغیر تجزیه می‌کنیم

$$\alpha_j(U) = \eta_j + g_j(U), \quad j = 1, \dots, p \quad (7.4)$$

که برای تضمین شناسایی‌پذیری<sup>۱۵</sup> فرض می‌شود  $E\{g_j(U)\} = 0$ . به وضوح تجزیه (۷.۴) یکتاست و

$$\eta_j = E\{\alpha_j(U)\}, \quad g_j(U) = \alpha_j(U) - E\{\alpha_j(U)\}.$$

<sup>15</sup>Identifiability

اکنون می‌توانیم از اسپلاین‌های چند جمله‌ای برای تقریب توابع ناپارامتری  $g_j(\cdot)$  استفاده کنیم. فرض کنید  $0 = \tau_0 < \tau_1 < \dots < \tau_{K_n} < \tau_{K_n+1} = 1$  افزایشی از بازه  $[0, 1]$  به  $K_n + 1$  زیر بازه است. مجموعه توابع اسپلاین درجه  $q$ ، دارای پایه‌های  $B$  - اسپلاین  $\mathbf{B}(\cdot) = (B_1(\cdot), \dots, B_{K'}(\cdot))^T$  می‌باشند که  $K' = K_n + q + 1$ . با در نظر گرفتن شرط محدودکننده  $E\{g_j(U)\} = 0$ ، تنها فضای توابع اسپلاین  $S_j^0 = \{s : s = \sum_{k=1}^{K'} \gamma_{jk} B_k(x), \sum_{i=1}^n s(X_{ij}) = 0\}$  با پایه‌های مرکزی شده  $\{B_{jk}(x) = B_k(x) - \sum_{i=1}^n B_k(X_{ij})/n, k = 1, 2, \dots, K\}$  را در نظر می‌گیریم که  $K = K' - 1$ . لذا تابع  $g_j(U)$  را به صورت

$$g_j(U) \approx \sum_{k=1}^K B_{jk}(U) \gamma_{j,k} = \mathbf{B}(U)^T \boldsymbol{\gamma}_j, \quad j = 1, \dots, p \quad (8.4)$$

تقریب می‌زنیم که  $\boldsymbol{\gamma}_j = (\gamma_{j,1}, \dots, \gamma_{j,K})^T$ . فرض کنید  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_p)^T$  و  $\boldsymbol{\gamma} = (\boldsymbol{\gamma}_1^T, \dots, \boldsymbol{\gamma}_p^T)^T$ . با استفاده از تقریب (۸.۴)، مدل (۶.۴) را می‌توان به صورت

$$Y_i = \mathbf{X}_i^T \boldsymbol{\eta} + \boldsymbol{\Pi}_i^T \boldsymbol{\gamma} + \varepsilon \quad (9.4)$$

نوشت. مدل (۹.۴) یک مدل خطی استاندارد است و هر تابع ناپارامتری  $\alpha_j(\cdot)$  در مدل (۶.۴)، با استفاده از  $\eta_j$  و  $\gamma_j$  در مدل (۹.۴) مشخص می‌شود. حال برای برآورد نیرومند پارامترهای  $\boldsymbol{\eta}$  و  $\boldsymbol{\gamma}$  از رگرسیون نما استفاده می‌کنیم. این روش علاوه بر نیرومند بودن، برآوردی را نتیجه می‌دهد که بسیار کارا است. بنابراین، با ترکیب رگرسیون نما و توابع اسپلاین، تابع هدف زیر را در نظر می‌گیریم

$$Q(\boldsymbol{\eta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \phi_h(Y_i - \mathbf{X}_i^T \boldsymbol{\eta} - \boldsymbol{\Pi}_i^T \boldsymbol{\gamma}) \quad (10.4)$$

که  $\phi_h(t) = \phi(t/h)/h$  و  $\phi(\cdot)$  تابع چگالی هسته با پهنای باند  $h$  است. برای سهولت در محاسبات، سرتاسر این فصل، از تابع چگالی نرمال استاندارد به عنوان تابع چگالی هسته استفاده می‌شود. با ماکزیمم کردن تابع هدف (۱۰.۴)، برآوردهای ضرایب  $\hat{\boldsymbol{\eta}}$  و  $\hat{\boldsymbol{\gamma}}$  به دست می‌آیند. در نتیجه ضرایب  $\alpha_j(U)$  به صورت

$$\hat{\alpha}_j(U) = \hat{\eta}_j + \sum_{k=1}^K B_{jk}(U) \hat{\gamma}_{j,k}$$

برآورد می‌شود. به خاطر داشته باشید که هدف اصلی ما تعیین متغیرهای بی‌اهمیت ( $\alpha_j(U) \equiv 0$ ) و مؤلفه‌های پارامتری ( $\alpha_j(U) \equiv \alpha_j$ ) است. اگر  $\|\boldsymbol{\gamma}_j\|_1 = \sum_{k=1}^K |\gamma_{j,k}| \neq 0$ ، آنگاه  $X_j$  دارای اثر متغیر است. اگر  $\|\boldsymbol{\gamma}_j\|_1 = 0$  و  $\eta_k \neq 0$  باشد،  $X_j$  دارای اثر ثابت غیر صفر است. در غیر این صورت  $X_j$  بی‌اثر است.

تابع هدف (۱۰.۴) قابلیت انتخاب متغیرهای مهم و تشخیص اثرات ثابت و متغیر را ندارد. لذا، به جای (۱۰.۴) از تابع هدف جریمه‌شده با دو تابع جریمه به صورت

$$L_\lambda(\boldsymbol{\eta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \phi_h(Y_i - \mathbf{X}_i^T \boldsymbol{\eta} - \boldsymbol{\Pi}_i^T \boldsymbol{\gamma}) - n \sum_{j=1}^p p_{\lambda_1}(|\eta_j|) - n \sum_{j=1}^p p_{\lambda_2}(\|\boldsymbol{\gamma}_j\|_{\mathbb{R}}) \quad (11.4)$$

استفاده می‌شود که  $\mathbf{R} = \int_{[0,1]} \mathbf{B}(u)\mathbf{B}(u)^T du$  و  $p_\lambda(\cdot)$  تابع جریمه SCAD (فن و لی، ۲۰۰۱) و  $\lambda_1, \lambda_2 \geq 0$  پارامترهای جریمه هستند. تابع جریمه اول در (۱۱.۴)، ضرایب ثابت  $\eta_j$  در  $\alpha_j(U)$  را به سمت صفر منقبض می‌کند. تابع جریمه دوم منقبض کردن بخش تابعی،  $g_j(U)$  را انجام می‌دهد. بنابراین اگر  $X_j$  دارای اثر ثابت باشد، همه مولفه های  $\gamma_j$  دقیقاً صفر برآورد می‌شوند. از طرفی اگر  $X_j$  دارای هیچ اثری روی متغیر پاسخ نباشد،  $\eta_j$  و همه مولفه های  $\gamma_j$  به سمت صفر منقبض می‌شوند. در غیر از این دو حالت،  $X_j$  دارای اثر متغیر است.

فرض کنید

$$(\hat{\eta}, \hat{\gamma}) = \arg \max_{\eta, \gamma} L_\lambda(\eta, \gamma).$$

بنابراین، مجموعه اندیس اثرات متغیر، اثرات ثابت و اثرات صفر، به ترتیب، به صورت

$$\hat{\mathcal{I}}_V = \{j : \|\hat{\gamma}_j\|_1 \neq 0, j = 1, \dots, p\}$$

$$\hat{\mathcal{I}}_C = \{j : \|\hat{\gamma}_j\|_1 = 0, \hat{\eta}_j \neq 0, j = 1, \dots, p\}$$

$$\hat{\mathcal{I}}_Z = \{j : \|\hat{\gamma}_j\|_1 = 0, \hat{\eta}_j = 0, j = 1, \dots, p\}$$

تعیین می‌شوند. در نتیجه اگر  $\alpha_j(U)$  ضریب متغیر باشد، به صورت  $\hat{\alpha}_j(U) = \hat{\eta}_j + \mathbf{B}(U)^T \hat{\gamma}_j$  و در صورت ثابت بودن به صورت  $\hat{\eta}_j$  برآورد می‌شود.

### ۳.۴ ویژگی‌های جانبی

در این بخش، به بررسی برخی ویژگی‌های جانبی برآوردگرهای رگرسیون نمای جریمه‌شده در رابطه (۲.۴)،  $(\hat{\eta}, \hat{\gamma})$ ، می‌پردازیم. فرض کنید  $\alpha_{\circ j}(\cdot)$  مقدار درست  $\alpha_j(\cdot)$  باشد. با توجه به تجزیه (۷.۴)،  $\eta_{\circ j} = E\{\alpha_{\circ j}(U)\}$  و  $g_{\circ j}(\cdot) = \alpha_{\circ j}(\cdot) - E\{\alpha_{\circ j}(U)\}$ ، به ترتیب، مقادیر درست  $\eta_j$  و  $g_j(\cdot)$  هستند. فرض کنید بردار  $\gamma_{\circ} = (\gamma_{\circ 1}^T, \dots, \gamma_{\circ p}^T)^T$  ضریب بهترین تقریب برای  $g_j(U)$  در فضای  $B$ -اسپلاین باشد، به طوری که

$$\|\mathbf{B}(U)^T \gamma_{\circ j} - g_{\circ j}(U)\|_\infty = O(K_n^{-r}), j \in \mathcal{I}_V, \quad \|\gamma_{\circ j}\|_1 = 0, j \in \mathcal{I}_C \cup \mathcal{I}_Z$$

که  $r$  در شرط C۲ زیر تعیین می‌شود. تعریف کنید

$$G(x, u, h) = E\{\phi_h'(\varepsilon) | \mathbf{X} = x, U = u\}$$

$$F(x, u, h) = E\{\phi_h''(\varepsilon) | \mathbf{X} = x, U = u\}$$

9

$$a_n = \max_{k,j} \{ |p'_{\lambda_1}(|\eta_{\circ k}|)|, |p'_{\lambda_2}(\|\gamma_{\circ j}\|_{\mathbf{R}})| : |\eta_{\circ k}| \neq 0, \|\gamma_{\circ j}\|_1 \neq 0 \}$$

$$b_n = \max_{k,j} \{ |p''_{\lambda_1}(|\eta_{\circ k}|)|, |p''_{\lambda_2}(\|\gamma_{\circ j}\|_{\mathbf{R}})| : |\eta_{\circ k}| \neq 0, \|\gamma_{\circ j}\|_1 \neq 0 \}.$$

برای اثبات ویژگی‌های مجانبی برآوردگرهای ارائه‌شده، به شرایط نظم زیر نیاز داریم:

C۱. متغیر شاخص  $U$  دارای تکیه‌گاه کراندار  $\Omega$  و تابع چگالی آن،  $f_U(\cdot)$ ، مثبت و دارای مشتق دوم پیوسته است. بدون از دست دادن کلیت، فرض می‌کنیم  $\Omega$  روی فاصله واحد  $[0, 1]$  باشد.

C۲. به‌ازای هر  $j \in \mathcal{I}_v$  داریم  $\alpha_{\circ j}(u) \in \mathcal{H}_r$ ، که مجموعه همه توابع  $g$  روی  $[0, 1]$  است، به‌طوری که مشتق مرتبه  $\nu$  ام  $g$ ،  $g^\nu$ ، موجود و در شرط لیب‌شیتز مرتبه  $\alpha$  صدق می‌کند، یعنی

$$|g^\nu(s) - g^\nu(t)| \leq C|s - t|^\alpha, \quad 0 \leq s, t \leq 1$$

که  $C$  یک مقدار ثابت مثبت،  $\nu$  یک عدد صحیح مثبت و  $\alpha \in (0, 1]$  است به‌طوری که  $r = \nu + \alpha > 2$ .

C۳. فرض کنید  $\Sigma(u) = E\{\mathbf{X}\mathbf{X}^T | U = u\}$  نسبت به  $u$  پیوسته است. علاوه براین، برای  $u$  داده‌شده،  $\Sigma(u)$  ماتریس معین مثبت است و مقادیر ویژه آن کراندارند. همچنین فرض می‌کنیم  $\max_i \|\mathbf{X}_i\| / \sqrt{n} = o_p(1)$ .

C۴. فرض کنید  $\xi_j = \tau_j - \tau_{j-1}$  و  $\xi = \max_j \{\xi_j\}$ . یک مقدار ثابت مانند  $C_0$  وجود دارد به‌طوری که  $\max_j \{|\xi_{j+1} - \xi_j|\} = o(K_n^{-1})$  و  $\frac{\xi}{\min_j \{\xi_j\}} \leq C_0$ .

C۵.  $F(x, u, h)$  و  $G(x, u, h)$  نسبت به  $(x, u)$  پیوسته هستند.

C۶. برای هر  $h > 0$ ،  $F(x, u, h) < 0$ .

C۷.  $E\{\phi'_h(\varepsilon) | \mathbf{X} = x, U = u\} = 0$ . همچنین میانگین‌های شرطی  $E\{\phi''_h(\varepsilon) | \mathbf{X} = x, U = u\}$  و  $E\{\phi'''_h(\varepsilon) | \mathbf{X} = x, U = u\}$  نسبت به  $x$  پیوسته هستند.

C۸. به‌ازای  $j \in \mathcal{I}_Z \cup \{l : E\alpha_{\circ l}(U) = 0, l \in \mathcal{I}_V\}$  داریم  $\liminf_{n \rightarrow \infty} \liminf_{\eta_j \rightarrow 0^+} p'_{\lambda_1}(\eta_j) / \lambda_1 > 0$  و به‌ازای  $j \in \mathcal{I}_C \cup \mathcal{I}_Z$  داریم  $\liminf_{n \rightarrow \infty} \liminf_{\|\gamma_j\|_{\mathbf{R}} \rightarrow 0^+} p'_{\lambda_2}(\|\gamma_j\|_{\mathbf{R}}) / \lambda_2 > 0$ .

در ادامه، قضیه ۱.۳.۴ نرخ همگرایی برآوردگرهای ضرایب متغیر و قضیه ۲.۳.۴ توزیع مجانبی برآوردگرهای ضرایب ثابت غیر صفر را ارائه می‌دهد.

قضیه ۱.۳.۴. فرض کنید شرایط C۱ - C۸ برقرار هستند و تعداد گره‌ها  $K_n = O(n^{\frac{1}{r+1}})$  باشد. اگر وقتی  $n \rightarrow \infty$ ، داشته باشیم  $b_n \rightarrow 0$ ، آنگاه

$$\|\hat{\alpha}_k(u) - \alpha_{\circ k}(u)\|_2 = O_p(n^{\frac{-r}{r+1}} + a_n), \quad k \in \mathcal{I}_V.$$

**قضیه ۲.۳.۴.** تحت شرایط  $C\lambda - C1$ ، فرض کنید تعداد گره‌ها  $K_n = O(n^{\frac{1}{r+1}})$ ،  $\lambda_{\max} = \max\{\lambda_1, \lambda_2\}$  و  $\lambda_{\min} = \min\{\lambda_1, \lambda_2\}$  باشد. اگر  $\lambda_{\max} \rightarrow 0$  و  $\lambda_{\min} \rightarrow \infty$  وقتی  $n \rightarrow \infty$ ، آنگاه

$$\sqrt{n}(\hat{\eta}_{\mathcal{I}_C} - \eta_{\circ\mathcal{I}_C}) \xrightarrow{d} \mathcal{N}(0, G(h)F^{-2}(h)\Sigma_{\mathcal{I}_C}^{-1})$$

که  $\eta_{\circ\mathcal{I}_C}$  مقدار درست بردار ضرایب ثابت غیر صفر،  $\hat{\eta}_{\mathcal{I}_C}$  برآوردگر آن،  $G(h)$  و  $F(h)$  به صورت

$$G(h) = \frac{1}{2\pi h^\epsilon} E\{\varepsilon \exp(-\varepsilon^2/2h^2)\}^2, F(h) = \frac{1}{\sqrt{2\pi}h^3} E\{\exp(-\varepsilon^2/2h^2)(\varepsilon^2/h^2 - 1)\}$$

و  $\Sigma_{\mathcal{I}_C}$  در رابطه (۲۹.آ) تعریف می‌شود.

## ۴.۴ انتخاب پهنای باند و الگوریتم پیاده‌سازی

در این بخش، ابتدا به موضوع انتخاب پهنای باند،  $h$ ، می‌پردازیم. سپس از الگوریتم MEM (لی و همکاران، ۲۰۰۷) و تقریب LQA (فن و لی، ۲۰۰۱) برای پیاده‌سازی روش ارائه‌شده استفاده می‌کنیم.

### ۱.۴.۴ انتخاب پهنای باند بهینه

نیرومندی و کارایی برآوردگرهای نما، به پهنای باند  $h$  بستگی دارد. بنابراین، مقدار بهینه آن ( $h_{opt}$ ) را طوری انتخاب می‌کنیم که واریانس مجانبی برآوردگرهای ارائه‌شده را، در مقایسه با روش کمترین توان‌های دوم، تا حد ممکن کوچک کند. برای سهولت، فرض می‌شود که  $\varepsilon$  از  $\mathbf{X}$  و  $U$  مستقل است. ژائو و زوئی (۲۰۰۹) نشان دادند که برآوردگر کمترین توان‌های دوم ضرایب رگرسیونی ثابت در مدل ضریب متغیر خطی- جزئی دارای توزیع مجانبی زیر است

$$\sqrt{n}(\hat{\eta}_{\mathcal{I}_C} - \eta_{\circ\mathcal{I}_C}) \xrightarrow{d} \mathcal{N}(0, \sigma^2 \Sigma_{\mathcal{I}_C}^{-1})$$

که  $\Sigma_{\mathcal{I}_C}$  در قضیه ۲.۳.۴ تعریف شده‌است. بنابراین نسبت واریانس مجانبی برآورد رگرسیون نما به واریانس مجانبی برآوردگر کمترین توان‌های دوم به صورت

$$r(h) = \frac{G(h)F^{-2}(h)}{\sigma^2}$$

به دست می‌آید که  $\sigma^2 = E(\varepsilon^2)$  و  $G(h)$  و  $F(h)$  در قضیه ۲.۳.۴ داده شده‌اند. نسبت  $r(h)$  فقط به  $h$  وابسته است. بنابراین، مقدار بهینه  $h$  به صورت

$$h_{opt} = \arg \min_h r(h) = \arg \min_h G(h)F^{-2}(h) \quad (۱۲.۴)$$

حاصل می‌شود. از رابطه (۱۲.۴) می‌توان دید که  $h_{opt}$  به  $n$  بستگی نداشته و تنها به توزیع شرطی  $\varepsilon$  وابسته است.



در عمل، به دلیل نامعلوم بودن توزیع خطاها، محاسبه مستقیم  $G(h)$  و  $F(h)$  امکان‌پذیر نیست. از این رو، ابتدا باقیمانده‌های  $\hat{\varepsilon}_i = Y_i - \sum_{j=1}^p \hat{\alpha}_j(U_i) X_{ij}$  را محاسبه می‌کنیم که ضرایب  $\hat{\alpha}_j(\cdot)$  با یک روش ساده برآورد می‌شوند. بهتر است برای محاسبه این ضرایب از یک روش نیرومند ساده مانند رگرسیون چندکی استفاده شود. پس از محاسبه باقیمانده‌ها،  $G(h)$  و  $F(h)$  را به صورت

$$\hat{G}(h) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{2\pi}h^2} \left\{ \hat{\varepsilon}_i \exp(-\hat{\varepsilon}_i^2 / 2h^2) \right\}^2$$

$$\hat{F}(h) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\sqrt{2\pi}h^3} \left\{ \exp(-\hat{\varepsilon}_i^2 / 2h^2) (\hat{\varepsilon}_i^2 / h^2 - 1) \right\}$$

برآورد می‌کنیم. بنابراین،  $r(h)$  به صورت

$$\hat{r}(h) = \frac{\hat{G}(h)\hat{F}^{-2}(h)}{\hat{\sigma}^2}$$

برآورد می‌شود. با این حال، چون هیچ جواب صریحی برای  $h$  وجود ندارد، از روش جستجوی شبکه‌ای<sup>۱۶</sup> استفاده می‌کنیم. براساس پیشنهاد یائو و همکاران (۲۰۱۲) و ژانگ و همکاران (۲۰۱۳)، نقاط ممکن برای  $h$  را به صورت

$$h = \hat{\sigma} \times 10^{2j}, \quad j = 0, \dots, k$$

در نظر می‌گیریم که  $k$  یک مقدار ثابت مانند  $k = 5$  یا  $k = 10$  است. در نتیجه، با استفاده از روش جستجوی شبکه‌ای، پهنای باند  $h_{opt}$  براساس (۱۲.۴) به دست می‌آید.

## ۲.۴.۴ الگوریتم MEM

برای حل مسئله بهینه‌سازی (۱۱.۴)، ابتدا از تقریب LQA (فن و لی، ۲۰۰۱)، برای تقریب توابع جریمه  $p_{\lambda_1}(\cdot)$  و  $p_{\lambda_2}(\cdot)$  استفاده می‌کنیم. براساس این تقریب، برای مقادیر اولیه داده شده  $\eta_j^\circ$  و  $\gamma_j^\circ$ ، توابع جریمه را به صورت

$$p_{\lambda_1}(|\eta_j|) \approx p_{\lambda_1}(|\eta_j^\circ|) + \frac{1}{2} \frac{p'_{\lambda_1}(|\eta_j^\circ|)}{|\eta_j^\circ|} \left\{ |\eta_j|^2 - |\eta_j^\circ|^2 \right\}$$

و

$$p_{\lambda_2}(\|\gamma_j\|_{\mathbf{R}}) \approx p_{\lambda_2}(\|\gamma_j^\circ\|_{\mathbf{R}}) + \frac{1}{2} \frac{p'_{\lambda_2}(\|\gamma_j^\circ\|_{\mathbf{R}})}{\|\gamma_j^\circ\|_{\mathbf{R}}} \left\{ \|\gamma_j\|_{\mathbf{R}}^2 - \|\gamma_j^\circ\|_{\mathbf{R}}^2 \right\}$$

تقریب می‌زنیم. فرض کنید  $\mathbf{Z}_i = (\mathbf{X}_i^T, \mathbf{\Pi}_i^T)^T$ . قرار دهید  $m = 0$  و  $\boldsymbol{\theta}^{(m)} = (\boldsymbol{\eta}^{(m)T}, \boldsymbol{\gamma}^{(m)T})^T$  برای  $\varepsilon_1 > 0$  و  $\varepsilon_2 > 0$ ، ماتریس بلوکی  $\boldsymbol{\Sigma}_\lambda(\boldsymbol{\theta}^{(m)})$  را به صورت زیر تعریف کنید

$$\boldsymbol{\Sigma}_\lambda(\boldsymbol{\theta}^{(m)}) = \text{diag} \left( \frac{p'_{\lambda_1}(|\eta_1^{(m)}|)}{|\eta_1^{(m)}| + \varepsilon_1}, \dots, \frac{p'_{\lambda_1}(|\eta_p^{(m)}|)}{|\eta_p^{(m)}| + \varepsilon_1}, \frac{p'_{\lambda_2}(\|\gamma_1^{(m)}\|_{\mathbf{R}})}{\|\gamma_1^{(m)}\|_{\mathbf{R}} + \varepsilon_2}, \dots, \frac{p'_{\lambda_2}(\|\gamma_p^{(m)}\|_{\mathbf{R}})}{\|\gamma_p^{(m)}\|_{\mathbf{R}} + \varepsilon_2} \mathbf{R} \right).$$

<sup>16</sup>Grid search method

براساس پیشنهاد هانتز و لی (۲۰۰۵)،  $\varepsilon_1$  و  $\varepsilon_2$  را به صورت

$$\varepsilon_1 = \frac{\tau}{\sqrt{2n\lambda_1}} \min\{|\eta_j^{(0)}|, j = 1, \dots, p\}, \quad \varepsilon_2 = \frac{\tau}{\sqrt{2n\lambda_2}} \min\{\|\gamma_j^{(0)}\|_{\mathbf{R}}, j = 1, \dots, p\}$$

در نظر می‌گیریم که  $\tau = 10^{-8}$ . اکنون می‌توانیم با در نظر گرفتن تابع چگالی هسته نرمال استاندارد و با استفاده از الگوریتم MEM، پارامترهای مجهول را به صورت زیر برآورد کنیم:

• گام ۱ (E):  $\pi(i|\theta^{(m)})$  را به صورت زیر به روز رسانی کنید

$$\pi(i|\theta^{(m)}) = \frac{\phi_h(Y_i - \mathbf{Z}_i^T \theta^{(m)})}{\sum_{i=1}^n \phi_h(Y_i - \mathbf{Z}_i^T \theta^{(m)})} \propto \phi_h(Y_i - \mathbf{Z}_i^T \theta^{(m)}), \quad i = 1, \dots, n.$$

• گام ۲ (M): مقدار  $\theta$  را به صورت زیر به روز رسانی کرده تا  $\theta^{(m+1)}$  به دست آید

$$\begin{aligned} \theta^{(m+1)} &= \arg \max_{\theta} \sum_{i=1}^n \pi(i|\theta^{(m)}) \log \phi_h(Y_i - \mathbf{Z}_i^T \theta^{(m)}) + \frac{n}{\sqrt{2}} \theta^T \Sigma_{\lambda}(\theta^{(m)}) \theta \\ &= (\mathbf{Z}^T \mathbf{W} \mathbf{Z} + n \Sigma_{\lambda}(\theta^{(m)}))^{-1} \mathbf{Z}^T \mathbf{W} \mathbf{Y} \end{aligned}$$

که  $\mathbf{W}$  یک ماتریس قطری به صورت  $\mathbf{W} = \text{diag}\{\pi(1|\theta^{(m)}), \dots, \pi(n|\theta^{(m)})\}$  است.

• گام ۳: قرار دهید  $m = m + 1$  و گام‌های اول و دوم را تا رسیدن به همگرایی تکرار کنید.

در این الگوریتم، شرط توقف تکرار،  $\|\hat{\theta}^{(m+1)} - \hat{\theta}^{(m)}\| < 10^{-8}$  است. برآورد حاصل را با  $\hat{\theta} = (\hat{\eta}^T, \hat{\gamma}^T)^T$  نشان می‌دهیم. در عمل مؤلفه‌های  $\hat{\theta}$  دقیقاً برابر با صفر نیستند. لذا با استفاده از یک مقدار آستانه‌ای، مؤلفه‌های کوچک را برابر صفر قرار می‌دهیم. در مطالعات شبیه‌سازی، هنگامی که  $|\hat{\eta}_j| < 10^{-4}$ ، آنگاه  $\hat{\eta}_j$  را برابر صفر، و هنگامی که  $\|\hat{\gamma}_j\|_2 < 10^{-4}$ ، آنگاه قرار می‌دهیم  $\hat{\gamma}_j = 0$ . علاوه بر این از برآورد رگرسیون چندکی به عنوان برآورد اولیه ضرایب استفاده می‌شود. در الگوریتم فوق، وقتی از هسته نرمال برای  $\phi(\cdot)$  استفاده می‌شود، تابع بهینه‌شده در گام ۲ مجموع وزنی لگاریتم تابع درست‌نمایی متناظر با رگرسیون خطی معمولی است. در این حالت، یک عبارت بسته-شکل برای ماکزیمم‌کننده حاصل می‌شود. اگر در گام ۲ هسته دیگری، غیر از نرمال، به کار رود، برای ماکزیمم‌کردن لازم است از برخی الگوریتم‌های بهینه‌سازی استفاده شود.

### ۳.۴.۴ انتخاب پارامترهای تنظیم‌کننده

برای پیاده‌سازی روش ارائه‌شده، لازم است مقدار مناسبی برای پارامترهای تنظیم‌کننده انتخاب شود. این پارامترها عبارتند از: مرتبه  $B$  - اسپلاین  $q$ ، تعداد گره‌های داخلی  $K_n$  و پارامترهای جریمه  $\lambda_1$  و  $\lambda_2$ . در مطالعات شبیه‌سازی از اسپلاین مکعبی ( $q = 3$ ) و  $K_n = \lceil n^{1/(2q+3)} \rceil$  گره داخلی با فواصل یکسان استفاده می‌شود. این روش انتخاب گره‌ها توسط زوئی و کیو<sup>۱۷</sup> (۲۰۱۲) پیشنهاد شده است.

همچنین مقدار بهینه  $(\lambda_1, \lambda_2)$  را با اعتبارسنجی متقابل ۵ گروهی انتخاب می‌کنیم. بطور واضح‌تر، داده‌های  $\mathbf{D}_n$  را بطور تصادفی به ۵ گروه با اندازه تقریباً یکسان افراز می‌کنیم. این گروه‌ها را، برای  $j = 1, 2, 3, 4, 5$  به صورت  $\{\mathbf{D}_i, i \in \mathcal{I}_j\}$  نشان می‌دهیم که  $\mathcal{I}_j$  مجموعه اندیس مشاهدات متناظر با  $j$  امین گروه می‌باشد. برای هر  $j$ ، مجموعه  $\{\mathbf{D}_i, i \in \mathcal{I}_j\}$  را به عنوان داده‌های آزمون و ۴ گروه دیگر را به عنوان داده‌های آموزشی در نظر می‌گیریم. برای هر زوج مشخص  $(\lambda_1, \lambda_2)$  و هر  $\mathcal{I}_j$ ، ابتدا از داده‌های آموزشی برای برآورد ضرایب متغیر براساس (۱۱.۴) استفاده کرده و متغیر پاسخ را برای داده‌های آزمون به صورت  $\hat{Y}_i = \mathbf{X}_i^T \hat{\boldsymbol{\eta}} + \mathbf{\Pi}_i^T \hat{\boldsymbol{\gamma}}, i \in \mathcal{I}_j$  پیشگویی می‌کنیم. سپس امتیاز اعتبارسنجی متقابل متناظر با این زوج  $(\lambda_1, \lambda_2)$  را به صورت

$$CV(\lambda_1, \lambda_2) = \sum_{j=1}^5 \sum_{i \in \mathcal{I}_j} \phi_{h_{opt}}(Y_i - \hat{Y}_i) \quad (13.4)$$

محاسبه می‌کنیم. در پایان،  $\lambda_1$  و  $\lambda_2$  بهینه با ماکزیمم کردن  $CV(\lambda_1, \lambda_2)$  به دست می‌آید.

## ۵.۴ مطالعات شبیه‌سازی

در این بخش، با استفاده از شبیه‌سازی به بررسی عملکرد رگرسیون نمای جریمه‌شده می‌پردازیم. برای سادگی، رگرسیون نمای جریمه‌شده را با PMR نشان می‌دهیم. عملکرد این روش را با دو روش تکراری ارائه‌شده توسط تانگ و همکاران (۲۰۱۲)، رگرسیون کمترین توان‌های دوم جریمه‌شده (PLS) و رگرسیون چندکی جریمه‌شده با چندک  $\tau = 0.5$  (PQR)، مقایسه می‌کنیم. تانگ و همکاران (۲۰۱۲) از روش‌های تکراری دومرحله‌ای برای انتخاب متغیر و تشخیص ساختار در مدل ضریب متغیر خطی-جزئی استفاده کردند. آنها در مرحله اول به تشخیص ضرایب ثابت و متغیر، و در مرحله دوم به تشخیص ضرایب ثابت غیر صفر و صفر پرداختند. برای بررسی نیرومندی و کارایی روش ارائه‌شده، سه نوع توزیع مختلف برای خطا در نظر می‌گیریم: توزیع نرمال استاندارد، توزیع تی-استودنت با ۳ درجه آزادی و توزیع آمیخته نرمال (MN) که به صورت  $0.9\mathcal{N}(0, 1) + 0.1\mathcal{N}(0, 25)$  MN تعریف می‌شود. برای بررسی میزان دقت روش ارائه شده در انتخاب مدل درست، از معیارهای زیر استفاده می‌کنیم:

- NV: متوسط تعداد ضرایب انتخاب شده به‌عنوان ضرایب متغیر
- NC: متوسط تعداد ضرایب انتخاب شده به‌عنوان ضرایب ثابت
- NVT: متوسط تعداد ضرایب متغیر انتخاب شده که واقعاً متغیر هستند
- NCT: متوسط تعداد ضرایب ثابت غیر صفر انتخاب شده که واقعاً ثابت غیر صفر هستند
- CF: درصد دفعاتی که ساختار مدل، دقیقاً به درستی تشخیص داده می‌شود.

برای بررسی دقت برآوردگر  $\hat{\alpha}(\cdot)$ ، از ریشه دوم میانگین توان‌های دوم خطا (RMSE) استفاده می‌کنیم که به صورت

$$\text{RMSE} = \left\{ \frac{1}{n_{\text{grid}}} \sum_{i=1}^{n_{\text{grid}}} \sum_{j=1}^p (\hat{\alpha}_j(u_i) - \alpha_j(u_i))^2 \right\}^{\frac{1}{2}} \quad (14.4)$$

تعریف می‌شود و  $\{u_1, \dots, u_{\text{grid}}\}$  نقاط شبکه‌ای هستند که تابع  $\hat{\alpha}_j(\cdot)$  در این نقاط ارزیابی می‌شود.

**مثال ۱۵.۴.** فرض کنید  $p = 15$  و داده‌ها را از مدل زیر تولید کنید

$$Y_i = \sum_{j=1}^p \alpha_j(U_i) X_{ij} + \varepsilon_i, \quad i = 1, \dots, n \quad (15.4)$$

که  $\alpha_5(u) = 1/5$ ،  $\alpha_4(u) = 1$ ،  $\alpha_3(u) = 2/5$ ،  $\alpha_2(u) = 8u(1-u)$ ،  $\alpha_1(u) = 2 \sin(2\pi u)$  و برای  $6 \leq j \leq p$ ،  $\alpha_j(u) = 0$ . لذا در این مدل  $NVT = 2$  و  $NCT = 3$  است. متغیر  $U$  از توزیع  $U[0, 1]$  و متغیرهای توضیحی  $(X_1, \dots, X_p)$  از توزیع نرمال چندمتغیره با میانگین صفر و ماتریس کوواریانس  $\Sigma = (\sigma_{ij})_{p \times p}$  تولید می‌شود که  $\sigma_{ij} = 0.5^{|i-j|}$ . نتایج شبیه‌سازی برای  $n = 200, 400$ ، پس از 500 تکرار، در جدول ۱.۴ گزارش شده‌اند.

در جدول ۱.۴ ملاحظه می‌شود که عملکرد هر سه روش تحت تاثیر توزیع خطاها قرار می‌گیرد. همانطور که انتظار داشتیم، برای هر نوع توزیع خطا، عملکرد هر یک از روش‌ها با افزایش حجم نمونه بهبود می‌یابد. برای هر دو مقدار حجم نمونه، رگرسیون نما، در مقایسه با دو روش دیگر، بسیار بهتر عمل می‌کند. هنگامی که توزیع خطاها نرمال است، مقادیر  $NVT = 1/63, 1/82$  و  $NCT = 2/99, 2/99$  برای  $n = 200$  و  $n = 400$  نشان می‌دهد که رگرسیون نمای جریمه‌شده در تشخیص اثرات ثابت و متغیر نسبت به دو روش دیگر کارا تر است. علاوه بر این، مقدار CF متناظر با رگرسیون نما بسیار بزرگتر از مقادیر CF متناظر با دو روش دیگر است. مشاهدات فوق از جدول ۱.۴، دلالت بر قابلیت بالای روش ارائه شده در شناسایی ساختار مدل درست دارند.

با توجه به مقادیر RMSE در جدول ۱.۴، می‌توان دید که در حضور داده‌های پرت یا توزیع دم‌سنگین خطاها، برآوردگرهای حاصل از رگرسیون نمای جریمه‌شده نسبت به رگرسیون کمترین توان‌های دوم جریمه‌شده و رگرسیون چندکی جریمه‌شده از دقت برآورد بیشتری برخوردار است، اما برای خطای نرمال، رگرسیون کمترین توان‌های دوم از رگرسیون چندکی و رگرسیون نما اندکی بهتر است.

برای بررسی بیشتر دقت برآوردگرها، ریشه دوم میانگین توان‌های دوم خطای متناظر با هر روش را با ریشه دوم میانگین توان‌های دوم برآوردگر پیشگو مقایسه کرده‌ایم. برآوردگر پیشگو برآوردگری است که با استفاده از مدل درست محاسبه می‌شود. منظور از مدل درست، مدلی است که ساختار خطی- جزئی آن از پیش معلوم است، یعنی ضرایب متغیر، ضرایب ثابت غیر صفر و صفر قبل از تحلیل داده‌ها معلوم‌اند و لذا برای به‌دست آوردن برآوردگر پیشگو از هیچ تابع جریمه‌ای استفاده نمی‌شود. نتایج شبیه‌سازی نشان می‌دهد که مقادیر RMSE برآوردگرهای جریمه‌شده و برآوردگر پیشگو بسیار نزدیک است.

جدول ۱.۴: نتایج تشخیص ساختار روش‌های مختلف در مثال ۱.۵.۴

RMSE	CF%	NCT(۳)*	NC(۳)*	NVT(۲)*	NV(۲)*	روش	خطا	n
جریمه‌شده پیشگو								
۰٫۲۹۷	۰٫۳۹۲	۵۱	۲٫۷۳	۲٫۸۴	۱٫۸۸	۲٫۵۷	PLS	$\mathcal{N}(0, 1)$ ۲۰۰
۰٫۳۴۱	۰٫۴۵۵	۴۵	۲٫۸۵	۳٫۲۹	۱٫۸۹	۲٫۳۷	PQR	
۰٫۳۱۶	۰٫۴۲۰	۹۲	۳٫۰۰	۳٫۰۵	۱٫۹۴	۱٫۹۴	PMR	
۰٫۴۲۴	۱٫۰۸۶	۲۲	۲٫۵۴	۲٫۹۷	۱٫۷۲	۳٫۰۹	PLS	t(۳)
۰٫۴۱۹	۰٫۶۸۳	۳۹	۲٫۷۹	۳٫۲۹	۱٫۷۳	۲٫۲۸	PQR	
۰٫۳۷۱	۰٫۴۷۷	۶۳	۲٫۹۹	۳٫۳۶	۱٫۶۳	۱٫۶۴	PMR	
۰٫۴۹۳	۱٫۰۹۰	۱۷	۲٫۴۷	۲٫۹۲	۱٫۷۳	۳٫۳۳	PLS	MN
۰٫۳۶۴	۰٫۶۴۲	۴۸	۲٫۸۴	۳٫۲۹	۱٫۷۸	۲٫۲۶	PQR	
۰٫۳۳۲	۰٫۴۵۸	۸۱	۲٫۹۹	۳٫۱۷	۱٫۸۲	۱٫۸۳	PMR	
۰٫۲۲۲	۰٫۲۴۱	۹۱	۲٫۹۹	۳٫۰۷	۱٫۹۲	۱٫۹۳	PLS	$\mathcal{N}(0, 1)$ ۴۰۰
۰٫۲۴۶	۰٫۲۹۶	۸۲	۲٫۹۷	۳٫۱۶	۱٫۹۶	۱٫۹۸	PQR	
۰٫۲۰۱	۰٫۲۵۹	۹۴	۳٫۰۰	۳٫۲۵	۲٫۰۰	۲٫۰۰	PMR	
۰٫۳۲۳	۰٫۵۱۰	۶۰	۲٫۷۹	۲٫۹۲	۱٫۸۶	۲٫۳۳	PLS	t(۳)
۰٫۳۰۸	۰٫۳۸۷	۷۸	۲٫۹۶	۳٫۱۴	۱٫۹۱	۱٫۹۷	PQR	
۰٫۲۵۴	۰٫۳۱۳	۸۷	۲٫۹۷	۳٫۳۵	۱٫۹۸	۲٫۰۰	PMR	
۰٫۴۰۷	۰٫۸۲۴	۵۵	۲٫۷۸	۲٫۸۷	۱٫۸۹	۲٫۴۴	PLS	MN
۰٫۳۱۲	۰٫۳۷۷	۸۲	۲٫۹۶	۳٫۱۰	۱٫۹۳	۲٫۰۱	PQR	
۰٫۲۴۵	۰٫۳۰۱	۹۱	۲٫۹۹	۳٫۱۸	۱٫۹۸	۱٫۹۸	PMR	

\* اعداد داخل پرانتز، مقدار درست آن معیار در مدل واقعی است.

هرچند که برای خطاهای نرمال، دقت برآورد رگرسیون کمترین توان‌های دوم جریمه‌شده اندکی بیشتر از دقت رگرسیون نمای جریمه‌شده است، اما مقادیر  $NVT = 1.94, 2$ ،  $NCT = 3, 3$  و  $CF = 0.92, 0.94$  برای  $n = 200, 400$  نشان می‌دهد که رگرسیون نمای جریمه‌شده، در مقایسه با دو روش رقیب، بهترین روش برای تشخیص ساختار مدل است.

**مثال ۲.۵.۴.** در این مثال، روش انتخاب متغیر معرفی شده را برای یک مدل ضریب متغیر خطی - جزئی با بعد بالاتر به کار می‌بریم. بدین منظور، مدل زیر را با ابعاد  $p = 30, 50$  در نظر بگیرید

$$Y_i = \sum_{j=1}^p \alpha_j(U_i) X_{ij} + \varepsilon_i, \quad i = 1, \dots, n \quad (16.4)$$

که  $\alpha_3(u) = u^2 \exp(1 + u)$ ،  $\alpha_2(u) = -2u(2 - 3u)^2$ ،  $\alpha_1(u) = 3 \cos(\pi(6u - 5)/3)$

جدول ۲.۴: درصد انتخاب درست ضرایب توسط رگرسیون نما در مثال ۲.۵.۴.

RMSE	CF%	$\alpha_6$	$\alpha_5$	$\alpha_4$	$\alpha_3$	$\alpha_2$	$\alpha_1$	خطا	p
جریمه شده									
پیشگو									
۰,۳۳۶۷	۰,۳۴۲۶	۹۸	۱۰۰	۱۰۰	۱۰۰	۹۸	۱۰۰	$\mathcal{N}(0, 1)$	۳۰
۰,۳۸۹۲	۰,۴۲۸۲	۸۹	۹۵	۱۰۰	۱۰۰	۹۰	۱۰۰	t(۳)	
۰,۳۵۸۳	۰,۴۳۴۱	۸۸	۹۷	۱۰۰	۱۰۰	۹۱	۱۰۰	MN	
۰,۳۳۶۱	۰,۳۶۵۲	۹۶	۱۰۰	۱۰۰	۱۰۰	۹۶	۱۰۰	$\mathcal{N}(0, 1)$	۵۰
۰,۳۹۳۵	۰,۴۴۶۷	۸۵	۹۵	۱۰۰	۱۰۰	۸۸	۱۰۰	t(۳)	
۰,۳۶۶۳	۰,۴۱۳۰	۸۷	۹۴	۱۰۰	۱۰۰	۸۸	۱۰۰	MN	

همانند مثال ۱.۵.۴ تولید می‌شوند. حجم نمونه تصادفی را  $n = 400$  در نظر گرفته و شبیه‌سازی را ۲۰۰ بار تکرار می‌کنیم. در این مثال، بعد ضرایب اسپلاین و بعد ماتریس طرح، به ترتیب، برابر با  $p(K_n + q)$  و  $p + p(K_n + q + 1)$  می‌باشند. درصد انتخاب درست هر یک از ضرایب متغیر و ثابت غیر صفر براساس ۲۰۰ تکرار در جدول ۲.۴ گزارش شده‌است. به‌عنوان مثال، عدد ۹۸ متناظر با  $\alpha_1$  بدین معنی است که ضریب  $\alpha_1(u)$  در ۹۸٪ اوقات به عنوان ضریب متغیر انتخاب می‌شود، یا عدد ۱۰۰ متناظر با  $\alpha_6$  بدین معنی است که در ۱۰۰٪ اوقات ضریب  $\alpha_6(u)$  به عنوان ثابت غیر صفر انتخاب می‌شود. همچنین مقدار RMSE و CF متناظر با هر روش را نیز در جدول ۲.۴ ارائه داده‌ایم.

نتایج شبیه‌سازی نشان می‌دهد که حتی برای مدل‌های ضریب متغیر با ابعاد بزرگتر، رگرسیون نمای جریمه‌شده با احتمال بالا می‌تواند ضرایب متغیر و ثابت غیر صفر را به درستی شناسایی کند.

## ۶.۴ تحلیل داده‌های قیمت مسکن

در این بخش، روش معرفی‌شده را برای تحلیل داده‌های مسکن بوستون<sup>۱۸</sup>، که توسط فن و هوانگ (۲۰۰۵) تحلیل شده‌است، به‌کار می‌بریم. این مجموعه داده در بسته “mlbench” نرم‌افزار R قابل دسترسی است و شامل ۱۴ متغیر زیر است:

– MEDV: میانه قیمت خانه

– LSTAT: درصد وضعیت پایین‌تر جامعه

– CRIM: نرخ سرانه جرم

– ZN: نسبت سطح ناحیه‌بندی شده با قطعات بزرگ

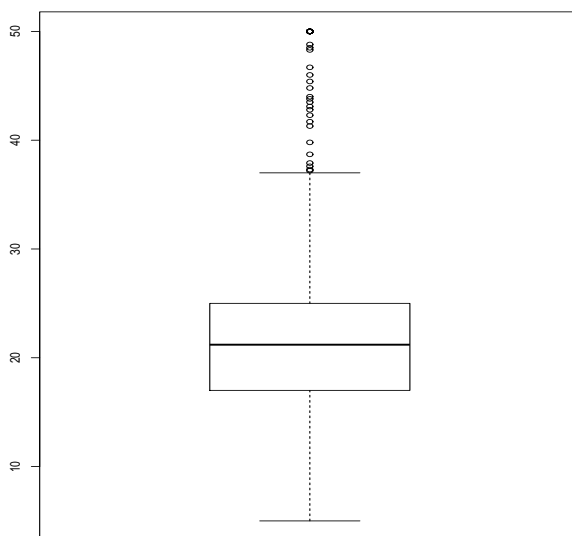
<sup>18</sup>Boston housing data

- INDUS: نسبت زمین‌های غیرتجاری
- CHAS: متغیر ساختگی رودخانه چارلز: اگر رودخانه مرزی است ۱، در غیر این صورت ۰
- NOX: غلظت اکسید نیتروژن
- RM: متوسط تعداد اتاق در هر خانه
- AGE: نسبت واحدهای ساخته‌شده قبل از سال ۱۹۴۰، که توسط مالک استفاده می‌شود
- DIS: میانگین وزنی فاصله تا پنج مرکز اشتغال بوستون
- RAD: شاخص دسترسی به بزرگراه‌های محوری
- TAX: نرخ مالیاتی
- PTRATIO: نسبت دانش آموز به معلم
- Black: نسبت جمعیت سیاه‌پوست.

این مجموعه داده شامل ۵۰۶ مشاهده است. هدف از تحلیل این داده‌ها، بررسی ارتباط بین متغیر پاسخ MEDV و سایر متغیرها است. پیرو پیشنهاد فن و هوانگ (۲۰۰۵)،  $U = \sqrt{LSTAT}$  را به عنوان متغیر شاخص در نظر می‌گیریم. این انتخاب متغیر  $U$  این امکان را به ما می‌دهد که برای درصدهای مختلف از وضعیت درآمد پایین جامعه، مدل‌های خطی مختلفی را برازش کنیم. علاوه بر این اثر سایر متغیرهای توضیحی با متغیر  $U$  در نظر گرفته می‌شود. مدل زیر را در نظر بگیرید

$$\begin{aligned}
 MEDV = & \alpha_0(U) + \alpha_1(U)CRIM + \alpha_2(U)ZN + \alpha_3(U)INDUS + \alpha_4(U)CHAS \\
 & + \alpha_5(U)NOX + \alpha_6(U)RM + \alpha_7(U)AGE + \alpha_8(U)DIS + \alpha_9(U)RAD \\
 & + \alpha_{10}(U)TAX + \alpha_{11}(U)PTRATIO + \alpha_{12}(U)Black + \epsilon.
 \end{aligned} \tag{۱۷.۴}$$

نمودار جعبه‌ای متغیر پاسخ MEDV در شکل ۳.۴ رسم شده است. این نمودار، وجود داده‌های پرت را در متغیر پاسخ نشان می‌دهد. از این رو، روش ارائه‌شده، به عنوان یک روش انتخاب متغیر نیرومند، برای تحلیل این داده‌ها به کار می‌رود. برای برازش مدل ضریب متغیر خطی-جزئی با روش پیشنهادی، از توابع اسپلاین مکعبی و ۴ گره با فواصل یکسان استفاده می‌شود. نتایج انتخاب متغیر و تشخیص ساختار در جدول ۳.۴ آورده شده است. جدول ۳.۴ نشان می‌دهد که متغیرهای INTERCEPT، CHAS، CRIM، NOX و RM دارای اثر متغیر بر پاسخ (نماد  $V$  در جدول) و متغیرهای DIS، RAD، PTRATIO، Black دارای اثر ثابت غیر صفر می‌باشند. همچنین متغیرهای ZN، INDUS، AGE و TAX به عنوان متغیرهای بی‌اثر بر پاسخ (نماد  $\circ$  در جدول) انتخاب شده‌اند. به منظور مقایسه، نتایج انتخاب متغیر و تشخیص ساختار با رگرسیون میانگین جرمی‌شده (PLS) و رگرسیون چندکی جرمی‌شده (PQR) نیز در جدول ۳.۴ گزارش شده است.



شکل ۲.۴: نمودار جعبه‌ای متغیر پاسخ MDEV.

همانطور که در بخش ۲.۴ بیان شد، پس از تقریب ضرایب با توابع B- اسپلاین، مدل ضریب متغیر به مدل خطی تبدیل می‌شود. لذا برای ارزیابی عملکرد روش پیشنهادی از ضریب تعیین تعدیل شده استفاده می‌کنیم. این ضریب به صورت  $R_{adj}^2 = 1 - [(1 - R^2)(n - 1)/(n - k - 1)]$  تعریف می‌شود که  $k$  تعداد متغیرهای توضیحی موجود در مدل و  $R^2 = 1 - RSS/\sum_{i=1}^n (y_i - \bar{y})^2$  که  $RSS$  مجموع توان‌های دوم باقیمانده‌هاست. در جدول ۳.۴، ضریب تعیین تعدیل شده،  $R_{adj}^2$  و انحراف معیار باقیمانده‌های متناظر با هر روش،  $\hat{\sigma}$ ، گزارش شده‌است. با توجه به مقادیر  $R_{adj}^2$  و  $\hat{\sigma}$ ، مدل برازش شده با استفاده از رگرسیون نما می‌تواند ۸۳/۶۲٪ از تغییرات کل قیمت مسکن را تبیین کند و انحراف استاندارد باقیمانده‌ها برابر با  $\hat{\sigma} = ۳/۸۰۷$  است. مقدار  $R_{adj}^2$  متناظر با رگرسیون نما بزرگتر از مقدار  $R_{adj}^2$  متناظر با دو روش دیگر است، یعنی مدل حاصل از رگرسیون نما تغییرات بیشتری از قیمت مسکن را تبیین می‌کند. حال می‌خواهیم بدانیم آیا متغیرهای انتخاب شده با ضرایب متغیر واقعاً با متغیر پاسخ مرتبط هستند. بدین منظور، برآوردهای ضرایب متغیر برازش شده و فواصل اطمینان نقطه به نقطه ۹۵٪ آنها، که با روش بوت‌استرپ به دست آمده، در شکل ۳.۴ نشان داده شده‌است. به وضوح، همه نمودارها نشان می‌دهند که صفر بودن این ضرایب بسیار نامحتمل است، زیرا برای هر یک از این متغیرها، خط  $\alpha(u) = 0$  کاملاً در داخل فاصله اطمینان قرار نمی‌گیرد. با توجه به شکل ۳.۴، ملاحظه می‌شود هنگامی که  $\sqrt{LSTAT} > 0.3$ ، ضریب متغیر CRIM بسیار نزدیک به صفر است، یعنی تاثیر میزان جرم و جنایت بر قیمت مسکن، وابسته به درصد وضعیت درآمد پایین جامعه است. این تعبیر، مانند تعبیر ضرایب رگرسیون چند متغیره، با فرض ثابت بودن سایر متغیرهای توضیحی به دست می‌آید. لازم به ذکر است، برای محاسبه فواصل اطمینان نقطه به نقطه هر یک از ضرایب متغیر، ابتدا با استفاده از  $B$  نمونه بوت‌استرپ مستقل،  $B$  برآورد  $\tilde{\alpha}_j(u)$  به دست می‌آوریم. حال یک فاصله اطمینان



جدول ۳.۴: نتایج انتخاب متغیر و برآورد ضرایب ثابت غیر صفر در داده‌های مسکن بوستون. نمادهای  $\circ$  و  $V$ ، به ترتیب، بیانگر اثر صفر و اثر متغیر است. مقادیر سطر آخر جدول، ضریب تعیین تعدیل شده و انحراف معیار باقیمانده‌ها را نشان می‌دهند.

متغیرهای توضیحی	PLS	PQR	PMR
INTERCEPT	V	V	V
CRIM	$\circ$	$\circ$	V
ZN	V	$\circ$	$\circ$
INDUS	$\circ$	-۰٫۰۳۹	$\circ$
CHAS	V	V	V
NOX	V	-۱۰٫۳۹	V
RM	V	V	V
AGE	-۰٫۲۱۲	-۰٫۳۰۷	$\circ$
DIS	-۰٫۸۵۴	-۱٫۱۰۸	-۰٫۵۶۸
RAD	V	V	۰٫۲۴۵
TAX	V	$\circ$	$\circ$
PTRATIO	-۱٫۰۳۲	-۰٫۷۸۴	-۰٫۶۷۸
Black	V	V	۰٫۰۱۹

( $R_{adj}^2, \hat{\sigma}$ ) (۷۹٫۴۴، ۴٫۳۵۷) (%۸۱٫۳۴، ۴٫۱۱۲) (۸۳٫۶۲، ۳٫۸۰۷)

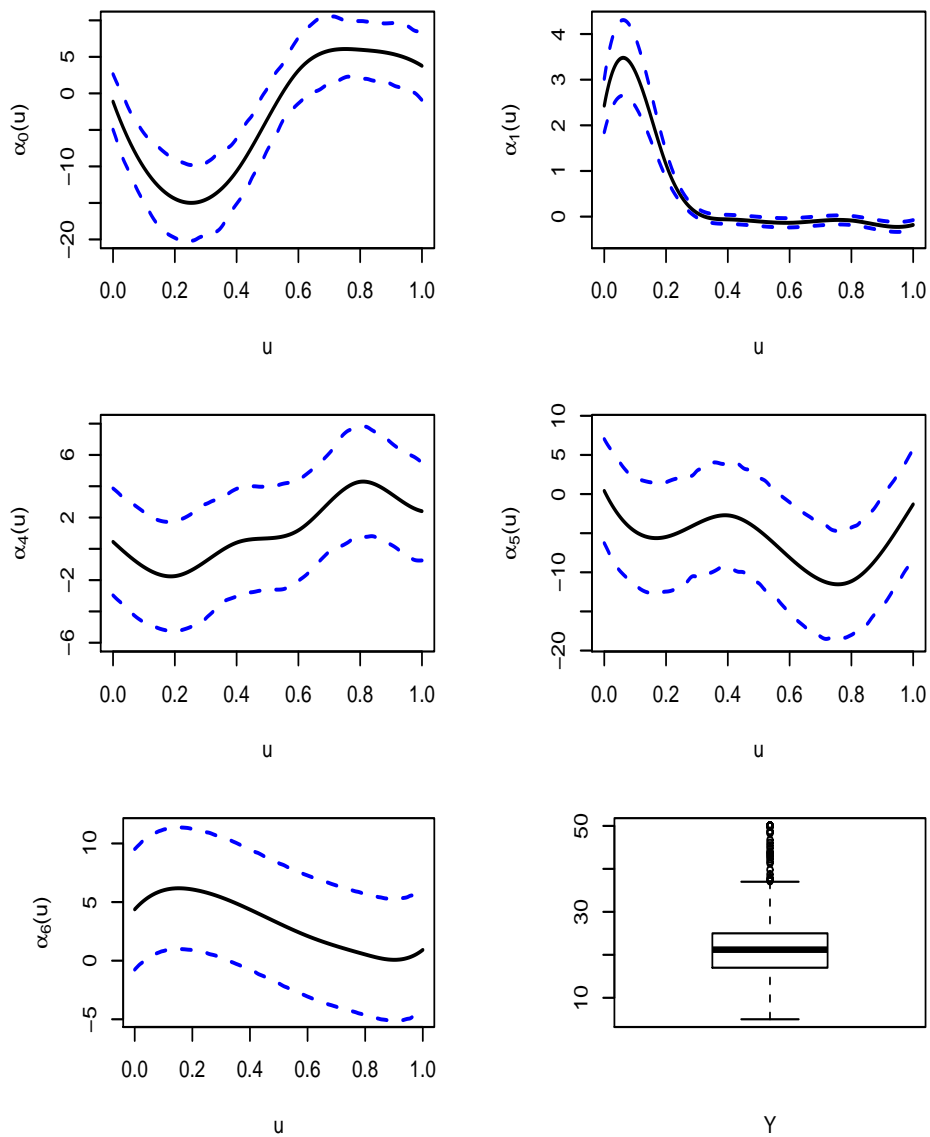
$1 - \alpha$  برای  $\alpha_j(u)$  براساس صدک‌های بوت‌استرپ به صورت

$$(L_{j,\alpha/2}(u), U_{j,\alpha/2}(u)), \quad j \in \hat{I}_V$$

به دست می‌آید که  $L_{j,\alpha/2}(u)$  و  $U_{j,\alpha/2}(u)$ ، به ترتیب، صدک‌های  $100 \times \alpha/2$  و  $100 \times (1 - \alpha/2)$  برآوردهای بوت‌استرپ  $\tilde{\alpha}_j(u)$  می‌باشند. علاوه بر این، براساس تقریب نرمال، فاصله

$$\hat{\alpha}_j(u) \pm Z_{1-\alpha/2} \hat{\sigma}_{j,B}(u), \quad j \in \hat{I}_V$$

یک فاصله اطمینان در سطح  $1 - \alpha$  برای  $\alpha_j(u)$  است که  $\hat{\sigma}_{j,B}(u)$  خطای استاندارد نمونه‌ای برآوردهای  $\tilde{\alpha}_j(u)$  و  $Z_{(1-\alpha/2)}$ ، صدک  $100 \times (1 - \alpha/2)$  توزیع نرمال استاندارد است. اگر اربیبی  $\hat{\alpha}_j(u)$  نسبت به واریانس آن بطور مجانبی قابل اغماض باشد، آنگاه  $\hat{\alpha}_j(u) \pm Z_{1-\alpha/2} \hat{\sigma}_{j,B}(u)$  یک فاصله اطمینان مجانبی در سطح  $1 - \alpha$  برای  $\alpha_j(u)$  است. برای اطلاعات بیشتر می‌توان به هوانگ و همکاران (۲۰۰۲) و وانگ و همکاران (۲۰۰۸) مراجعه کرد.



شکل ۳.۴: برآوردهای ضرایب متغیر مدل (۱۷.۴) و فواصل اطمینان نقطه‌ای ۹۵٪ این ضرایب با استفاده از رگرسیون نمای جریمه‌شده (PMR).

## بحث و نتیجه‌گیری

در این فصل، با به‌کارگیری رگرسیون نما، یک روش نیرومند برای برآورد ضرایب در مدل ضریب متغیر ارائه شد. این روش قادر است علاوه بر برآورد ضرایب، اثرات صفر، ثابت غیر صفر و متغیر را نیز بطور همزمان تشخیص دهد. کارایی روش ارائه‌شده با مثال شبیه‌سازی و تحلیل مجموعه داده واقعی مورد بررسی قرار گرفت. نتایج مطالعات عددی عملکرد رضایت‌بخش روش پیشنهادی را در حضور داده‌های پرت و خطاهای دم‌سنگین نشان می‌دهند.

اگرچه در این فصل مدل ضریب متغیر مورد بررسی قرار گرفته است، اما می‌توان روش ارائه‌شده را به مدل‌های دیگر مانند مدل ضریب متغیر تعمیم‌یافته، مدل تک-شاخص و مدل تک-شاخص ضریب متغیر تعمیم داد.

این روش پیشنهادی، مشابه سایر روش‌های انتخاب متغیر مبتنی بر جریمه، ممکن است برای داده‌های با بعد بسیار بالا دارای عملکرد مناسبی نباشد. برای رفع این چالش، می‌توان ابتدا با استفاده از روش‌های غربالگری بعد مسئله را از مقیاس بسیار بزرگ  $p$  به یک بعد کوچکتر از  $n$  کاهش داده و سپس از روش نیرومند ارائه شده برای انتخاب متغیر و تشخیص ساختار استفاده کرد. تاکنون چندین روش غربالگری برای کاهش بعد در مدل‌های ضریب متغیر معرفی شده است. به عنوان مثال، لیو و همکاران (۲۰۱۴) یک روش غربالگری برای مدل‌های ضریب متغیر براساس همبستگی شرطی (CC-SIS) معرفی کردند. همچنین، فن و همکاران (۲۰۱۴) و سانگ و همکاران (۲۰۱۴) روش‌های غربالگری دیگری برای این مدل‌ها ارائه دادند. علاوه بر این تاکنون چندین روش آزاد-مدل برای غربالگری معرفی شده است که در فصل ۲ به معرفی برخی از این روش‌ها پرداختیم. این روش‌ها می‌توانند برای غربالگری در مدل‌های ضریب متغیر مورد استفاده قرار گیرند.



# پیوست آ

## اثبات قضایا

### ۱.آ قضایای فصل ۳

اثبات قضیه ۱.۲.۳ در اثبات این قضیه،  $c_1$  و  $c_2$  مقادیر ثابت مثبت هستند که ممکن است، در جاهای مختلف، دارای مقادیر متفاوت باشند.

ابتدا فرض کنید  $F(y)$  معلوم باشد. بنابراین،  $\widehat{dcov}^* \{X_j, F(Y)\} = \hat{S}_{j,1}^* + \hat{S}_{j,2}^* - 2\hat{S}_{j,3}^*$  که در آن

$$\begin{aligned}\hat{S}_{j,1}^* &= \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| |F(Y_i) - F(Y_k)| \\ \hat{S}_{j,2}^* &= \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n |F(Y_i) - F(Y_k)| \\ \hat{S}_{j,3}^* &= \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n \sum_{l=1}^n |X_{ij} - X_{lj}| |F(Y_k) - F(Y_l)|.\end{aligned}$$

تعریف کنید

$$\hat{\omega}_k^* = \widehat{dcorr}^* \{X_k, F(Y)\}.$$

قضیه ۱ لی، ژانگ و ژو (۲۰۱۲) بیان می‌کند که برای هر  $0 < \gamma < 1/2 - \kappa$ ، ثابت‌های مثبت  $c_1$  و  $c_2$  وجود دارند به طوری که

$$P(|\hat{\omega}_j^* - \omega_j| \geq cn^{-\kappa}) \leq O\left(\left[\exp\{-c_1 n^{1-2(\kappa+\gamma)}\} + n \exp(-c_2 n^\gamma)\right]\right). \quad (1.A)$$

بنابراین برای اثبات قضیه، کافی است نشان دهیم که برای مقادیر بزرگ  $n$ ، اختلاف بین  $\hat{\omega}_j^*$  و  $\hat{\omega}_j$  قابل اغماض است، که معادل با بررسی اختلاف بین  $\hat{S}_{j,m}^*$  و  $\hat{S}_{j,m}$  به ازای  $m = 1, 2, 3$  است. در اینجا

قضیه را تنها برای  $m = ۱$  ثابت می‌کنیم، زیرا برای  $m = ۲, ۳$  نیز به‌طور مشابه ثابت می‌شود. به خاطر بیاورید که

$$\hat{S}_{j,1}^* = \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| |F(Y_i) - F(Y_k)|$$

۹

$$\hat{S}_{j,1} = \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| |F_n(Y_i) - F_n(Y_k)|.$$

بنابراین

$$\begin{aligned} & P\left(|\hat{S}_{j,1}^* - \hat{S}_{j,1}| \geq \varepsilon\right) \\ &= P\left(\frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n |X_{ij} - X_{kj}| \left| |F(Y_i) - F(Y_k)| - |F_n(Y_i) - F_n(Y_k)| \right| \geq \varepsilon\right) \\ &\leq P\left((A_n B_n)^{\frac{1}{\gamma}} \geq \varepsilon\right) \\ &\leq P\left((A_n B_n)^{\frac{1}{\gamma}} \geq \varepsilon, |X_j| \leq M\right) + P\left((A_n B_n)^{\frac{1}{\gamma}} \geq \varepsilon, |X_j| > M\right) \\ &=: I_1 + I_2 \end{aligned}$$

که  $M$  یک عدد مثبت بوده و بعداً مقدار آن مشخص می‌شود،  $A_n = \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n (X_{ij} - X_{kj})^2$  و  $B_n = \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n \{|F(Y_i) - F(Y_k)| - |F_n(Y_i) - F_n(Y_k)|\}^2$  حال با استفاده از نامساوی  $||x| - |y|| \leq |x - y| \leq |x| + |y|$  داریم

$$\begin{aligned} & \left| |F(Y_i) - F(Y_k)| - |F_n(Y_i) - F_n(Y_k)| \right| \\ &\leq |F_n(Y_i) - F(Y_i)| + |F_n(Y_k) - F(Y_k)| \leq 2 \max_{1 \leq i \leq n} |F_n(Y_i) - F(Y_i)|. \end{aligned}$$

علاوه براین، چون  $A_n \leq \frac{1}{n^\gamma} \sum_{i=1}^n \sum_{k=1}^n 2(X_{ij}^2 + X_{kj}^2) \leq 4M^2$ ، می‌توان نتیجه گرفت

$$\begin{aligned} I_1 &\leq P\left(2Mn^{-\gamma} \left\{ \sum_{i=1}^n \sum_{k=1}^n (|F(Y_i) - F(Y_k)| - |F_n(Y_i) - F_n(Y_k)|)^2 \right\}^{\frac{1}{\gamma}} \geq \varepsilon\right) \\ &\leq P\left(4M \max_{1 \leq i \leq n} |F_n(Y_i) - F(Y_i)| \geq \varepsilon\right) \leq P\left(\max_{y \in \mathbb{R}} |F_n(Y_i) - F(Y_i)| \geq \varepsilon/4M\right) \\ &\leq 2 \exp\{-2n(\varepsilon/4M)^\gamma\} = 2 \exp(-n\varepsilon^\gamma/8M^\gamma) \end{aligned} \quad (۲.آ)$$

که نامساوی آخر از نامساوی دورتزکی-کیفر-ولفویتز<sup>۱</sup> نتیجه می‌شود (پیوست پ را ببینید). از طرفی، به ازای هر  $0 < t \leq t_0$ ، که در شرط C۱ تعریف شده‌است، داریم

$$\begin{aligned} I_2 &\leq P(|X_j| > M) = P(\exp(t|X_j|) > \exp(tM)) \\ &\leq E\{\exp(t|X_j|)\} \exp(-tM) \leq C \exp(-tM) \end{aligned} \quad (۳.آ)$$

<sup>۱</sup>Dvoretzky-Kiefer-Wolfowitz inequality

که  $C$  یک ثابت مثبت است. نامساوی دوم از نامساوی مارکف و نامساوی آخر از شرط  $C1$  نتیجه می‌شود. پس با انتخاب  $M = O(n^\gamma)$  که  $0 < \gamma < \frac{1}{4} - \kappa$ ، و نامساوی‌های (۲.آ) و (۳.آ) برای برخی مقادیر ثابت مثبت  $c_1$  و  $c_2$  داریم

$$\begin{aligned} P\left(|\hat{S}_{j,1}^* - \hat{S}_{j,1}| \geq \varepsilon\right) &\leq 2 \exp(-n\varepsilon^2/\lambda M^2) + C \exp(-tM) \\ &\leq 2 \exp(-c_1 \varepsilon^2 n^{1-2\gamma}) + C \exp(-c_2 n^\gamma). \end{aligned} \quad (4.1)$$

بنابراین، نرخ همگرایی صورت  $\hat{\omega}_j^*$  بدست آمد. به‌طور مشابه، می‌توان نرخ همگرایی مخرج آن را نیز بدست آورد. در پایان، می‌توان نشان داد که نرخ همگرایی  $\hat{\omega}_j^*$  همان (۴.آ) است. بنابراین، با فرض  $\varepsilon = \nu n^{-\kappa}$  داریم

$$P\left(|\hat{\omega}_j^* - \hat{\omega}_j| \geq \nu n^{-\kappa}\right) \leq O\left(\exp\{-c_1 n^{1-2(\kappa+\gamma)}\} + \exp(-c_2 n^\gamma)\right). \quad (5.1)$$

با ترکیب (۱.آ) و (۵.آ) نتیجه می‌شود

$$\begin{aligned} P\left(|\hat{\omega}_j - \omega_j| \geq 2\nu n^{-\kappa}\right) &\leq P\left(|\hat{\omega}_j^* - \hat{\omega}_j| \geq \nu n^{-\kappa}\right) + P\left(|\hat{\omega}_j - \omega_j| \geq \nu n^{-\kappa}\right) \\ &\leq O\left(\exp\{-c_1 n^{1-2(\kappa+\gamma)}\} + n \exp(-c_2 n^\gamma)\right). \end{aligned} \quad (6.1)$$

بنابراین می‌توان نوشت

$$\begin{aligned} P\left(\max_{1 \leq j \leq p} |\hat{\omega}_j - \omega_j| \geq 2\nu n^{-\kappa}\right) &= P\left(\bigcup_{j=1}^p \left\{|\hat{\omega}_j - \omega_j| \geq 2\nu n^{-\kappa}\right\}\right) \\ &\leq \sum_{j=1}^p P\left(|\hat{\omega}_j - \omega_j| \geq 2\nu n^{-\kappa}\right) \\ &\leq p \max_{1 \leq j \leq p} P\left(|\hat{\omega}_j - \omega_j| \geq 2\nu n^{-\kappa}\right) \\ &\leq O\left(p \left[\exp\{-c_1 n^{1-2(\kappa+\gamma)}\} + n \exp(-c_2 n^\gamma)\right]\right). \end{aligned}$$

با انتخاب  $c = 2\nu$  اثبات کامل است.

**اثبات قضیه ۲.۲.۳** اگر  $\mathbb{M} \subsetneq \hat{\mathbb{M}}$ ، پس به ازای برخی مقادیر  $j \in \mathbb{M}$  داریم  $\hat{\omega}_j < cn^{-\kappa}$ . بنابراین از شرط  $C2$  نتیجه می‌شود که برای برخی مقادیر  $j \in \mathbb{M}$  داریم  $|\hat{r}_j - r_j| > cn^{-\alpha}$ . بنابراین

$$\{\mathbb{M} \subsetneq \hat{\mathbb{M}}\} \subseteq \{|\hat{r}_j - r_j| > cn^{-\alpha} : j \in \mathbb{M}\}. \quad (7.1)$$

از رابطه (۷.آ) نتیجه می‌شود  $\{\max_{j \in \mathbb{M}} |\hat{r}_j - r_j| \leq cn^{-\alpha}\} \subseteq \{\mathbb{M} \subseteq \hat{\mathbb{M}}\}$ . با استفاده از قضیه ۱.۲.۳ داریم

$$\begin{aligned} P(\mathbb{M} \subseteq \hat{\mathbb{M}}) &\geq P\left(\max_{j \in \mathbb{M}} |\hat{\omega}_j - \omega_j| \leq cn^{-\kappa}\right) \\ &\geq 1 - O\left(s_n \left[\exp\{-c_1 n^{1-2(\kappa+\gamma)}\} + n \exp(-c_2 n^\gamma)\right]\right). \end{aligned}$$

اثبات قضیه کامل است.

## ۲.آ قضایای فصل ۴

فرض کنید  $C$  یک ثابت مثبت باشد که در جاهای مختلف این بخش، ممکن است دارای مقدار مختلف باشد. لم‌های ۱.۲.آ و ۲.۲.آ در اثبات قضایای فصل ۴ استفاده می‌شود.

لم ۱.۲.آ. فرض کنید شرایط  $C1-C4$  برقرار است و  $K_n = O(n^{\frac{1}{r+1}})$  باشد. اگر وقتی  $n \rightarrow \infty$  داشته باشیم  $b_n \rightarrow 0$  آنگاه

$$\|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_{\mathcal{V}} = O_p(n^{\frac{-r}{r+1}} + a_n) \quad (\text{i})$$

$$\|\hat{g}_j(u) - g_{\circ j}(u)\|_{\mathcal{V}} = O_p(n^{\frac{-r}{r+1}} + a_n), \quad j = 1, \dots, p \quad (\text{ii})$$

اثبات. فرض کنید  $\boldsymbol{u} = (\boldsymbol{u}_1^T, \boldsymbol{u}_2^T)^T$  و  $\boldsymbol{\gamma} = \boldsymbol{\gamma}_0 + \delta_n \boldsymbol{u}_2$ ،  $\boldsymbol{\eta} = \boldsymbol{\eta}_0 + \delta_n \boldsymbol{u}_1$ ،  $\delta_n = n^{\frac{-r}{r+1}} + a_n$  فرض کنید که  $\boldsymbol{u}_1$  و  $\boldsymbol{u}_2$  دو بردار، به ترتیب، با اندازه‌های  $p$  و  $p \times K$  هستند. ابتدا نشان می‌دهیم که برای هر  $\varepsilon > 0$  یک ثابت بزرگ  $C$  وجود دارد، به طوری که

$$P\left(\sup_{\|\boldsymbol{u}\|_{\mathcal{V}}=C} L_{\lambda}(\boldsymbol{\eta}, \boldsymbol{\gamma}) < L_{\lambda}(\boldsymbol{\eta}_0, \boldsymbol{\gamma}_0)\right) \geq 1 - \varepsilon. \quad (\text{A.1})$$

با استفاده از بسط تیلور حول  $\varepsilon_i + \boldsymbol{X}_i^T \boldsymbol{R}(U_i)$  داریم

$$\begin{aligned} K_n^{-1} \{L_{\lambda}(\boldsymbol{\eta}, \boldsymbol{\gamma}) - L_{\lambda}(\boldsymbol{\eta}_0, \boldsymbol{\gamma}_0)\} &\leq \frac{\delta_n}{K_n} \sum_{i=1}^n \phi'_h(\varepsilon_i + \boldsymbol{X}_i^T \boldsymbol{R}(U_i)) (\boldsymbol{X}_i^T \boldsymbol{u}_1 + \boldsymbol{\Pi}_i^T \boldsymbol{u}_2) \\ &\quad + \frac{\delta_n^2}{2K_n} \sum_{i=1}^n \phi''_h(\varepsilon_i + \boldsymbol{X}_i^T \boldsymbol{R}(U_i)) (\boldsymbol{X}_i^T \boldsymbol{u}_1 + \boldsymbol{\Pi}_i^T \boldsymbol{u}_2)^2 \\ &\quad + \frac{\delta_n^3}{6K_n} \sum_{i=1}^n \phi'''_h(\zeta_i) (\boldsymbol{X}_i^T \boldsymbol{u}_1 + \boldsymbol{\Pi}_i^T \boldsymbol{u}_2)^3 \\ &\quad - \frac{n}{K_n} \sum_{j \in \mathcal{S}} \{p_{\lambda_1}(|\eta_j|) - p_{\lambda_1}(|\eta_{0j}|)\} \\ &\quad - \frac{n}{K_n} \sum_{j \in \mathcal{I}_V} \{p_{\lambda_2}(\|\boldsymbol{\gamma}_j\|_{\mathbf{R}}) - p_{\lambda_2}(\|\boldsymbol{\gamma}_{0j}\|_{\mathbf{R}})\} \\ &:= \Delta_1 + \Delta_2 + \Delta_3 + \Delta_4 + \Delta_5 \end{aligned}$$

که در آن  $\varepsilon_i + \boldsymbol{X}_i^T \boldsymbol{R}(U_i)$  بین  $\zeta_i$  و  $\mathcal{S} = \mathcal{I}_C \cup \{l : E\alpha_{\circ l}(U) \neq 0, l \in \mathcal{I}_V\}$  دارای مقداری بین  $\varepsilon_i + \boldsymbol{X}_i^T \boldsymbol{R}(U_i) - \delta_n (\boldsymbol{X}_i^T \boldsymbol{u}_1 + \boldsymbol{\Pi}_i^T \boldsymbol{u}_2)$  و  $\boldsymbol{R}(U_i) = (R_1(U_i), \dots, R_p(U_i))^T$ ، همچنین، است. مؤلفه  $j$ ام آن به صورت  $R_j(U_i) = g_{\circ j}(U_i) - \mathbf{B}(U_i)^T \boldsymbol{\gamma}_{0j}$  تعریف می‌شود. حال با استفاده از بسط تیلور حول  $\varepsilon_i$ ، داریم

$$\begin{aligned} &\sum_{i=1}^n \phi'_h(\varepsilon_i + \boldsymbol{X}_i^T \boldsymbol{R}(U_i)) (\boldsymbol{X}_i^T \boldsymbol{u}_1 + \boldsymbol{\Pi}_i^T \boldsymbol{u}_2) \\ &= \sum_{i=1}^n \left\{ \phi'_h(\varepsilon_i) + \frac{1}{2} \phi''_h(\varepsilon_i) \boldsymbol{X}_i^T \boldsymbol{R}(U_i) + \frac{1}{6} \phi'''_h(\zeta_i^*) [\boldsymbol{X}_i^T \boldsymbol{R}(U_i)]^2 \right\} (\boldsymbol{X}_i^T \boldsymbol{u}_1 + \boldsymbol{\Pi}_i^T \boldsymbol{u}_2) \end{aligned}$$



که  $\zeta_i^*$  بین  $\varepsilon_i$  و  $\mathbf{X}_i^T \mathbf{R}(U_i)$  قرار دارد. از طرف دیگر،

$$\begin{aligned} \|\mathbf{X}_i^T \mathbf{R}(U_i) \mathbf{R}(U_i)^T \mathbf{X}_i\| &\leq \|\mathbf{X}_i^T \mathbf{X}_i\| \|\mathbf{R}(U_i)\|_\infty \\ &= \|\mathbf{X}_i^T \mathbf{X}_i\| O(K_n^{-r}). \end{aligned} \quad (9.1)$$

بنابراین نتیجه می‌شود

$$\begin{aligned} \sum_{i=1}^n \phi'_h(\varepsilon_i + \mathbf{X}_i^T \mathbf{R}(U_i)) (\mathbf{X}_i^T \mathbf{u}_\lambda + \mathbf{\Pi}_i^T \mathbf{u}_\nu) &= \sum_{i=1}^n (\mathbf{X}_i^T \mathbf{u}_\lambda + \mathbf{\Pi}_i^T \mathbf{u}_\nu) \|\mathbf{X}_i^T \mathbf{X}_i\| O(K_n^{-r}) \\ &= O_p(nK_n^{-r} \|\mathbf{u}\|_\nu). \end{aligned}$$

در نتیجه، داریم  $\Delta_1 = O_p(n\delta_n K_n^{-(r+1)} \|\mathbf{u}\|_\nu) = O_p(n\delta_n^\nu K_n^{-1} \|\mathbf{u}\|_\nu)$  به‌طور مشابه، می‌توان ثابت کرد که

$$\Delta_2 = E(F(\mathbf{X}, U, h)) O_p(nK_n^{-1} \delta_n^\nu \|\mathbf{u}\|_\nu), \quad \Delta_3 = O_p(nK_n^{-1} \delta_n^\nu \|\mathbf{u}\|_\nu).$$

لذا با انتخاب یک ثابت  $C$ ، که بقدر کافی بزرگ است،  $\Delta_2$  به‌طور یکنواخت  $\Delta_1$  را در  $\|\mathbf{u}\|_\nu = C$  مغلوب می‌کند. با توجه به شرط  $a_n \rightarrow 0$ ، می‌توان نوشت  $\delta_n \rightarrow 0$ . بنابراین  $\delta_n \|\mathbf{u}\|_\nu \rightarrow 0$  که نتیجه می‌دهد  $\Delta_3 = o_p(\Delta_2)$ . پس  $\Delta_3$  به‌طور یکنواخت توسط  $\Delta_2$  در  $\|\mathbf{u}\|_\nu = C$  مغلوب می‌شود. همچنین با استفاده از بسط تیلور داریم

$$\begin{aligned} p_{\lambda_1}(|\eta_j|) - p_{\lambda_1}(|\eta_{\circ j}|) &= (|\eta_j| - |\eta_{\circ j}|) p'_{\lambda_1}(|\eta_{\circ j}|) + \frac{1}{\nu} (|\eta_j| - |\eta_{\circ j}|)^\nu p''_{\lambda_1}(|\eta_{\circ j}|) \\ &\leq (|\eta_j| - |\eta_{\circ j}|) \max_{j \in \mathcal{S}} \{p'_{\lambda_1}(|\eta_{\circ j}|)\} + \frac{1}{\nu} (|\eta_j| - |\eta_{\circ j}|)^\nu \max_{j \in \mathcal{S}} \{p''_{\lambda_1}(|\eta_{\circ j}|)\} \\ &= (|\eta_{\circ j} + \delta_n \mathbf{u}_\lambda| - |\eta_{\circ j}|) a_n + \frac{1}{\nu} (|\eta_{\circ j} + \delta_n \mathbf{u}_\lambda| - |\eta_{\circ j}|)^\nu b_n \\ &\leq |\delta_n \mathbf{u}_\lambda| a_n + \frac{1}{\nu} |\delta_n \mathbf{u}_\lambda|^\nu b_n \\ &= O(\delta_n a_n \|\mathbf{u}\|_\nu) + O(\delta_n^\nu b_n \|\mathbf{u}\|_\nu). \end{aligned}$$

بنابراین نتیجه می‌شود

$$\Delta_4 = O_p(nK_n^{-1} \delta_n a_n \|\mathbf{u}\|_\nu + nK_n^{-1} \delta_n^\nu b_n \|\mathbf{u}\|_\nu).$$

سپس با توجه به شرط  $b_n \rightarrow 0$ ، به سادگی می‌توان نشان داد که  $\Delta_4$  توسط  $\Delta_2$  به‌طور یکنواخت در  $\|\mathbf{u}\|_\nu = C$  مغلوب می‌شود.

به‌طور مشابه، برای  $\Delta_\delta$ ، داریم

$$\begin{aligned}
 p_{\lambda_\nu}(\|\gamma_j\|_{\mathbf{R}}) - p_{\lambda_\nu}(\|\gamma_{\circ j}\|_{\mathbf{R}}) &= (\|\gamma_j\|_{\mathbf{R}} - \|\gamma_{\circ j}\|_{\mathbf{R}})p'_{\lambda_\nu}(\|\gamma_{\circ j}\|_{\mathbf{R}}) + \\
 &\quad \frac{1}{\nu}(\|\gamma_j\|_{\mathbf{R}} - \|\gamma_{\circ j}\|_{\mathbf{R}})^2 p''_{\lambda_\nu}(\|\gamma_{\circ j}\|_{\mathbf{R}}) \\
 &\leq (\|\gamma_j\|_{\mathbf{R}} - \|\gamma_{\circ j}\|_{\mathbf{R}}) \max_{j \in \mathcal{I}_V} \{p'_{\lambda_\nu}(\|\gamma_{\circ j}\|_{\mathbf{R}})\} + \\
 &\quad \frac{1}{\nu}(\|\gamma_j\|_{\mathbf{R}} - \|\gamma_{\circ j}\|_{\mathbf{R}})^2 \max_{j \in \mathcal{I}_V} \{p''_{\lambda_\nu}(\|\gamma_{\circ j}\|_{\mathbf{R}})\} \\
 &= (\|\gamma_{\circ j} + \delta_n \mathbf{u}_\nu\|_{\mathbf{R}} - \|\gamma_{\circ j}\|_{\mathbf{R}}) a_n + \\
 &\quad \frac{1}{\nu}(\|\gamma_{\circ j} + \delta_n \mathbf{u}_\nu\|_{\mathbf{R}} - \|\gamma_{\circ j}\|_{\mathbf{R}})^2 b_n \\
 &\leq \|\delta_n \mathbf{u}_\nu\|_{\mathbf{R}} a_n + \frac{1}{\nu} \|\delta_n \mathbf{u}_\nu\|_{\mathbf{R}}^2 b_n \\
 &= a_n \delta_n (\mathbf{u}_\nu^T \mathbf{R} \mathbf{u}_\nu)^{\frac{1}{\nu}} + \frac{1}{\nu} b_n \delta_n^2 (\mathbf{u}_\nu^T \mathbf{R} \mathbf{u}_\nu) \\
 &= O(a_n \delta_n \|\mathbf{u}\|_\nu) + O(b_n \delta_n^2 \|\mathbf{u}\|_\nu^2).
 \end{aligned}$$

تساوی آخر از این واقعیت نتیجه می‌شود که  $\|\mathbf{R}\| = O(1)$ . بنابراین داریم

$$\Delta_\delta = O_p(nK_n^{-1} a_n \delta_n \|\mathbf{u}\|_\nu + nK_n^{-1} b_n \delta_n^2 \|\mathbf{u}\|_\nu^2).$$

حال با شرط  $b_n \rightarrow 0$ ، می‌توان نتیجه گرفت که  $\Delta_\delta$  نیز توسط  $\Delta_\nu$  مغلوب می‌شود. با توجه به شرط  $C_6$ ، می‌دانیم که  $F(x, u, h) < 0$ . لذا با انتخاب یک ثابت  $C$  به‌قدر کافی بزرگ، نامساوی  $K_n^{-1} \{L_\lambda(\boldsymbol{\eta}, \boldsymbol{\gamma}) - L_\lambda(\boldsymbol{\eta}_\circ, \boldsymbol{\gamma}_\circ)\} < 0$  با احتمال نزدیک به یک برقرار است، که برقرار بودن (۸.آ) را نتیجه می‌دهد. در نتیجه، با احتمال نزدیک به یک، یک ماکزیمم موضعی در گوی  $\{(\hat{\boldsymbol{\eta}}^T, \hat{\boldsymbol{\gamma}}^T)^T + \delta_n \mathbf{u} : \|\mathbf{u}\|_\nu \leq C\}$  وجود دارد، یعنی یک ماکزیمم‌کننده موضعی مانند  $(\hat{\boldsymbol{\eta}}^T, \hat{\boldsymbol{\gamma}}^T)^T$  موجود است به‌طوری که

$$(\hat{\boldsymbol{\eta}}^T, \hat{\boldsymbol{\gamma}}^T)^T = (\boldsymbol{\eta}_\circ^T, \boldsymbol{\gamma}_\circ^T)^T + \delta_n \mathbf{u}. \quad (۱۰.آ)$$

از رابطه (۱۰.آ) نتیجه می‌شود

$$\left\| \frac{\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_\circ}{\delta_n} \right\|_\nu = \|\mathbf{u}_1\|_\nu < C, \quad \left\| \frac{\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_\circ}{\delta_n} \right\|_\nu = \|\mathbf{u}_2\|_\nu < C. \quad (۱۱.آ)$$

اکنون با جایگذاری  $\delta_n = n^{-\frac{r}{r+1}} + a_n$  در (۱۱.آ) داریم

$$\|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_\circ\|_\nu = O_p(n^{-\frac{r}{r+1}} + a_n), \quad \|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_\circ\|_\nu = O_p(n^{-\frac{r}{r+1}} + a_n). \quad (۱۲.آ)$$

بنابراین اثبات (i) کامل است.

حال به اثبات (ii) می‌پردازیم. در ابتدا، توجه کنید که

$$\begin{aligned} \|\hat{g}_j(u) - g_{\circ j}(u)\|_{\Psi} &= \int_{[0,1]} [\hat{g}_j(u) - g_{\circ j}(u)]_{\Psi}^2 du \\ &= \int_{[0,1]} [\mathbf{B}(u)^T \hat{\gamma}_j - \mathbf{B}(u)^T \gamma_{\circ j} + R_j(u)]_{\Psi}^2 du \\ &\leq 2 \int_{[0,1]} [\mathbf{B}(u)^T \hat{\gamma}_j - \mathbf{B}(u)^T \gamma_{\circ j}]_{\Psi}^2 du + 2 \int_{[0,1]} [R_j(u)]_{\Psi}^2 du \\ &= 2(\hat{\gamma}_j - \gamma_{\circ j})^T \mathbf{R}(\hat{\gamma}_j - \gamma_{\circ j}) + 2 \int_{[0,1]} [R_j(u)]_{\Psi}^2 du \quad (13.1) \end{aligned}$$

که  $\mathbf{R} = \int_{[0,1]} \mathbf{B}(u)\mathbf{B}(u)^T du$ . اکنون با در نظر گرفتن  $\|\mathbf{R}\| = O(1)$  (ژائو و همکاران، ۲۰۱۴) و  $\|R_j(u)\| = O(K_n^r) = O(n^{\frac{-r}{r+1}})$  داریم

$$(\hat{\gamma}_j - \gamma_{\circ j})^T \mathbf{R}(\hat{\gamma}_j - \gamma_{\circ j}) = O_p(n^{\frac{-r}{r+1}} + a_n^{\Psi}), \quad \int_{[0,1]} [R_j(u)]_{\Psi}^2 du = O_p(n^{\frac{-2r}{r+1}}). \quad (14.1)$$

بنابراین، با ترکیب روابط (13.1) و (14.1)، می‌توان (ii) را نتیجه گرفت. اثبات لم کامل است.

**لم ۲.۲.۴.** تحت شرایط بیان شده در قضیه ۲.۳.۴، داریم

(i) برای  $j \in \{1, \dots, p\} \setminus \mathcal{S}$ ، رابطه  $\hat{\eta}_j = 0$  با احتمال نزدیک به یک برقرار است.

(ii) برای  $j \in \mathcal{I}_C \cup \mathcal{I}_Z$ ، رابطه  $\hat{g}_j(u) \equiv 0$  با احتمال نزدیک به یک برقرار است.

**اثبات.** ابتدا (i) را ثابت می‌کنیم. بنا به ویژگی‌های تابع جریمه SCAD، برای مقادیر بزرگ  $n$ ، وقتی  $\lambda^{max} \rightarrow 0$  داریم  $a_n = 0$ . پس با توجه به لم ۱.۲.۴، کافی است نشان دهیم، برای هر  $\gamma$  که در شرط  $\|\gamma - \gamma_{\circ}\|_{\Psi} = O_p(n^{\frac{-r}{r+1}})$  و  $\eta_j = 0$  که در شرط  $|\eta_j - \eta_{\circ j}| = O_p(n^{\frac{-r}{r+1}})$ ،  $j \in \mathcal{S}$  صدق کند، نامساوی‌های

$$\frac{\partial L_{\lambda}(\boldsymbol{\eta}, \boldsymbol{\gamma})}{\partial \eta_j} < 0 \quad \text{for } 0 < \eta_j < \varepsilon, \quad j \in \{1, \dots, p\} \setminus \mathcal{S} \quad (15.1)$$

و

$$\frac{\partial L_{\lambda}(\boldsymbol{\eta}, \boldsymbol{\gamma})}{\partial \eta_j} > 0 \quad \text{for } -\varepsilon < \eta_j < 0, \quad j \in \{1, \dots, p\} \setminus \mathcal{S} \quad (16.1)$$

برای برخی مقادیر کوچک  $\varepsilon = Cn^{-r/(r+1)}$ ، با احتمال نزدیک به یک برقرار است. در نتیجه جواب ماکزیمم‌کننده  $L_{\lambda}(\boldsymbol{\eta}, \boldsymbol{\gamma})$  در  $\hat{\eta}_j = 0$  اتفاق می‌افتد که  $j \in \{1, \dots, p\} \setminus \mathcal{S}$ .

به طریق مشابه مباحث بیان شده در اثبات لم ۱.۲.آ، داریم

$$\begin{aligned} \frac{\partial L_{\lambda}(\boldsymbol{\eta}, \boldsymbol{\gamma})}{\partial \eta_j} &= \frac{\partial Q(\boldsymbol{\eta}, \boldsymbol{\gamma})}{\partial \eta_j} - np'_{\lambda_1}(|\eta_j|)\text{sgn}(\eta_j) \\ &= - \sum_{i=1}^n X_{ij} \phi'_h(Y_i - \mathbf{X}_i^T \boldsymbol{\eta} - \boldsymbol{\Pi}_i^T \boldsymbol{\gamma}) - np'_{\lambda_1}(|\eta_j|)\text{sgn}(\eta_j) \\ &= - \sum_{i=1}^n \left\{ X_{ij} \phi'_h(\varepsilon_i + \mathbf{X}_i^T \mathbf{R}(U_i)) - X_{ij} \phi''_h(\varepsilon_i + \mathbf{X}_i^T \mathbf{R}(U_i)) [\boldsymbol{\Pi}_i^T (\boldsymbol{\gamma} - \boldsymbol{\gamma}_\circ) \right. \\ &\quad \left. + \mathbf{X}_i^T (\boldsymbol{\eta} - \boldsymbol{\eta}_\circ)] + X_{ij} \phi'''_h(\zeta_i^*) [\boldsymbol{\Pi}_i^T (\boldsymbol{\gamma} - \boldsymbol{\gamma}_\circ) + \mathbf{X}_i^T (\boldsymbol{\eta} - \boldsymbol{\eta}_\circ)]^2 \right\} - np'_{\lambda_1}(|\eta_j|)\text{sgn}(\eta_j) \\ &= n\lambda_1 \left\{ \frac{-1}{\lambda_1} p'_{\lambda_1}(|\eta_j|)\text{sgn}(\eta_j) + O_p(\lambda_1^{-1} n^{-\frac{r}{r+1}}) \right\} \quad (17.آ) \end{aligned}$$

که  $\zeta_i^*$  مقداری بین  $Y_i - \mathbf{X}_i^T \boldsymbol{\eta} - \boldsymbol{\Pi}_i^T \boldsymbol{\gamma}$  و  $\varepsilon_i + \mathbf{X}_i^T \mathbf{R}(U_i)$

بنا به شرط C۸،  $\liminf_{n \rightarrow \infty} \liminf_{\eta_j \rightarrow 0^+} p'_{\lambda_1}(\eta_j)/\lambda_1 > 0$  و  $\lambda_1 n^{\frac{r}{r+1}} > \lambda_{\min} n^{\frac{r}{r+1}} \rightarrow \infty$ ، لذا علامت

مشتق فوق، کاملاً با علامت  $\eta_j$  تعیین می‌شود. پس (۱۵.آ) و (۱۶.آ) برقرار است. بنابراین اثبات (i) کامل است.

با استفاده از روش به کار رفته در اثبات قسمت (i) این لم، برای  $j \in \mathcal{I}_C \cup \mathcal{I}_Z$  با احتمال نزدیک به یک داریم  $\|\hat{\gamma}_j\|_1 = 0$ . سپس با در نظر گرفتن  $\sup_u \|\mathbf{B}(u)\|_2 = O(1)$ ، قسمت (ii) به سادگی از  $\hat{g}_j(u) = \mathbf{B}(u)^T \hat{\gamma}_j$  نتیجه می‌شود.

**اثبات قضیه ۱.۳.۴** برای  $j \in \mathcal{I}_V$ ، با تعریف  $\hat{G}_j(u)$  و لم ۱.۲.آ، داریم

$$\begin{aligned} \|\hat{\alpha}_j(u) - \alpha_{\circ j}(u)\|_2^2 &= \int_{[0,1]} [\hat{\alpha}_j(u) - \alpha_{\circ j}(u)]^2 du \\ &\leq 2|\hat{\eta}_j - \eta_{\circ j}|^2 + 2\|\hat{g}_j(u) - g_{\circ j}(u)\|_2^2 \\ &= O_p(n^{-\frac{r}{r+1}} + a_n^2). \quad (18.آ) \end{aligned}$$

اثبات کامل است.

**اثبات قضیه ۲.۳.۴** فرض کنید  $\boldsymbol{\gamma}_{\mathcal{I}_V} = ((\eta_j, \gamma_j^T), j \in \mathcal{I}_V)^T$ ،  $\boldsymbol{\gamma}_{\mathcal{I}_V} = ((\eta_{\circ j}, \gamma_{\circ j}^T), j \in \mathcal{I}_V)^T$  و  $\boldsymbol{\eta}_{\mathcal{I}_C} = (\eta_j, j \in \mathcal{I}_C)^T$  با استفاده از لم ۲.۲.آ، با احتمال نزدیک به یک،  $(\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}, \hat{\boldsymbol{\eta}}_{\mathcal{I}_C})$  ماکزیمم کننده تابع

$$\mathcal{L}(\boldsymbol{\eta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \phi_h(Y_i - \mathbf{X}_i^{cT} \boldsymbol{\eta}_{\mathcal{I}_C} - \boldsymbol{\Pi}_i^{vT} \boldsymbol{\gamma}_{\mathcal{I}_V}) - n \sum_{j \in \mathcal{S}} p_{\lambda_1}(|\eta_j|) - n \sum_{j \in \mathcal{I}_V} p_{\lambda_2}(\|\boldsymbol{\gamma}_j\|_{\mathbf{R}}) \quad (19.آ)$$

است که در آن  $\mathbf{X}_i^c = (X_{ij}, j \in \mathcal{I}_C)^T$  و  $\boldsymbol{\Pi}_i^v = ((X_{ij}, X_{ij} \mathbf{B}(U_i)^T), j \in \mathcal{I}_V)^T$

فرض کنید  $\mathcal{L}_1(\boldsymbol{\eta}_{\mathcal{I}_C}, \boldsymbol{\gamma}_{\mathcal{I}_V}) = \frac{\partial \mathcal{L}(\boldsymbol{\eta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\eta}_{\mathcal{I}_C}}$  و  $\mathcal{L}_2(\boldsymbol{\eta}_{\mathcal{I}_C}, \boldsymbol{\gamma}_{\mathcal{I}_V}) = \frac{\partial \mathcal{L}(\boldsymbol{\eta}, \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}_{\mathcal{I}_V}}$ . پس  $\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}$  و  $\hat{\boldsymbol{\eta}}_{\mathcal{I}_C}$  باید در معادلات زیر صدق کنند

(۲۰.آ)

$$\frac{1}{n} \mathcal{L}_1(\hat{\boldsymbol{\eta}}_{\mathcal{I}_C}, \hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}) = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^c \phi'_h(Y_i - \mathbf{X}_i^{cT} \hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\Pi}_i^{vT} \hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}) - p'_{\lambda_1}(|\hat{\boldsymbol{\eta}}_{\mathcal{I}_C}|) \mathbf{o} \operatorname{sgn}(\hat{\boldsymbol{\eta}}_{\mathcal{I}_C}) = \mathbf{o}$$

و

$$\frac{1}{n} \mathcal{L}_2(\hat{\boldsymbol{\eta}}_{\mathcal{I}_C}, \hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}) = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\Pi}_i^v \phi'_h(Y_i - \mathbf{X}_i^{cT} \hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\Pi}_i^{vT} \hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}) - \kappa = \mathbf{o} \quad (21. \bar{A})$$

که  $p'_{\lambda_1}(|\hat{\boldsymbol{\eta}}|)$  یک زیر-بردار از  $p'_{\lambda_1}(|\hat{\boldsymbol{\eta}}|) \mathbf{o} \operatorname{sgn}(\hat{\boldsymbol{\eta}})$  است و  $j$ -امین مؤلفه از  $p'_{\lambda_1}(|\hat{\boldsymbol{\eta}}|)$  و  $\operatorname{sgn}(\hat{\boldsymbol{\eta}})$ ، به ترتیب، به صورت  $p'_{\lambda_1}(|\hat{\eta}_j|)$  و  $\operatorname{sgn}(\hat{\eta}_j)$  می‌باشند که  $j = 1, \dots, p$  و “ $\mathbf{o}$ ” نشان‌دهنده ضرب مؤلفه به مؤلفه است. همچنین  $\kappa$  یک بردار با طول  $K \times |\mathcal{I}_V|$  است که  $j$ -امین زیربردار بلوکی آن به صورت  $\frac{\hat{\gamma}_j}{\|\hat{\gamma}_j\|_{\mathbf{R}}} p'_{\lambda_2}(\|\hat{\gamma}_j\|_{\mathbf{R}})$  می‌باشد. با استفاده از بسط تیلور  $p'_{\lambda_1}(|\hat{\eta}_j|)$  می‌توان نوشت

$$p'_{\lambda_1}(|\hat{\eta}_j|) = p'_{\lambda_1}(|\eta_{\circ j}|) + (p''_{\lambda_1}(|\eta_{\circ j}|) + o_p(1))(\hat{\eta}_j - \eta_{\circ j}), \quad j \in \mathcal{I}_C. \quad (22. \bar{A})$$

توجه کنید که براساس یادآوری ۱ فن و لی (۲۰۰۱)، وقتی  $\lambda_{\max} \rightarrow \mathbf{o}$ ، آنگاه  $p'_{\lambda_1}(|\eta_{\circ j}|) = \mathbf{o}$ . بنابراین با توجه به شرط  $b_n \rightarrow \mathbf{o}$ ، با ترکیب (۲۰.آ) و (۲۲.آ) داریم

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^c \left\{ \phi'_h(\varepsilon_i) + \phi''_h(\varepsilon_i) \left[ \mathbf{X}_i^{vT} \mathbf{R}^v(U_i) - [\mathbf{X}_i^{cT} (\hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\eta}_{\circ \mathcal{I}_C}) + \boldsymbol{\Pi}_i^{vT} (\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V} - \boldsymbol{\gamma}_{\circ \mathcal{I}_V})] \right] \right. \\ & \quad \left. + \phi'''_h(\varepsilon_i^*) \left[ \mathbf{X}_i^{vT} \mathbf{R}^v(U_i) - [\mathbf{X}_i^{cT} (\hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\eta}_{\circ \mathcal{I}_C}) + \boldsymbol{\Pi}_i^{vT} (\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V} - \boldsymbol{\gamma}_{\circ \mathcal{I}_V})] \right]^2 \right\} \\ & \quad + o_p(\hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\eta}_{\circ \mathcal{I}_C}) = \mathbf{o} \end{aligned} \quad (23. \bar{A})$$

که  $\mathbf{R}^v(U_i) = (R_j(U_i), j \in \mathcal{I}_V)^T$  و  $\varepsilon_i^*$  بین  $\varepsilon_i$  و  $Y_i - \mathbf{X}_i^{cT} \hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\Pi}_i^{vT} \hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}$  قرار دارد. مشابه معادله (۲۳.آ)، با استفاده از بسط تیلور  $p'_{\lambda_2}(\|\hat{\gamma}_j\|_{\mathbf{R}})$  و رابطه (۲۱.آ) می‌توان نوشت

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \boldsymbol{\Pi}_i^v \left\{ \phi'_h(\varepsilon_i) + \phi''_h(\varepsilon_i) \left[ \mathbf{X}_i^{vT} \mathbf{R}^v(U_i) - [\mathbf{X}_i^{cT} (\hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\eta}_{\circ \mathcal{I}_C}) + \boldsymbol{\Pi}_i^{vT} (\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V} - \boldsymbol{\gamma}_{\circ \mathcal{I}_V})] \right] \right. \\ & \quad \left. + \phi'''_h(\varepsilon_i^{**}) \left[ \mathbf{X}_i^{vT} \mathbf{R}^v(U_i) - [\mathbf{X}_i^{cT} (\hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\eta}_{\circ \mathcal{I}_C}) + \boldsymbol{\Pi}_i^{vT} (\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V} - \boldsymbol{\gamma}_{\circ \mathcal{I}_V})] \right]^2 \right\} \\ & \quad + o_p(\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V} - \boldsymbol{\gamma}_{\circ \mathcal{I}_V}) = \mathbf{o} \end{aligned} \quad (24. \bar{A})$$

که  $\varepsilon_i^{**}$  بین  $\varepsilon_i$  و  $Y_i - \mathbf{X}_i^{cT} \hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\Pi}_i^{vT} \hat{\boldsymbol{\gamma}}_{\mathcal{I}_V}$  قرار دارد. فرض کنید  $\Phi_n = \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \boldsymbol{\Pi}_i^v \boldsymbol{\Pi}_i^{vT}$  و  $\Psi_n = \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \boldsymbol{\Pi}_i^v \mathbf{X}_i^{cT}$  داریم

$$\hat{\boldsymbol{\gamma}}_{\mathcal{I}_V} - \boldsymbol{\gamma}_{\circ \mathcal{I}_V} = (\Phi_n + o_p(1))^{-1} \{ -\Psi_n (\hat{\boldsymbol{\eta}}_{\mathcal{I}_C} - \boldsymbol{\eta}_{\circ \mathcal{I}_C}) + \Lambda_n \} \quad (25. \bar{A})$$

که  $\Lambda_n = \frac{1}{n} \sum_{i=1}^n \Pi_i^v [\phi'_h(\varepsilon_i) + \phi''_h(\varepsilon_i) \mathbf{X}_i^{vT} \mathbf{R}^v(U_i)]$  با جایگزینی آ.۲۵ در (آ.۲۳)، می توان نوشت

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \mathbf{X}_i^c \{ \mathbf{X}_i^c - \Psi_n^T \Phi(n)^{-1} \Pi_i^v \}^T (\hat{\eta}_{\mathcal{I}_C} - \eta_{\mathcal{I}_C}) + o_p(\hat{\eta}_{\mathcal{I}_C} - \eta_{\mathcal{I}_C}) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i^c \left\{ \phi'_h(\varepsilon_i) + \phi''_h(\varepsilon_i) \mathbf{X}_i^{cT} \mathbf{R}^v(U_i) - \phi''_h(\varepsilon_i) \Pi_i^{vT} \left[ \frac{1}{n} \sum_{k=1}^n \phi'_h(\varepsilon_k) \Phi_n^{-1} \Pi_k^v \right] \right\} \\ & \quad - \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \mathbf{X}_i^c \Pi_i^{vT} \left[ \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k^{vT} \mathbf{R}^v(U_k) \right]. \end{aligned} \quad (\text{آ.۲۶})$$

علاوه بر این داریم

$$\begin{aligned} E \left\{ \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \Psi_n^T \Phi_n^{-1} \Pi_i^v [\mathbf{X}_i^{cT} - \Pi_i^{vT} \Phi_n^{-1} \Psi_n] \right\} &= 0 \\ \text{Var} \left\{ \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \Psi_n^T \Phi_n^{-1} \Pi_i^v [\mathbf{X}_i^{cT} - \Pi_i^{vT} \Phi_n^{-1} \Psi_n] \right\} &= o_p\left(\frac{1}{n}\right). \end{aligned}$$

حال با استفاده از (آ.۲۶)، به سادگی می توان نشان داد

$$\left[ \frac{1}{n} \sum_{i=1}^n \phi''_h(\varepsilon_i) \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^T + o_p(1) \right] \sqrt{n} (\hat{\eta}_{\mathcal{I}_C} - \eta_{\mathcal{I}_C}) = I_1 + I_2 \quad (\text{آ.۲۷})$$

که در آن

$$\begin{aligned} I_1 &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi'_h(\varepsilon_i) \tilde{\mathbf{X}}_i \\ I_2 &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi''_h(\varepsilon_i) \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^{vT} \mathbf{R}^v(U_i) \end{aligned}$$

و  $\tilde{\mathbf{X}}_i = \mathbf{X}_i^c - \Psi_n^T \Phi_n^{-1} \Pi_i^v$  بنا به قضیه حد مرکزی، داریم

$$I_1 \xrightarrow{D} N(0, G(h) \Sigma_{\mathcal{I}_C}) \quad (\text{آ.۲۸})$$

که در آن

$$\Sigma_{\mathcal{I}_C} = E(\mathbf{X}_i^c \mathbf{X}_i^{cT}) - E(\mathbf{X}_i^c \Pi_i^{vT}) \{ E(\Pi_i^v \Pi_i^{vT}) \}^{-1} E(\Pi_i^v \mathbf{X}_i^{cT}). \quad (\text{آ.۲۹})$$

توجه کنید که در رابطه (آ.۲۸) برای محاسبه کوواریانس  $I_1$  از این واقعیات استفاده شده است که

$$\begin{aligned} \Phi_n^{-1} &\rightarrow E(\phi''_h(\varepsilon_i) \Pi_i^v \Pi_i^{vT})^{-1} \\ \Psi_n &\rightarrow E(\phi''_h(\varepsilon_i) \Pi_i^v \mathbf{X}_i^{cT}) \\ \Psi_n \Phi_n^{-1} &\rightarrow E[\phi''_h(\varepsilon_i)^{-1} (\Pi_i^v \Pi_i^{vT})^{-1} \phi''_h(\varepsilon_i) \Pi_i^v \mathbf{X}_i^{cT}] = E(\Pi_i^v \Pi_i^{vT})^{-1} E(\Pi_i^v \mathbf{X}_i^{cT}) \\ \Psi_n^T \Phi_n^{-1} \Pi_i^v \Pi_i^{vT} &= \Psi_n^T E(\Pi_i^v \Pi_i^{vT})^{-1} \Pi_i^v \Pi_i^{vT} = \Psi_n^T. \end{aligned}$$

با استفاده از روابط فوق نتیجه می‌شود

$$\begin{aligned}
 Cov(I_{\Psi}) &= \frac{1}{n} \sum_{i=1}^n Cov(\phi'_h(\varepsilon_i) \tilde{\mathbf{X}}_i) \\
 &= \frac{1}{n} \sum_{i=1}^n G(h) Cov(\tilde{\mathbf{X}}_i) \\
 &= E[\tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^T] - E(\tilde{\mathbf{X}}_i) E(\tilde{\mathbf{X}}_i^T) \\
 &= E[\tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^T] \\
 &= E[(\mathbf{X}_i^c - \Psi_n^T \Phi_n^{-1} \Pi_i^v)(\mathbf{X}_i^c - \Psi_n^T \Phi_n^{-1} \Pi_i^v)^T] \\
 &= E[\mathbf{X}_i^c \mathbf{X}_i^{cT} - \mathbf{X}_i^c \Pi_i^{vT} \Phi_n^{-1} \Psi_n - \Psi_n^T \Phi_n^{-1} \Pi_i^v \mathbf{X}_i^{cT} \\
 &\quad + \Psi_n^T \Phi_n^{-1} \Pi_i^v \Pi_i^{vT} \Phi_n^{-1} \Psi_n] \\
 &= E[\mathbf{X}_i^c \mathbf{X}_i^{cT} - \Psi_n^T \Phi_n^{-1} \Psi_n - \Psi_n^T \Phi_n^{-1} \Psi_n + \Psi_n^T \Phi_n^{-1} \Psi_n] \\
 &= E[\mathbf{X}_i^c \mathbf{X}_i^{cT}] - E(\mathbf{X}_i^c \Pi_i^{vT}) E(\Pi_i^v \Pi_i^{vT})^{-1} E(\Pi_i^v \mathbf{X}_i^c).
 \end{aligned}$$

با استفاده از تعریف  $\mathbf{R}^v(U_i)$ ، می‌توان ثابت کرد  $I_{\Psi} = o_p(1)$  (ژائو و همکاران، ۲۰۱۴). همچنین با استفاده از قانون اعداد بزرگ، داریم

$$\frac{1}{n} \sum_{i=1}^n \phi_h''(\varepsilon_i) \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^T \xrightarrow{p} E(\phi_h''(\varepsilon) \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T) = F(h) \Sigma_{I_C}. \quad (30.A)$$

از طرفی با استفاده از قضیه اسلاتسکی می‌دانیم که اگر  $A_n \xrightarrow{P} A$  و  $X_n \xrightarrow{D} X$ ، آنگاه  $\frac{X_n}{A_n} \xrightarrow{D} \frac{X}{A}$ . بنابراین با در نظر گرفتن روابط (27.A)–(30.A) و قضیه اسلاتسکی، اثبات کامل می‌شود.





# پیوست ب

## توابع اسپلاین و $B$ - اسپلاین

مدل ناپارامتری زیر را در نظر بگیرید

$$E(Y|X) = f(X)$$

که در آن  $f(\cdot)$  یک تابع هموار نامعلوم است و باید بر اساس یک نمونه تصادفی از  $(X, Y)$  برآورد شود. تقریب این تابع رگرسیونی "هموارسازی" نامیده می‌شود. بنابراین یک هموارسازی، تابعی برای خلاصه‌سازی متغیر پاسخ  $Y$  به‌عنوان تابعی از متغیر توضیحی  $X$  است. تاکنون روش‌های متعددی برای هموارسازی معرفی شده است. یکی از متداول‌ترین این روش‌ها، اسپلاین‌ها هستند. اسپلاین توسط شخصی به نام شوئنبرگ<sup>۱</sup> در سال ۱۹۴۶ معرفی شد. وی برای نخستین بار قضیه‌ای برای اثبات وجود اسپلاین‌های درون‌یاب بیان کرد و سپس مطالعه چندجمله‌ای‌های تکه‌ای برای توابع اسپلاین آغاز شد. در ادامه به بیان جزئیات این روش می‌پردازیم.

### ب.۱ اسپلاین

#### اسپلاین مرتبه یک

بازه  $[a, b]$  را به‌صورت اجتماع دو بازه  $[a, t] \cup [t, b]$  در نظر بگیرید. یک اسپلاین از مرتبه یک، تابعی است که در هر یک از بازه‌های مذکور، یک تابع خطی و در نقطه  $t$  نیز پیوسته باشد. این تابع را می‌توان

<sup>۱</sup>Schoenberg

به صورت زیر نوشت

$$f(x) = \begin{cases} \beta_0 + \beta_1 x & x \leq t \\ \beta'_0 + \beta'_1 x & x > t. \end{cases} \quad (1.ب)$$

با اعمال شرط پیوستگی تابع  $f$  در نقطه‌ی  $x = t$  داریم

$$\begin{aligned} \beta_0 + \beta_1 t &= \beta'_0 + \beta'_1 t \\ \beta'_0 &= \beta_0 + \beta_1 t - \beta'_1 t. \end{aligned}$$

با جایگذاری  $\beta'_0$  در ضابطه‌ی دوم  $f$  داریم

$$\beta'_0 + \beta'_1 x = \beta_0 + \beta_1 t - \beta'_1 t + \beta'_1 x = \beta_0 + \beta_1 t + \beta'_1 (x - t).$$

لذا

$$f(x) = \begin{cases} \beta_0 + \beta_1 x & x \leq t \\ \beta_0 + \beta_1 t + \beta'_1 (x - t) & x > t \end{cases}$$

که می‌توان این تابع دو ضابطه‌ای را به صورت زیر نوشت

$$f(x) = \beta_0 + \beta_1 x + \beta_2 (x - t)_+ \quad (2.ب)$$

که  $(x - t)_+ = (x - t)I(x > t)$  و  $\beta_2 = \beta'_1 - \beta_1$  تفاوت بین شیب خطوط اول و دوم رابطه (1.ب) را بیان می‌کند. در واقع می‌توان گفت که تابع  $f$ ، ترکیبی خطی از توابع پایه‌ای ۱،  $x$  و  $(x - t)_+$  است. اسپلاین مرتبه اول را اسپلاین خطی نیز گویند.

به همین ترتیب، فرض کنید که بتوان بازه  $[a, b]$  را به صورت اجتماع بازه‌های

$$[a, t_1), [t_1, t_2), \dots, [t_k, b]$$

نوشت به طوری که  $t_1 < \dots < t_k$ . نقاط  $\{t_i\}_{i=1}^k$  را گره نامند. در این صورت اسپلاین مرتبه اول را می‌توان به صورت

$$f(x) = \beta_0 + \beta_1 x + \sum_{j=1}^k \beta_{1+j} (x - t_j)_+$$

بیان کرد که ترکیبی خطی از توابع پایه‌ای ۱،  $x$ ،  $(x - t_1)_+$ ،  $\dots$  و  $(x - t_k)_+$  به شمار می‌آید.

### اسپلاین مرتبه $p$

توابع اسپلاین برای مراتب بالاتر را نیز می‌توان به‌طور مشابه تعریف کرد. فرض کنید بازه  $[a, b]$  به  $k + 1$  زیربازه به‌صورت  $[t_k, b], \dots, [t_1, t_2), [a, t_1)$  افزایش شده باشد به‌طوری که  $a < t_1 < \dots < t_k < b$ . تابع  $f(\cdot)$  را یک اسپلاین مرتبه‌ی  $p$  نامیم اگر در هر یک از زیر بازه‌های فوق یک چندجمله‌ای از مرتبه  $p$  باشد به‌طوری که  $f(\cdot)$  و مشتق‌های آن از مراتب ۱ تا  $p - 1$ ، در گره‌ها پیوسته باشند. لذا تابع اسپلاین مرتبه  $p$  را می‌توان به‌صورت زیر بیان کرد

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_p x^p + \sum_{j=1}^k \beta_{p+j} (x - t_j)^p \quad (3.ب)$$

که در آن

$$(x - t)_+^p = (x - t)^p I(x > t)$$

تابع توانی بریده‌شده<sup>۲</sup> از مرتبه  $p$  است. یعنی تابع اسپلاین مرتبه  $p$ ، ترکیبی خطی از یک چندجمله‌ای مرتبه‌ی  $p$  و توابع توانی بریده‌شده است. به‌عبارت دیگر مجموعه

$$\{1, x, x^2, \dots, x^p, (x - t_1)^p, \dots, (x - t_k)^p\}$$

یک پایه برای فضای اسپلاین‌های مرتبه‌ی  $p$  با گره‌های  $t_1, t_2, \dots, t_k$  است. این مجموعه را پایه‌های اسپلاین مرتبه  $p$  گویند.

### اسپلاین مکعبی

اگر  $p = 3$ ، آن‌گاه تابع مورد اشاره در رابطه (3.ب) را یک اسپلاین مکعبی<sup>۳</sup> نامند که به‌صورت

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \sum_{j=1}^k \beta_{3+j} (x - t_j)^3,$$

تعریف می‌شود. اسپلاین مکعبی یکی از متداول‌ترین نوع اسپلاین است که به دلیل انعطاف‌پذیری بالایی که دارد، در بسیاری از علوم مورد استفاده قرار می‌گیرد.

## ۲.ب رگرسیون اسپلاین

مدل رگرسیونی

$$y = f(x) + \epsilon$$

<sup>2</sup>Truncated Power Function

<sup>3</sup>Cubic Spline

را در نظر بگیرید. اگر تابع  $f(\cdot)$  را با تابع اسپلاین تعریف شده در رابطه (ب.۳) جایگزین کنیم، نتیجه حاصل را رگرسیون اسپلاین نامند.

یک برتری رگرسیون اسپلاین، برای روش‌های ناپارامتری و نیمه پارامتری این است که متعلق به رده رگرسیون خطی است. به عبارت دیگر، فرم تابع  $f(\cdot)$  نسبت به پارامترها خطی است و از این رو هزینه‌های محاسباتی برای برآورد منحنی رگرسیون مربوطه، ناچیز است. بنابراین برازش یک اسپلاین به داده‌ها می‌تواند با استفاده از روش کمترین توان‌های دوم خطا صورت گیرد. استفاده از توابع پایه‌ای بریده‌شده، منجر به برازش مدلی به صورت زیر می‌شود

$$y_i = \beta_0 + \beta_1 x_i + \dots + \beta_p x_i^p + \beta_{p+1} (x_i - t_1)_+^p + \dots + \beta_{p+k} (x_i - t_k)_+^p + \epsilon_i, \quad i = 1, \dots, n$$

که  $\epsilon_i$  ها در پذیره‌های زیربنایی رگرسیون صدق می‌کنند. شکل ماتریسی این عبارت را می‌توان به صورت

$$\mathbf{y} = \mathbf{T}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

بیان کرد که در آن  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_{p+k})^T$  و  $\mathbf{T}$  ماتریس طرح با اندازه  $n \times (p+k+1)$  است که درایه‌های ستون اول آن برابر ۱، درایه‌های  $p$  ستون بعدی، مربوط به مشاهدات متغیرهای  $x, x^2, \dots, x^p$  و درایه‌های ستون‌های  $p+2$  تا  $p+k+1$  مربوط به مشاهدات حاصل از متغیرهای  $(x-t_1)^p, \dots, (x-t_k)^p$  به ازای گره‌های از قبل تعیین شده  $t_1, t_2, \dots, t_k$  هستند، یعنی

$$\mathbf{T} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p & (x_1 - t_1)^p & (x_1 - t_2)^p & \dots & (x_1 - t_k)^p \\ 1 & x_2 & x_2^2 & \dots & x_2^p & (x_2 - t_1)^p & (x_2 - t_2)^p & \dots & (x_2 - t_k)^p \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p & (x_n - t_1)^p & (x_n - t_2)^p & \dots & (x_n - t_k)^p \end{bmatrix}$$

با به کار بردن روش کمترین توان‌های دوم، برآوردگر زیر حاصل می‌شود

$$\hat{\boldsymbol{\beta}} = (\mathbf{T}'\mathbf{T})^{-1} \mathbf{T}'\mathbf{y}. \quad (4.ب)$$

مشکلی که در برآوردگر (ب.۴) وجود دارد، همخطی بین پایه‌های توانی بریده‌شده

$$(x - t_1)_+^p, (x - t_2)_+^p, \dots, (x - t_k)^p$$

است که باعث عدم وارون‌پذیری یا بدشرطیدگی<sup>۴</sup> ماتریس  $\mathbf{T}'\mathbf{T}$  خواهد شد که منجر به ناپایداری ضرایب رگرسیونی می‌شود. برای اطلاعات بیشتر به اریه (۱۳۹۴) و راپرت (۲۰۰۳) مراجعه کنید.

یک روش مناسب برای حل مشکل فوق، استفاده از پایه‌های  $B$ -اسپلاین به جای پایه‌های توانی بریده‌شده، است که منجر به همان برآورد در بازه  $[t_1, t_k]$  می‌شود. با استفاده از پایه‌های  $B$ -اسپلاین می‌توان چندجمله‌ای‌های متعامدی به دست آورد. در نتیجه، ضرایب رگرسیونی به دست آمده پایدار خواهند بود.

<sup>4</sup>III-Condition

### ۳.ب -B اسپلاین

بازه  $[a, b]$  و تعداد  $q$  گره

$$a = t_0 < t_1 < \dots < t_q < t_{q+1} = b$$

را در نظر بگیرید. پایه‌های  $B$ -اسپلاین، چندجمله‌ای‌های تکه‌ای در بازه  $[a, b]$  هستند و دارای تکیه‌گاه کوچکی می‌باشند. متداول‌ترین روش محاسبه پایه‌های  $B$ -اسپلاین مرتبه  $p$ ، به‌وسیله پایه‌های  $B$ -اسپلاین مرتبه  $p-1$  و بر اساس رابطه بازگشتی زیر است

$$B_{i,p}(x) = \left(\frac{x - t_i}{t_{i+p} - t_i}\right)B_{i,p-1}(x) + \left(\frac{t_{i+p+1} - x}{t_{i+p+1} - t_{i+1}}\right)B_{i+1,p-1}(x), \quad i = -p, -p+1, \dots, q \quad (5.ب)$$

که در آن  $a = t_0 = t_{-1} = t_{-2} = \dots = t_{-p+1} = t_{-p}$  و  $t_{q+1} = b$ . بنابراین تعداد پایه‌های  $B$ -اسپلاین مرتبه  $p$  برابر با  $p+q+1$  است.

#### پایه‌های $B$ -اسپلاین مرتبه صفر

پایه‌های  $B$ -اسپلاین مرتبه صفر ( $p=0$ )، برای هر یک از گره‌ها، به‌صورت زیر تعریف می‌شوند

$$B_{i,0}(x) = \begin{cases} 1 & x \in [t_i, t_{i+1}] \\ 0 & \text{سایر نقاط} \end{cases}$$

که در آن، اندیس صفر بیانگر مرتبه و اندیس  $i$  ( $i = 0, 1, 2, \dots, q$ ) معرف  $i$ -امین تابع متناظر با بازه  $[t_i, t_{i+1}]$  است.

به‌عنوان مثال، این پایه‌ها برای گره‌های  $t_1 = 1, t_2 = 2, t_3 = 3$  در بازه  $[0, 4]$  به‌صورت زیر تعریف می‌شوند

$$B_{0,0}(x) = \begin{cases} 1 & 0 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases} \quad B_{1,0}(x) = \begin{cases} 1 & 1 \leq x < 2 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{2,0}(x) = \begin{cases} 1 & 2 \leq x < 3 \\ 0 & \text{سایر نقاط} \end{cases} \quad B_{3,0}(x) = \begin{cases} 1 & 3 \leq x < 4 \\ 0 & \text{سایر نقاط} \end{cases}$$

### پایه‌های $B$ -اسپلاین مرتبه یک

پایه‌های  $B$ -اسپلاین خطی یا مرتبه یک ( $p = 1$ )، با توجه به رابطه (ب.۵) بر حسب پایه‌های مرتبه صفر به صورت زیر تعریف می‌شوند

$$B_{i,1}(x) = \frac{x - t_i}{t_{i+1} - t_i} B_{i,0}(x) + \frac{t_{i+2} - x}{t_{i+2} - t_{i+1}} B_{i+1,0}(x), \quad i = -1, 0, 1, \dots, q.$$

که در آن،  $t_{-1} = t_0 = a$  و  $t_{q+1} = b$ . همچنین تعداد این پایه‌ها برابر  $q + 2 = p + q + 1$  است. با جایگذاری پایه‌های  $B$ -اسپلاین مرتبه صفر، تابع زیر حاصل می‌شود

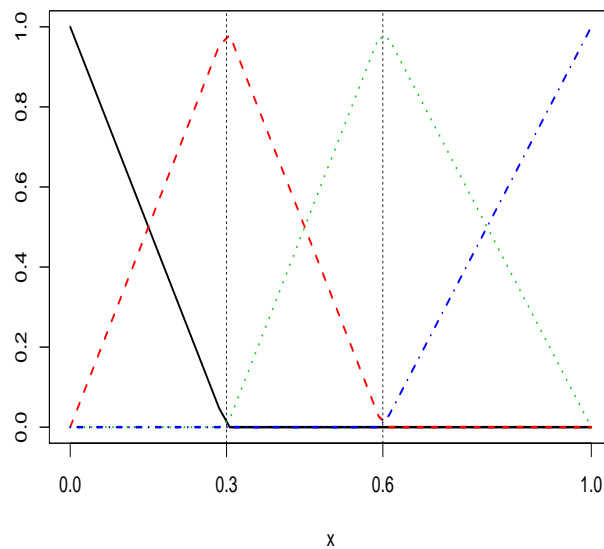
$$B_{i,1}(x) = \begin{cases} \frac{x-t_i}{t_{i+1}-t_i} & t_i \leq x < t_{i+1} \\ \frac{t_{i+2}-x}{t_{i+2}-t_{i+1}} & t_{i+1} \leq x < t_{i+2} \\ 0 & \text{سایر نقاط} \end{cases}$$

به عنوان مثال، اگر در بازه  $[0, 1]$ ، دو گره  $0.3$  و  $0.6$  را در نظر بگیریم، آن‌گاه پایه‌های  $B$ -اسپلاین مرتبه اول، برای  $i = -1, 0, 1, 2, 3$ ، که تعداد آن‌ها برابر  $4 = p + q + 1$  است، به صورت زیر به دست می‌آیند

$$B_{-1,1}(x) = \begin{cases} \frac{0.3-x}{0.3} & 0 \leq x < 0.3 \\ 0 & \text{سایر نقاط} \end{cases} \quad B_{0,1}(x) = \begin{cases} \frac{x}{0.3} & 0 \leq x < 0.3 \\ \frac{0.6-x}{0.3} & 0.3 \leq x < 0.6 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{1,1}(x) = \begin{cases} \frac{x-0.3}{0.3} & 0.3 \leq x < 0.6 \\ \frac{1-x}{0.4} & 0.6 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases} \quad B_{2,1}(x) = \begin{cases} \frac{x-0.6}{0.4} & 0.6 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$

نمودار این پایه‌ها در شکل ب.۱ رسم شده‌اند. دو پایه وسطی دارای شکلی شبیه کلاه (۸) هستند، در حالی که پایه‌های اولی و آخری، کلاه‌های ناقصی هستند. توجه کنید که چون پایه‌های  $B$ -اسپلاین در زیربازه‌ها غیر صفر هستند، لذا برخی از این پایه‌ها دو به دو متعامد هستند.



شکل ب.۱: نمودار پایه‌های  $B$ -اسپلاین مرتبه یک که به ترتیب از سمت چپ نمودارهای  $B_{\circ, 1}(x)$ ،  $B_{-1, 1}(x)$ ،  $B_{1, 1}(x)$  و  $B_{2, 1}(x)$  هستند.

### پایه‌های $B$ -اسپلاین مرتبه دو

پایه‌های  $B$ -اسپلاین مرتبه دو ( $p = 2$ ) بر حسب مرتبه یک تعریف شده و به صورت زیر بیان می‌شوند

$$\begin{aligned}
 B_{i, 2}(x) &= \frac{x - t_i}{t_{i+2} - t_i} B_{i, 1}(x) + \frac{t_{i+3} - x}{t_{i+3} - t_{i+1}} B_{i+1, 1}(x) \\
 &= \frac{x - t_i}{t_{i+2} - t_i} \left[ \frac{x - t_i}{t_{i+1} - t_i} B_{i, \circ}(x) + \frac{t_{i+2} - x}{t_{i+2} - t_{i+1}} B_{i+1, \circ}(x) \right] \\
 &\quad + \frac{t_{i+3} - x}{t_{i+3} - t_{i+1}} \left[ \frac{x - t_{i+1}}{t_{i+2} - t_{i+1}} B_{i+1, \circ}(x) + \frac{t_{i+3} - x}{t_{i+3} - t_{i+2}} B_{i+2, \circ}(x) \right], \quad i = -2, -1, \circ, 1, \dots, k
 \end{aligned}$$

که در آن  $t_{-2} = t_{-1} = t_{\circ} = a$  و  $t_{q+1} = b$ . با جایگذاری پایه‌های  $B$ -اسپلاین مرتبه اول، تابع زیر نتیجه می‌شود

$$B_{i, 2} = \begin{cases} \frac{(x-t_i)^2}{(t_{i+2}-t_i)(t_{i+1}-t_i)} & t_i \leq x < t_{i+1} \\ \frac{(x-t_i)(t_{i+2}-x)}{(t_{i+2}-t_i)(t_{i+2}-t_{i+1})} + \frac{(t_{i+3}-x)(x-t_{i+1})}{(t_{i+3}-t_{i+1})(t_{i+2}-t_{i+1})} & t_{i+1} \leq x < t_{i+2} \\ \frac{(t_{i+3}-x)^2}{(t_{i+3}-t_{i+1})(t_{i+3}-t_{i+2})} & t_{i+2} \leq x < t_{i+3} \\ \circ & \text{سایر نقاط} \end{cases}$$

به عنوان مثال، در بازه  $[0, 1]$  سه گره  $0, 1/3, 2/3, 1$  را در نظر بگیرید. تعداد پایه های  $B$  - اسپلاین مرتبه دوم با سه گره، برابر  $6 = p + q + 1$  است که به صورت زیر به دست می آید

$$B_{-2,2}(x) = \begin{cases} (1/3 - x)^2 & 0 \leq x < 1/3 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{-1,2}(x) = \begin{cases} \frac{(x)(1/3-x)}{(1/3)(1/3)} + \frac{(2/3-x)(x)}{(2/3)(1/3)} & 0 \leq x < 1/3 \\ \frac{(2/3-x)^2}{(2/3)(1/3)} & 1/3 \leq x < 2/3 \\ 0 & \text{سایر نقاط} \end{cases}$$

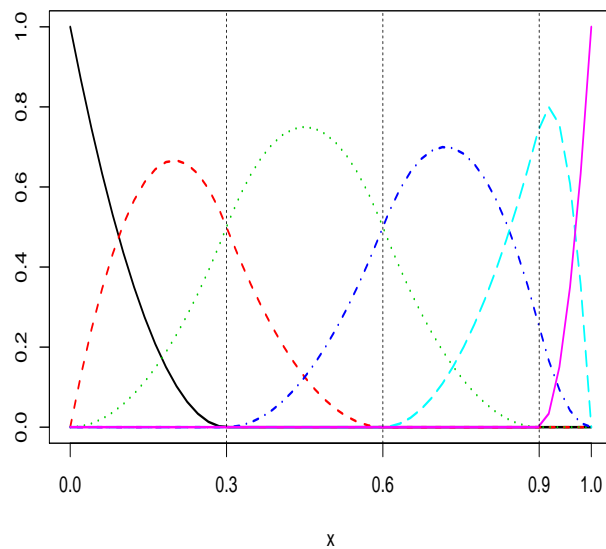
$$B_{0,2}(x) = \begin{cases} \frac{x^2}{(1/3)(1/3)} & 0 \leq x < 1/3 \\ \frac{(x)(2/3-x)}{(2/3)(1/3)} + \frac{(1-x)(x-1/3)}{(1/3)(1/3)} & 1/3 \leq x < 2/3 \\ \frac{(1-x)^2}{(2/3)(1/3)} & 2/3 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{1,2}(x) = \begin{cases} \frac{(x-1/3)^2}{(2/3)(1/3)} & 1/3 \leq x < 2/3 \\ \frac{(x-1/3)(1-x)}{(2/3)(1/3)} + \frac{(1-x)(x-2/3)}{(1/3)(1/3)} & 2/3 \leq x < 1 \\ \frac{(1-x)^2}{(1/3)(1/3)} & 0 \leq x < 1/3 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{2,2}(x) = \begin{cases} \frac{(x-2/3)^2}{(1/3)(1/3)} & 2/3 \leq x < 1 \\ \frac{(x-2/3)(1-x)}{(1/3)(1/3)} + \frac{(1-x)(x-1/3)}{(1/3)(1/3)} & 0 \leq x < 2/3 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{3,2}(x) = \begin{cases} \frac{(x-1/3)^2}{(1/3)(1/3)} & 0 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$





شکل ب.۲: نمودار پایه‌های  $B$ -اسپلاین مرتبه‌ی دو که به ترتیب از سمت چپ نمودارهای  $B_{-۲,۲}(x)$ ،  $B_{-۱,۲}(x)$ ،  $B_{۰,۲}(x)$ ،  $B_{۱,۲}(x)$ ،  $B_{۲,۲}(x)$  و  $B_{۳,۲}(x)$  هستند.

### پایه‌های $B$ -اسپلاین مرتبه‌ی سه

پایه‌های  $B$ -اسپلاین مرتبه سه ( $p = ۳$ ) به صورت سلسله‌مراتبی بر حسب درجات پایین‌تر تعریف شده و به صورت زیر بیان می‌شوند

$$\begin{aligned}
 B_{i,۳}(x) &= \frac{(x - t_i)}{(t_{i+۳} - t_i)} B_{i,۲}(x) + \frac{(t_{i+۴} - x)}{(t_{i+۴} - t_{i+۱})} B_{i+۱,۲}(x) \\
 &= \frac{(x - t_i)}{(t_{i+۳} - t_i)} \left[ \frac{x - t_i}{t_{i+۲} - t_i} B_{i,۱}(x) + \frac{t_{i+۳} - x}{t_{i+۳} - t_{i+۱}} B_{i+۱,۱}(x) \right] \\
 &+ \frac{(t_{i+۴} - x)}{(t_{i+۴} - t_{i+۱})} \left[ \frac{x - t_{i+۱}}{t_{i+۳} - t_{i+۱}} B_{i+۱,۱}(x) + \frac{t_{i+۴} - x}{t_{i+۴} - t_{i+۲}} B_{i+۲,۱}(x) \right] \\
 &= \frac{(x - t_i)}{(t_{i+۳} - t_i)} \left[ \frac{x - t_i}{t_{i+۲} - t_i} \left( \frac{x - t_i}{t_{i+۱} - t_i} B_{i,۰}(x) + \frac{t_{i+۲} - x}{t_{i+۲} - t_{i+۱}} B_{i+۱,۰}(x) \right) \right. \\
 &+ \left. \frac{t_{i+۳} - x}{t_{i+۳} - t_{i+۱}} \left( \frac{x - t_{i+۱}}{t_{i+۲} - t_{i+۱}} B_{i+۱,۰}(x) + \frac{t_{i+۳} - x}{t_{i+۳} - t_{i+۲}} B_{i+۲,۰}(x) \right) \right] \\
 &+ \frac{(t_{i+۴} - x)}{(t_{i+۴} - t_{i+۱})} \left[ \frac{x - t_{i+۱}}{t_{i+۳} - t_{i+۱}} \left( \frac{x - t_{i+۱}}{t_{i+۲} - t_{i+۱}} B_{i+۱,۰}(x) + \frac{t_{i+۳} - x}{t_{i+۳} - t_{i+۲}} B_{i+۲,۰}(x) \right) \right. \\
 &+ \left. \frac{t_{i+۴} - x}{t_{i+۴} - t_{i+۲}} \left( \frac{x - t_{i+۲}}{t_{i+۳} - t_{i+۲}} B_{i+۲,۰}(x) + \frac{t_{i+۴} - x}{t_{i+۴} - t_{i+۳}} B_{i+۳,۰}(x) \right) \right]
 \end{aligned}$$

$$B_{i,3}(x) = \begin{cases} \frac{(t_{i+4}-x)^3}{(t_{i+4}-t_{i+1})(t_{i+4}-t_{i+2})(t_{i+4}-t_{i+3})} & t_i \leq x < t_{i+1} \\ \frac{(x-t_i)^2(t_{i+2}-x)}{(t_{i+3}-t_i)(t_{i+2}-t_i)(t_{i+2}-t_{i+1})} + \frac{(x-t_i)(t_{i+3}-x)(x-t_{i+1})}{(t_{i+3}-t_i)(t_{i+3}-t_{i+1})(t_{i+2}-t_{i+1})} \\ + \frac{(t_{i+4}-x)(x-t_{i+1})^2}{(t_{i+4}-t_{i+1})(t_{i+3}-t_{i+1})(t_{i+2}-t_{i+1})} & t_{i+1} \leq x < t_{i+2} \\ \frac{(x-t_i)(t_{i+3}-x)^2}{(t_{i+3}-t_i)(t_{i+3}-t_{i+1})(t_{i+3}-t_{i+2})} + \frac{(t_{i+4}-x)(x-t_{i+1})(t_{i+3}-x)}{(t_{i+4}-t_{i+1})(t_{i+3}-t_{i+1})(t_{i+3}-t_{i+2})} \\ + \frac{(t_{i+4}-x)^2(x-t_{i+2})}{(t_{i+4}-t_{i+1})(t_{i+4}-t_{i+2})(t_{i+3}-t_{i+2})} & t_{i+2} \leq x < t_{i+3} \\ \frac{(x-t_i)^3}{(t_{i+3}-t_i)(t_{i+2}-t_i)(t_{i+1}-t_i)} & t_{i+3} \leq x < t_{i+4} \\ 0 & \text{سایر نقاط} \end{cases}$$

به عنوان مثال، در بازه  $[0, 1]$  چهار گره  $0, 1/2, 2/3, 1$  را در نظر بگیرید. تعداد پایه‌های  $B$  - اسپلاین مرتبه سه با چهار گره برابر  $1 + 2 + 3 = 6$  است که به صورت زیر به دست می‌آیند

$$B_{-3,3}(x) = \begin{cases} \frac{(0.5-x)^3}{(0.5)(0.5)(0.5)} & 0 \leq x < 0.5 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{-2,3}(x) = \begin{cases} \frac{(0.4-x)^3}{(0.4)(0.4)(0.4)} & 0 \leq x < 0.4 \\ \frac{(x)(0.5-x)^2}{(0.5)(0.5)(0.5)} + \frac{(0.4-x)(x)(0.5-x)}{(0.4)(0.4)(0.5)} + \frac{(0.4-x)^2(x)}{(0.4)(0.4)(0.4)} & 0.4 \leq x < 0.5 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{-1,3}(x) = \begin{cases} \frac{(0.6-x)^3}{(0.6)(0.4)(0.2)} & 0 \leq x < 0.2 \\ \frac{(x)^2(0.5-x)}{(0.4)(0.5)(0.5)} + \frac{(x)(0.4-x)(x)}{(0.4)(0.4)(0.2)} + \frac{(0.6-x)(x)^2}{(0.6)(0.4)(0.5)} & 0.2 \leq x < 0.4 \\ \frac{(x)(0.4-x)^2}{(0.4)(0.4)(0.4)} + \frac{(0.6-x)(x)(0.4-x)}{(0.6)(0.4)(0.4)} + \frac{(0.6-x)^2(x-0.5)}{(0.6)(0.4)(0.2)} & 0.4 \leq x < 0.6 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{0,3}(x) = \begin{cases} \frac{x^3}{(0.6)(0.4)(0.2)} & 0 \leq x < 0.2 \\ \frac{(x^3)(0.4-x)}{(0.6)(0.4)(0.2)} + \frac{(x)(0.6-x)(x-0.2)}{(0.6)(0.4)(0.2)} + \frac{(0.8-x)(x-0.2)^2}{(0.6)(0.4)(0.2)} & 0.2 \leq x < 0.4 \\ \frac{(x)(0.6-x)^2}{(0.6)(0.4)(0.2)} + \frac{(0.8-x)(x-0.2)(0.6-x)}{(0.6)(0.4)(0.2)} + \frac{(0.8-x)^2(x-0.4)}{(0.6)(0.4)(0.2)} & 0.4 \leq x < 0.6 \\ \frac{(0.8-x)^3}{(0.6)(0.4)(0.2)} & 0.6 \leq x < 0.8 \\ 0 & \text{سایر نقاط} \end{cases}$$

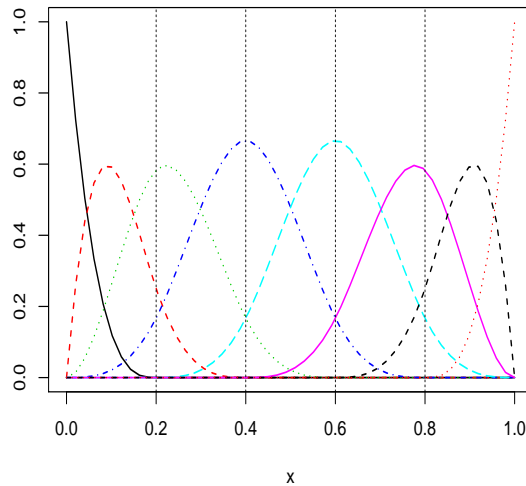
$$B_{1,3}(x) = \begin{cases} \frac{(1-x)^3}{(0.6)(0.4)(0.2)} & 0.2 \leq x < 0.4 \\ \frac{(x-0.2)^2(0.6-x)}{(0.6)(0.4)(0.2)} + \frac{(x-0.2)(0.8-x)(x-0.4)}{(0.6)(0.4)(0.2)} + \frac{(1-x)(x-0.4)^2}{(0.6)(0.4)(0.2)} & 0.4 \leq x < 0.6 \\ \frac{(x-0.2)(0.8-x)^2}{(0.6)(0.4)(0.2)} + \frac{(1-x)(x-0.4)(0.8-x)}{(0.6)(0.4)(0.2)} + \frac{(1-x)^2(x-0.6)}{(0.6)(0.4)(0.2)} & 0.6 \leq x < 0.8 \\ \frac{(x-0.2)^3}{(0.6)(0.4)(0.2)} & 0.8 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{2,3}(x) = \begin{cases} \frac{(x-0.4)^2(0.8-x)}{(0.6)(0.4)(0.2)} + \frac{(x-0.4)(1-x)(x-0.6)}{(0.6)(0.4)(0.2)} + \frac{(1-x)(x-0.6)^2}{(0.6)(0.4)(0.2)} & 0.4 \leq x < 0.6 \\ \frac{(x-0.4)(1-x)^2}{(0.6)(0.4)(0.2)} + \frac{(1-x)(x-0.6)(1-x)}{(0.4)(0.4)(0.2)} + \frac{(1-x)^2(x-0.8)}{(0.4)(0.2)(0.2)} & 0.6 \leq x < 0.8 \\ \frac{(x-0.4)^3}{(0.6)(0.4)(0.2)} & 0.8 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{3,3}(x) = \begin{cases} \frac{(x-0.6)^3}{(0.4)(0.4)(0.2)} & 0.6 \leq x < 0.8 \\ \frac{(x-0.6)^2(1-x)}{(0.4)(0.4)(0.2)} + \frac{(x-0.6)(1-x)(x-0.8)}{(0.4)(0.2)(0.2)} + \frac{(1-x)(x-0.8)}{(0.2)(0.2)(0.2)} & 0.8 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$

$$B_{4,3}(x) = \begin{cases} \frac{(x-0.8)^3}{(0.2)(0.2)(0.2)} & 0.8 \leq x < 1 \\ 0 & \text{سایر نقاط} \end{cases}$$

این توابع از جنس چندجمله‌ای‌های تکه‌ای مرتبه سه هستند و نمودارهای این توابع در شکل ۳.ب رسم شده‌اند.



شکل ب.۳: نمودار پایه‌های  $B$ -اسپلاین مرتبه سه که به ترتیب از سمت چپ نمودارهای  $B_{۰,۳}(x)$ ،  $B_{۱,۳}(x)$ ،  $B_{۲,۳}(x)$ ،  $B_{۳,۳}(x)$ ،  $B_{۴,۳}(x)$  و  $B_{۵,۳}(x)$  هستند.

## ب.۴ کمترین توان‌های دوم $B$ -اسپلاین

ابتدا با  $۱ + p + q$  پایه  $B$ -اسپلاین مرتبه  $p$ ، رابطه (ب.۵) را با اندیس‌گذاری مجدد، به صورت

$$B_{۰,p}, B_{۱,p}, \dots, B_{p+q,p}$$

در نظر می‌گیریم. مدل رگرسیونی  $B$ -اسپلاین به صورت زیر است

$$s(x) = \sum_{j=1}^{p+q} \beta_j B_{j,p}(x) + \epsilon. \quad (\text{ب.۶})$$

از آن جا که صورت ماتریسی این مدل به صورت

$$s = B\beta + \epsilon$$

است، لذا برآورد کمترین توان‌های دوم  $\beta$  عبارت است از

$$\hat{\beta} = (B'B)^{-1} B'y.$$

متعامد بودن پایه‌های  $B$ -اسپلاین موجب می‌شود که ماتریس  $B'B$  وارون پذیر باشد و لذا مشکل موجود در رابطه (ب.۴) مرتفع گردد. لازم به ذکر است که علیرغم به دست آوردن این مزیت، درک پایه‌های  $B$ -اسپلاین به آسانی درک توابع توانی بریده شده نیست و لذا ضرایب مدل مذکور غیر قابل تفسیر است. برای اطلاعات بیشتر در این باره می‌توانید به دی بور (۲۰۰۱)، راپرت (۲۰۰۳) و دایرکس<sup>۵</sup> (۱۹۹۵) مراجعه کنید.

<sup>5</sup>Dierckx

# پیوست

## تعاریف و نامساوی‌ها

**تعریف پ.۱.۰** (نامساوی دژرتزکی-کیفر-ولفوویتز). فرض کنید  $X_1, \dots, X_n$  یک نمونه تصادفی از تابع توزیع  $F(\cdot)$  بوده و  $F_n(\cdot)$  تابع توزیع تجربی متناظر باشد. این نامساوی یک کران بالا برای احتمال اینکه اختلاف بین  $F_n(\cdot)$  و  $F(\cdot)$  از یک مقدار ثابت داده شده  $\varepsilon > 0$  بزرگتر باشد، نتیجه می‌دهد. به‌طور دقیق‌تر، به ازای هر  $\varepsilon > 0$  داریم

$$P\left(\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \geq \varepsilon\right) \leq 2 \exp\{-2n\varepsilon^2\}.$$

**تعریف پ.۲.۰** (نامساوی مارکوف). اگر  $X$  یک متغیر تصادفی نامنفی باشد، آنگاه به‌ازای هر  $a > 0$  داریم

$$P(X \geq a) \leq \frac{E(X)}{a}.$$

**تعریف پ.۳.۰** (تجزیه چولسکی). هر ماتریس متقارن معین مثبت  $\mathbf{A}$  را می‌توان به‌صورت زیر نوشت

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T$$

که  $\mathbf{L}$  یک ماتریس پایین‌مثلثی و  $\mathbf{L}^T$  ترانپوز آن است. لازم به ذکر است که تجزیه چولسکی یکتاست.

به‌عنوان مثال، در یک ماتریس  $3 \times 3$  باید دستگاه معادلات زیر حل شود:

$$\begin{aligned} \mathbf{A} &= \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \\ &= \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix} \begin{pmatrix} l_{11} & l_{21} & l_{31} \\ 0 & l_{22} & l_{32} \\ 0 & 0 & l_{33} \end{pmatrix} = \mathbf{L}\mathbf{L}^T \\ &= \begin{pmatrix} l_{11}^2 & l_{21}l_{11} & l_{31}l_{11} \\ l_{21}l_{11} & l_{21}^2 + l_{22}^2 & l_{31}l_{21} + l_{32}l_{22} \\ l_{31}l_{11} & l_{31}l_{21} + l_{32}l_{22} & l_{31}^2 + l_{32}^2 + l_{33}^2 \end{pmatrix}. \end{aligned}$$

می‌توان نشان داد که برای اعضای قطر اصلی  $(l_{kk})$  و اعضای زیر قطر اصلی  $(l_{ik}, i > k)$  روابط زیر برقرار است

$$l_{kk} = a_k - \sum_{j=1}^{k-1} a_{kj}^2, \quad l_{ik} = \frac{1}{l_{kk}} \left( a_{ik} - \sum_{j=1}^{k-1} l_{ij}l_{kj} \right).$$

**تعریف پ.۴.۰ (پایه).** مجموعه بردارهای  $\mathbf{u}_1, \dots, \mathbf{u}_n$  در فضای برداری  $\mathbf{V}$  تشکیل یک پایه می‌دهند، اگر دو شرط زیر را داشته باشند:

۱. بردارهای  $\mathbf{u}_1, \dots, \mathbf{u}_n$  مستقل خطی باشند.
۲. هر کدام از بردارهای فضای  $\mathbf{V}$  را بتوان از ترکیب خطی این بردارها به دست آورد.

**تعریف پ.۵.۰ (همگرایی در توزیع).** دنباله متغیرهای تصادفی  $\{X_n\}$  در توزیع به متغیر تصادفی  $X$  همگراست اگر برای هر  $x \in \mathbb{R}$  که تابع  $F$  در آن پیوسته است، داشته باشیم

$$\lim_{n \rightarrow \infty} F_n(x) = F(x)$$

که در آن  $F_n$  تابع توزیع تجمعی  $\{X_n\}$  و  $F$  تابع توزیع تجمعی  $X$  است. این نوع همگرایی را به‌صورت  $X_n \xrightarrow{D} X$  نشان می‌دهند.

**تعریف پ.۶.۰ (همگرایی در احتمال).** دنباله متغیرهای تصادفی  $\{X_n\}$  در احتمال به متغیر تصادفی  $X$  همگراست اگر برای هر  $\varepsilon > 0$  داشته باشیم

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \varepsilon) = 0.$$

این نوع همگرایی را با نماد  $X_n \xrightarrow{P} X$  نشان می‌دهند.

**تعریف پ.۷.۰ (سازگاری).** دنباله برآوردگرهای  $\{T_n\}$  برای  $a(\theta)$  سازگار است، اگر  $\{T_n\}$  در احتمال به  $a(\theta)$  همگرا باشد.

**تعریف پ. ۸.۰.** گوییم  $f(n) = O(g(n))$  از مرتبه  $g(n)$  است، اگر وقتی  $n \rightarrow \infty$ ، آنگاه  $\frac{f(n)}{g(n)}$  متناهی باشد. به عبارت دیگر،  $f(n) = O(g(n))$ ، اگر عدد حقیقی مثبت  $M$  و عدد حقیقی  $n_0$  موجود باشند، به طوری که

$$\forall n > n_0. \quad |f(n)| \leq M|g(n)|$$

یعنی برای مقادیر به قدر کافی بزرگ  $n$ ، قدرمطلق  $f(n)$  حداکثر  $M$  برابر قدرمطلق  $g(n)$  باشد.

**تعریف پ. ۹.۰.** گوییم  $f(n) = o(g(n))$  از مرتبه کوچکتر از  $g(n)$  است، اگر وقتی  $n \rightarrow \infty$ ، آنگاه داشته باشیم

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0.$$

عبارت  $f(n) = o(g(n))$  بیان می کند که  $g(n)$  بسیار سریع تر از  $f(n)$  رشد می کند. بنابراین  $o$  حالت قوی تری از  $O$  را بیان می کند. هر عبارتی که  $o(g)$  باشد،  $O(g)$  هم هست، اما عکس آن برقرار نیست.

**تعریف پ. ۱۰.۰.** فرض کنید  $\{X_n\}$  یک دنباله از متغیرهای تصادفی و  $\{a_n\}$  یک دنباله از ثابت های متناظر باشد. نماد  $X_n = o_p(a(n))$  بدین معنی است که مجموعه مقادیر  $\frac{X_n}{a_n}$  در احتمال به صفر همگرا است، یعنی برای هر  $\varepsilon > 0$  داریم

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{X_n}{a_n}\right| > \varepsilon\right) = 0.$$

**تعریف پ. ۱۱.۰.** فرض کنید  $\{X_n\}$  یک دنباله از متغیرهای تصادفی و  $\{a_n\}$  یک دنباله از ثابت های متناظر باشد. نماد  $X_n = O_p(a(n))$  بدین معنی است که مقادیر  $\frac{X_n}{a_n}$  در احتمال کراندار است، یعنی برای هر  $\varepsilon > 0$ ، اعداد متناهی  $M, N > 0$  وجود دارند، به طوری که

$$\forall n > N \quad P\left(\left|\frac{X_n}{a_n}\right| > M\right) < \varepsilon.$$

**تعریف پ. ۱۲.۰.** گوییم برآوردگر  $\hat{\theta}_n$  یک برآوردگر  $\sqrt{n}$  - سازگار برای  $\theta$  است، اگر رابطه زیر برقرار باشد

$$\hat{\theta}_n - \theta = O_p(n^{-1/2}).$$

**تعریف پ. ۱۳.۰.** (بسط تیلور). فرض کنید تابع  $f(x)$  بر بازه  $[a, b]$  تعریف شده و  $n + 1$  بار مشتق پذیر باشد. نقطه  $x_0 \in (a, b)$  را در نظر بگیرید. برای هر  $x \in (a, b)$  که  $x = x_0 + \Delta x$ ، عددی مانند  $\xi$  بین دو نقطه  $x_0$  و  $x$  وجود دارد، به طوری که تساوی زیر برقرار است

$$f(x_0 + \Delta x) = f(x_0) + \Delta x f'(x_0) + \frac{\Delta x^2}{2!} f''(x_0) + \dots + \frac{\Delta x^n}{n!} f^{(n)}(x_0) + R_n(x_0 + \Delta x)$$

که

$$R_n(x_0 + \Delta x) = \frac{\Delta x^{n+1}}{(n+1)!} f^{(n+1)}(\xi), \quad x_0 \leq \xi \leq x.$$

منظور از  $f^{(n)}$  در این روابط مشتق مرتبه  $n$ ام می باشد. در عمل معمولاً از تعداد اندکی از جملات اول سری برای تقریب استفاده می شود. در این حالت، اگر مقدار  $\Delta x$  به قدر کافی کوچک باشد، تقریب

مناسبی از تابع در همسایگی  $x_0$  به دست خواهد آمد. در صورتی که تقریب تابع با سری تیلور برای  $n$  جمله اول انجام شود، مانده سری یا خطای تقریب متناسب با  $\Delta x^n$  خواهد بود و آن را با  $O(\Delta x^n)$  نشان می‌دهند

$$f(x_0 + \Delta x) = f(x_0) + \Delta x f'(x_0) + \frac{\Delta x^2}{2!} f''(x_0) + \dots + \frac{\Delta x^{n-1}}{(n-1)!} f^{(n-1)}(x_0) + O(\Delta x^n).$$



# پیوست

## گزیده‌ای از برنامه‌های رایانه‌ای

### جدول ۱.۱

```
library(scatterplot3d)
library(AER)
data("CPS1985")
fit <- lm(log(wage)~education+experience+I(experience^2),data=CPS1985)
summary(fit)
```

### شکل ۱.۱

```
b <- coef(fit)
m.school <- data.frame(CPS1985$education,CPS1985$education*b["education"])
m.school <- m.school[order(CPS1985$education),]
m.experi <- data.frame(CPS1985$experience,
CPS1985$experience*b["experience"]
+CPS1985$experience^2*b["I(experience^2)"])
m.experi <- m.experi[order(CPS1985$experience),]
par(mfrow=c(1,2))
par(mar = c(bottom=4, 2, top=2, 2))
```

```

plot(m.school,col="black",pch=19,cex=0.7,
xaxt="n",yaxt="n",xlab="school",ylab="")
axis(1,cex.axis=0.8)
axis(2,cex.axis=0.8)
lines(m.school,col="azure3")
plot(m.experi,col="black",pch=19,cex=0.7,xaxt="n",yaxt="n",
xlab="experience",ylab="")
axis(1,cex.axis=0.8)
axis(2,cex.axis=0.8)
lines(m.experi,col="azure3")
col <- rep("green",nrow(CPS1985)); col[CPS1985$education==12] <- "red"
pch <- rep(1,nrow(CPS1985)); pch[CPS1985$education==12] <- 8

```

## شکل ۲.۱

```

scatterplot3d(CPS1985$education,CPS1985$experience, fitted(fit),
xlab="school",ylab="experience",zlab="",
color=col,pch=pch,angle=130)

```

## شکل ۳.۱

```

library(np)
require(graphics)

education<- CPS1985$education
experience<-CPS1985$experience
X <- cbind(education,experience)
y <- log(CPS1985$wage)
bw <- npregbw(xdat=X, ydat=y)
model <- npreg(bws = bw, gradients = TRUE)
summary(model)

col <- rep("green",nrow(CPS1985)); col[CPS1985$education==12] <- "red"
pch <- rep(1,nrow(CPS1985)); pch[education==12] <- 8
scatterplot3d(education,experience, model$mean,
xlab="education",ylab="experience",zlab="",
color=col,pch=pch,angle=130)

```

```

library(mgcv)
b <- gam(y ~s(education) + s(experience))
plot(b,pages = 1)
fit3 <- fitted(b)
scatterplot3d(education,experience, fit3,
xlab="education",ylab="experience",zlab="",
color=col,pch=pch,angle=130)

```

```

S <- function(xx,l){
  (xx - sign(xx))*(abs(xx) > l)
}
scad <- function(theta,l,a=3.7){
  theta <- abs(theta)
  (theta <= l) * theta * l + (theta > l & theta <= a*l)*
  (a*l*theta-(theta^2+l^2)/2)/(a-1)+(theta>a*l)*(l^2*(a^2-1))/(2*(a-1))
}
Fs <- function(z,l,g=3.7){
  (abs(z)>g*l)*z+(abs(z)>2*l & abs(z)<=g*l)*(S(z,g*l/(g-1))/(1-1/(g-1)))+
  (abs(z) <= 2*l)*S(z,l)
}
mcp <- function(theta,l,a=3.7) {
  T <- length(theta)
  val <- numeric(T)
  for (i in 1:T) {
    x <- abs(theta[i])
    val[i] <- (x < a*l)*(1*x - x^2/(2*a)) + (x >= a*l)*(1/2)*a*l^2
  }
  val
}
Fm <- function(z,l,g=3.7){
  (abs(z) > g*l)*z + (abs(z) <= g*l)*(S(z,l)/(1-1/g))
}

```

```

}
lasso <- function(theta,l){l*abs(theta)}

res <- 101
x <- seq(-4,4,len=res)
xx <- seq(0.0001,4,len=res)
g <- 3
col <- c("#FF4E37FF", "#00B500FF", "#008DFFFF")
Y <- cbind(lasso(x, 1), scad(x, 1, g), mcp(x, 1, g))

require(latex2exp)
par(mar = c(bottom=4, 5, top=1, 2))
matplot(x,Y,type='l',xaxt="n",yaxt="n",lty=1:3,lwd=c(1,1,1.8),col=col,las=1,
xlab=expression(beta), ylab=TeX("$P_{\\lambda}(1\\beta 1)$"),cex.lab=0.85)
axis(1,cex.axis=0.7)
axis(2,cex.axis=0.7)
legend(-1.70, 4, legend=c("LASSO", "SCAD","MCP"),
col=c("red","green", "blue"), lty=1:3,lwd=c(1,1,1.8), cex=0.7)

```

## شکل ۲.۲

```

par(mfrow=c(2,2))
par(mar = c(bottom=2.2, 3, top=2, 2))
Y <- cbind(S(x,1), Fm(x,1,g), Fs(x,1,g))
plot(x,S(x,1),type='l',xaxt="n",yaxt="n",xlim=c(-4,4),
ylim=c(-4,4),lwd=2.5,xlab="",ylab="",main="LASSO",cex.main=0.8)
lines(x,x,col="red")
axis(1,cex.axis=0.7)
axis(2,cex.axis=0.7)
par(mar = c(bottom=2.2, 3, top=2, 2))
plot(x,0.5*x,type = 'l',xaxt="n",yaxt="n", xlim=c(-4,4), ylim=c(-4,4),
lwd=2.5,xlab="",ylab="",main="Ridge",cex.main=0.8)
lines(x,x,col="red")
axis(1,cex.axis=0.7)
axis(2,cex.axis=0.7)

```

```

plot(x,Fs(x,1,g),type='l',xaxt="n",yaxt="n",xlim=c(-4,4),ylim=c(-4,4),
lwd=2.5,xlab="",ylab="",main = "SCAD",cex.main=0.8)
lines(x,x,col="red")
axis(1,cex.axis=0.7)
axis(2,cex.axis=0.7)
plot(x,Fm(x,1,g),type='l',xaxt="n",yaxt="n",xlim=c(-4,4),ylim=c(-4,4),
lwd=2.5,xlab="",ylab="",main="MCP",cex.main=0.8)
lines(x,x,col="red")
axis(1,cex.axis=0.7)
axis(2,cex.axis=0.7)

```

## شکل ۴.۲

```

dLasso <- function(theta,l){rep(l,length(theta))}
dMCP <- function(theta,l,a=3.7) {
theta = abs(theta)
(theta <= a*l) * (1 - theta/a)}
dSCAD <- function(theta,l,a=3.7) {
theta = abs(theta)
(theta<l)*1+((theta>l)&(theta<a*l))*((a*l-theta)/(a-1))
}
Y<- cbind(dLasso(xx, 1), dSCAD(xx, 1, g), dMCP(xx, 1, g))
matplot(xx,Y,type='l',xaxt="n",yaxt="n",lty=1:3,lwd=c(1,1,1.8),
col=col, las=1,cex.lab=0.85, xlab=expression(beta),
ylab=TeX("$P'_{\lambda}(1\beta 1)$"))
axis(1,cex.axis=0.7)
axis(2,cex.axis=0.7)
legend(2.5, 0.9, legend=c("LASSO", "SCAD","MCP"),
col=c("red","green","blue"), lty=1:3, lwd=c(1,1,1.8),cex=0.7)

```

## شکل ۵.۲

```

library(splines)
library(grpreg)

n <- 2000; p <- 20; K < -2

```

```
X <- matrix(0,n,p)
for(j in 1:p ){
X[,j] <- runif(n)
}

obj0 <- function(xx,K)
{
n <- nrow(xx);p <- ncol(xx)
z <- c()
for (j in 1:p){
B1 <- splineDesign(knots=c(0,0,0,seq(0,1,length=K+2),1,1,1),x=xx[,j],
derivs=rep(0,n))
B1 <- B1[,1:(K+1)]
z <- cbind(z,B1)}
z
}
Z <- -obj0(X,K)
Z <- Z[,-3]
group<-c(c(1,1),rep(2:p,each=K+1))
Y <- -sqrt(2)*Z[,1]+sqrt(2)*Z[,2]+0.5*Z[,3]+Z[,4]-0.5*Z[,5]+rnorm(n,0,1)
fit1 <- grpreg(Z, Y, group, penalty="grLasso")
fit2 <- grpreg(Z, Y, group, penalty="grSCAD")
fit3 <- grpreg(Z, Y, group, penalty="grMCP",gamma=2.7)
par(mfrow=c(1,3))
myColors<-c("blue","red","green","black","yellow","purple","orange","brown")
plot(fit1,col=myColors,main="Group LASSO", xaxt="n",yaxt="n",cex.main=1)
axis(1,cex.axis=1.2)
axis(2,cex.axis=1.2)
abline(v=.04,col="black",lty=2,lwd=0.75)
plot(fit2,col=myColors,main="Group SCAD", xaxt="n",yaxt="n",cex.main=1)
axis(1,cex.axis=1.2)
axis(2,cex.axis=1.2)
abline(v=.04,col="black",lty=2,lwd=0.75)
plot(fit3,col=myColors,main="Group MCP", xaxt="n",yaxt="n",cex.main=1)
```

```
axis(1,cex.axis=1.2)
axis(2,cex.axis=1.2)
abline(v=.04,col="black",lty=2,lwd=0.75)
```

## شکل ۸.۲

```
res <- 201
x <- seq(-10,10,len=res)
xx <- seq(0.0001,4,len=res)
g <- 3
col <- c("#FF4E37FF", "#00B500FF", "#008DFFFF")
lambda <- 2
y <- scad(x, lambda) # a=3.7
#### SCAD
require(latex2exp)
par(mfrow=c(2,2))
par(mar = c(bottom=3, 4, top=2, 3))
# lambda=2,beta0=4
matplot(x, y,lwd=1.3, type='l', ylim=c(0,20), main=TeX("SCAD:
 $\beta^{(0)}=4$ "), xaxt="n", yaxt="n", cex.main=1.5, lty=1,
las=1, xlab=expression(beta), ylab=TeX("$P_{\lambda}(\beta)$"))
axis(1,cex.axis=0.9)
axis(2,cex.axis=0.9)
beta0<-4
LLAapprox<-function(Beta){
scad(beta0,lambda)+dSCAD(beta0,lambda)*(abs(Beta)-beta0)
}
LQAapprox<-function(Beta){
scad(beta0,lambda)+0.5*(dSCAD(beta0,lambda)/beta0)*(Beta^2-beta0^2)
}
lines(x,LQAapprox(x),lty=2,lwd=1.4,col="blue")
lines(x,LLAapprox(x),lty=3,lwd=1.4,col="red")
# lambda=2,beta0=1
par(mar = c(bottom=3, 4, top=2, 3))
matplot(x, y, type='l', ylim=c(0,20), main=TeX("SCAD:  $\beta^{(0)}=1$ "),
```

```

cex.main=1.5, xaxt="n", yaxt="n", lty=1, lwd=1.3, las=1,
xlab=expression(beta), ylab=TeX("$P_{\\lambda}(\\beta)$"))
axis(1,cex.axis=0.9)
axis(2,cex.axis=0.9)
beta0<-1
LLAapprox<-function(Beta){
scad(beta0,lambda)+dSCAD(beta0,lambda)*(abs(Beta)-beta0)
}
LQAapprox<-function(Beta){
scad(beta0,lambda)+0.5*(dSCAD(beta0,lambda)/beta0)*(Beta^2-beta0^2)
}
lines(x,LQAapprox(x),lty=2,lwd=1.4,col="blue")
lines(x,LLAapprox(x),lty=3,lwd=1.4,col="red")
##### Lq:q=0.5
# lambda=2,beta0=4
Lq <- function(theta,l){l*abs(theta)^0.5}
dLq <- function(theta,l){(1/2)*theta^(-0.5)}
lambda<-2
yLq <- Lq(x, lambda)
par(mar = c(bottom=3, 4, top=2, 3))
matplot(x, yLq,lwd=1.3, type='l', ylim=c(0,10), main=TeX("L$_{0.5}$:
\\beta^{(0)}$=4"), xaxt="n",yaxt="n", cex.main=1.5, lty=1, las=1,
xlab=expression(beta), ylab=TeX("$P_{\\lambda}(\\beta)$"))
axis(1,cex.axis=0.9)
axis(2,cex.axis=0.9)
beta0<-4
LLAapprox<-function(Beta){
Lq(beta0,lambda)+dLq(beta0,lambda)*(abs(Beta)-beta0)
}
LQAapprox<-function(Beta){
Lq(beta0,lambda)+0.5*(dLq(beta0,lambda)/beta0)*(x^2-beta0^2)
}
lines(x,LQAapprox(x),lty=2,lwd=1.4,col="blue")
lines(x,LLAapprox(x),lty=3,lwd=1.4,col="red")

```



```

# lambda=2,beta0=1
par(mar = c(bottom=3, 4, top=2, 3))
matplot(x, yLq,lwd=1.3, type='l', ylim=c(0,10),main=TeX("L_{0.5}:
\\beta^{(0)}=1"), xaxt="n", yaxt="n", cex.main=1.5, lty=1, las=1,
xlab=expression(beta), ylab=TeX("$P_{\\lambda}(\\beta)$"))
axis(1,cex.axis=0.9)
axis(2,cex.axis=0.9)
beta0<-1
LLAapprox<-function(Beta){
Lq(beta0,lambda)+dLq(beta0,lambda)*(abs(Beta)-beta0)
}
LQAapprox<-function(Beta){
Lq(beta0,lambda)+0.5*(dLq(beta0,lambda)/beta0)*(x^2-beta0^2)
}
lines(x,LQAapprox(x),lty=2,lwd=1.4,col="blue")
lines(x,LLAapprox(x),lty=3,lwd=1.4,col="red")

```

مثال ٢.٤.٣

```

##### Load R Packages
library(splines)
library(MASS)
library(sn)
library(energy)

##### Three Screening Procedures: RDC-SIS, DC-SIS, SIS
Screening<-function(X=X,Y=Y)
{
p<-dim(X)[2]
n<- length(Y)
s1<- matrix(0,nrow=p,ncol=1)
s2<- matrix(0,nrow=p,ncol=1)
s3<- matrix(0,nrow=p,ncol=1)
Y<- Y-mean(Y)
OLS<- t(X) %*% Y # for ranking based on SIS

```

```

G.y<- c() # Empirical distribution function of Y
for(i in 1:n){
G.y[i]<- sum(Y<=Y[i])
}
for(j in 1:p){
s1[j,] <- (dcor(X[,j],G.y)) # RDC-SIS
s2[j,] <- (dcor(X[,j],Y)) # DC-SIS
s3[j,] <- abs(OLS[j]) # SIS
}
rank1=(p+1)-rank(s1) # RDC-SIS
rank2=(p+1)-rank(s2) # DC-SIS
rank3=(p+1)-rank(s3) # SIS
list(rank.ROISIS=rank1,rank.DCSIS=rank2,rank.SIS=rank3)
}
### Sure Independent Ranking and Screening (SIRS)
SIRS<-function(X=X,Y=Y){
p<- dim(X)[2]
n<- length(Y)
w<- matrix(0, nrow=p, ncol=1)
X<- apply(X,2, function(x)(x-mean(x))/sd(x))
for(k in 1:p){
w.k.j<-NULL
for(j in 1:n){
s<- (t(X[,k])%*(1*(Y<Y[j])))/n)^2
w.k.j<- c(w.k.j,s)
}
w[k,]<- (n^3*(mean(w.k.j)))/(n*(n-1)*(n-2))
}
list(rank=(p+1)-rank(w))
}
#####
n<- 200 # number of data points
p<- 1000 # number of covariates
ss<- 500 # number of repetitions

```

---

```
gg<- 5      #number of important variables
c<- 2      # c=1 or 2
nm<- 4     # number of methodes (RDC-SIS,DC-SIS,SIS,SIRS)
R<- array(0,c(ss,gg,nm))
Sigma1<- diag(p) ## covariance matrix
for (i in 1:p){
  for (j in 1:p){
    if (i<j) {
      Sigma1[i,j]<- (0.5)^abs(i-j)
      Sigma1[j,i]<- Sigma1[i,j] }
    }
  }

### Beginning of the iteration
for( s in 1:ss)
{
  set.seed(12345+s)
  X<- mvrnorm(n,mu=rep(0,p),Sigma=Sigma1)
  beta<- c(1,0.8,0.6,0.4,0.2,rep(0,p-gg))
  # error<- rsn(n,0,1,2) # different error distributions
  # error<- rt(n,1)
  error<- rnorm(n)
  Y<- c*X%*%beta+sqrt(6.83)*error
  #####
  Q<- Screening(X,Y)
  ex1b<- SIRS(X,Y)
  R[s,,1]<- c(Q$rank.ROSIS[1:gg])
  R[s,,2]<- c(Q$rank.DCSIS[1:gg])
  R[s,,3]<- c(Q$rank.SIS[1:gg])
  R[s,,4]<- c(ex1b$rank[1:gg])
  print(s)
}

quantile(apply(R[, ,1],1,max),c(0.05,0.25,0.50,0.75,0.95)) # ROSIS
quantile(apply(R[, ,2],1,max),c(0.05,0.25,0.50,0.75,0.95)) # DC-SIS
```

```

quantile(apply(R[, ,3], 1, max), c(0.05, 0.25, 0.50, 0.75, 0.95)) # SIS
quantile(apply(R[, ,4], 1, max), c(0.05, 0.25, 0.50, 0.75, 0.95)) # SIRS
#####
h1<- rep(0, nm)
for (i in 1:nm)
{
for(s in 1:ss)
{
if (max(R[s, ,i])<=2*floor(n/log(n))){
h1[i]<- h1[i]+1 }
}
}
proportion1<- h1/ss
proportion1 # propotion S
#####
h2<- rep(0, gg)
for(s in 1:ss)
{
for(k in 1:gg){
if ((R[s, k, 3])<=2*floor(n/log(n))){
h2[k]<- h2[k]+1}
}
}
proportion2<- h2/ss
proportion2 # propotion P_{3}

```

#### مثال ۱.۵.۴

```

### Load R Packages
library(glmnet)
library(abind)
library(magic)
library(mvtnorm)
library(splines)

```

```
library(nor1mix)
library(MASS)
library(SparseM)
library(quantreg)
library(hqreg)

# First order derivative of SCAD penalty
scad.der<-function (x, lambda)
{
x<- abs(x)
a<- 3.7
u<- (x<=lambda)
penalty.derivative<- lambda*u+(pmax(a*lambda-x,0)/(a-1))*(1-u)
return(penalty.derivative)
}

##### Inputs
n<- 200 # sample size
p<- 15 # number of covariates
p1<- 2 # number of constant effects
p2<- 3 # number of varying effects
s<- p1+p2 # number of nonzero effects
degree<- 3 # degree of spline
K<- floor(n^{1/(2*degree+3)}) # number of internal knots

Sigma<- diag(p) ## covariance matrix
for (i in 1:p){
for (j in 1:p){
Sigma[i,j]<- 0.5^abs(i-j)
}
}

#### simulation
nsimu=500; NV=rep(0,nsimu); NC=rep(0,nsimu)
NVT=rep(0,nsimu);NCT=rep(0,nsimu);CF=rep(0,nsimu)
```

```
for (simu in 1:nsimu)
{
set.seed(1234+simu)
X<- mvrnorm(n,mu=rep(0,p),Sigma=Sigma)
u<- runif(n)
error<- rnorm(n) # normal error
# error<- rt(n,3) # t(3) error
# error<- rnormMix(n,norMix(mu=c(0,0),sigma=c(1,5),w=c(0.9,0.1))) # MN error
Y<- 2*sin(2*pi*u)*X[,1]+8*u*(1-u)*X[,2]+2.5*X[,3]+X[,4]+1.5*X[,5]+error

##### calculate R matrix
q=K+degree+1 # number of B-spline basis functions
grid<- seq(0.01,0.99,length=500) # grid points
evalspline<- function(K){
B<- bs(grid,df=q,intercept=TRUE,degree=degree)
R<- t(B)%*%B/length(grid)
list(R=R)
}
temp<- evalspline(K)
R<- temp$R

##### Z matrix
Zmatrix<- function(xx,y,B,K)
{
n<- nrow(xx);p<- ncol(xx)
z<- c()
for (j in 1:p)
{
Phi<- B*xx[,j]
z<- cbind(z,Phi)
}
Z<- cbind(xx,as.matrix(z))
return(Z)
```

---

```
}
Z<- Zmatrix(X,Y,B,K)

##### unpenalized quantile estimation as initial estimate
obj0=function(xx,y,B,K)
{
n<- nrow(xx);p<- ncol(xx)
Z<- Zmatrix(xx,y,B,K)
theta<- rq(y~Z-1)$coef
list(theta=theta)
}
object<- obj0(X,Y,B,K)
thetainit<- object$theta # initial estimate
etainit<- thetainit[1:p]
gammainit<- thetainit[(p+1):length(thetainit)]
pgaminit<- rep(0,p)
for (j in 1:p){
gammaj<- gammainit[((j-1)*(K+4)+1):(j*(K+4))]
pgaminit[j]<- sqrt(gammaj%%R%%gammaj)
}
##### Find optimal tuning parameters using 5-fold CV
cvfolds<- function(n,folds)
{ split(sample(1:n),rep(1:folds,length=n)) }
x<- X
y<- Y
nK<- 5 # number of folds
folds_ans<- cvfolds(n,nK)
mseld<-c ()

## Possible values of parameters
lambda1s<- seq(lambda1.min,lambda1.max,length=20)
lambda2s<- seq(lambda2.min,lambda2.max,length=20)
for (lambda1 in lambda1s)
{
```

```
for (lambda2 in lambda2s)
{
mse<- rep(0,nK)
for(i in 1:nK)
{
omit<- folds_ans[[i]]
xcv<- x[-omit,];ycv=y[-omit]
xt<- x[omit,];yt=y[omit]
Bcv<- B[-omit,];ucv=u[-omit]
Bt<- B[omit,];ut=u[omit]
zcv<- Z[-omit,];zt=Z[omit,]
thetam<- thetainit
iter<- 1; diff<- 1

while(diff>1e-6 & iter<50){
## Optimal h
e<- Y-Z%*%thetam; sigm<- sqrt(var(e))
rh<- rep(0,100)
for (j in 1:100)
{
h<- as.numeric(0.5*sigm*1.02^(j-1))
fh<- mean((e^2/h^2-1)/h^3*dnorm(e/h))
gh<- mean(e^2*(dnorm(e/h))^2/h^6)
rh[j]<- gh/(fh^2*sigm^2)
}
index<- which.min(rh);hopt=0.5*sigm*1.02^(index-1)
ecv<- ycv-zcv%*%thetam
w<- dnorm(ecv,mean=0,sd=hopt);W=as.vector(w/sum(w))
etam<- thetam[1:p]
gammam<- thetam[(p+1):length(thetam)]
pgam<- rep(0,p)
for (j in 1:p){
gammamj<- gammam[((j-1)*(K+4)+1):(j*(K+4))]
pgam[j]<- sqrt(gammamj%*%R%*%gammamj)
```



---

```

}
mpw1<- rep(0,p)
mpw2<- rep(0,p)
for(j in 1:p){
mpw1[j]<- scad.der(etam[j],lambda1)/(abs(etam[j])+1e-6)
mpw2[j]<- scad.der(pgam[j],lambda2)/(pgam[j]+1e-6)
}
omega1<- diag(mpw1)
omega2<- matrix(0,nrow=(K+4)*p,ncol=(K+4)*p)
for (j in 1:p){
omega2[((j-1)*(K+4)+1):(j*(K+4)),((j-1)*(K+4)+1):(j*(K+4))]<-mpw2[j]*R
}
Gm<- adiag(omega1,omega2)
thetam0=thetam
thetam<- ginv(t(zcv)%*%diag(W)%*%zcv+n*Gm)%*%t(zcv)%*%diag(W)%*%ycv
diff<- max(abs(thetam-thetam0))
iter<- iter+1
}
mse[i]<- sum(dnorm((yt-zt)%*%thetam),mean=0,sd=hopt))
}# for nK
mseld<- c(mseld,mean(mse))
} # for lambda2
} # for lambda1

mldindex<- which.max(mseld) # index of optimal tuning parameters
lambda1opt<- lambda1s[floor((mldindex-1)/length(lambda2s)+1)]
lambda2opt<- lambda2s[(mldindex-1)%%length(lambda2s)+1]

##### estimation by optimal tuning parameters #####
thetam<- thetainit
lambda1<- lambda1opt
lambda2<- lambda2opt

while(diff>1e-6 & iter<50){

```

```

## Find optimal h
e<- Y-Z%*%thetam; sigm=sqrt(var(e))
rh<- rep(0,100)
for (j in 1:100)
{
h<- as.numeric(0.5*sigm*1.02^(j-1))
fh<- mean((e^2/h^2-1)/h^3*dnorm(e/h))
gh<- mean(e^2*(dnorm(e/h))^2/h^6)
rh[j]<- gh/(fh^2*sigm^2)
}
index<- which.min(rh);hopt=0.5*sigm*1.02^(index-1)
w<- dnorm(e,mean=0,sd=hopt);W=as.vector(w/sum(w))
etam<- thetam[1:p]
gammam<- thetam[(p+1):length(thetam)]
pgam<- rep(0,p)
for (j in 1:p){
gammamj<- gammam[((j-1)*(K+4)+1):(j*(K+4))]
pgam[j]<- sqrt(gammamj%*%R%*%gammamj)
}
mpw1<- rep(0,p)
mpw2<- rep(0,p)
for(j in 1:p){
mpw1[j]<- scad.der(etam[j],lambda1)/(abs(etam[j])+1e-6)
mpw2[j]<-scad.der(pgam[j],lambda2)/(pgam[j]+1e-6)
}
omega1<- diag(mpw1)
omega2<- matrix(0,nrow=(K+4)*p,ncol=(K+4)*p)
for (j in 1:p){
omega2[((j-1)*(K+4)+1):(j*(K+4)),((j-1)*(K+4)+1):(j*(K+4))]<-mpw2[j]*R
}
Gm<- addiag(omega1,omega2)
thtam0<- thetam
thetam<- ginv(t(Z)%*%diag(W)%*%Z+n*Gm)%*%t(Z)%*%diag(W)%*%Y
diff<- max(abs(thetam-thetam0))

```

---

```

iter<- iter+1
}
etam<- thetam[1:p]
etam[abs(etam)<1e-4]<- 0 # small eta coefficients set to be zero

gammam<- thetam[(p+1):length(thetam)]
for (j in 1:p){ # small gamma coefficients set to be zero
  curgammam<-gammam[((j-1)*(K+4)+1):(j*(K+4))]
  if(sqrt(sum(curgammam^2))<1e-4)
  {gammam[((j-1)*(K+4)+1):(j*(K+4))]<-0}
}
pgam<- rep(0,p)
for (j in 1:p){
  gammamj<- gammam[((j-1)*(K+4)+1):(j*(K+4))]
  pgam[j]<- sqrt(gammamj%*%R%*%gammamj)
}

##### check the varying and constant effects
allk<-rep(K,p)
for (j in 1:p){
  curgammam<-gammam[((j-1)*(K+4)+1):(j*(K+4))]
  if(sum(abs(curgammam))<1e-4&abs(etam[j])<1e-4){allk[j]<- -K;next;}
  # zero effect
  if(sum(abs(curgammam))<1e-4&abs(etam[j])>=1e-4){allk[j]<- -2} else
  {allk[j]<-K} # constant effect
}

#####
# NV: number of varying effects
# NC: number of constant effects
NV<- sum(allk>0); NC<- sum(allk== -2)

# NVT: number of varying coefficients that are correctly identified
# NCT: number of constant coefficients that are correctly identified
u=0;v=0

```

```

for (j in 1:p1)
{ if (allk[j]>0) v=v+1 }
for (k in (p1+1):s)
{ if (allk[k]==-2) u=u+1 }
NVT<- v; NCT<- u

# CF: the proportion of correctly fitted models
r=0
{if(sum(allk[1:p1])==p1*K & sum(allk[(p1+1):s])==-(p2)*2 &
    sum(allk[(s+1):p])==-(p-s)*K) r=1}
CF[simu]<- r

### RASE
u.new<- runif(200) # new grid points
Beta_t<- matrix(0,200,p) # true coefficients
Beta_t[,1]<- (2*sin(2*pi*u.new))
Beta_t[,2]<- (8*u.new*(1-u.new))
Beta_t[,3]<- rep(2.5,200)
Beta_t[,4]<- rep(1.2,200)
Beta_t[,5]<- rep(1.5,200)
B<- bs(u.new,df=q,intercept=TRUE,degree=degree)
Beta_hat<- matrix(0,nrow=200, ncol=p)
for(j in 1:p){
gammaj<- gammam[((j-1)*(K+4)+1):(j*(K+4))]
Beta_hat[,j]<- etam[j]+B%*%gammaj
}
res<- (Beta_t-Beta_hat)^2
RASE[simu]<- mean(apply(res,1,sum))
cat("result:","number varying:", sum(allk>0)," number constant",sum(allk== -2),
    "true number varying:",NNT[simu],"true number constant:",NLT[simu],"CF:",
    CF[simu],"\n")
}

```

## مراجع

- [۱] اریه، م. (۱۳۹۴). برآورد مدل نیمه پارامتری به روش  $B$ -اسپلاین، پایان نامه کارشناسی ارشد، دانشگاه صنعتی شاهرود.
- [۲] روزبه، م. (۱۳۹۰). برآورد در مدل های خطی جزئی، رساله دکتری، دانشگاه فردوسی مشهد.
- [۳] کاظمی، م.، شاهسونی، د. و آرشی، م. (۱۳۹۷). انتخاب متغیر و تشخیص ساختار در بعد بالا برای مدل های جمعی خطی-جزئی، مجله علوم آماری، ۱۲(۲)، ۴۸۵-۵۱۲.
- [۴] نوروزی راد، م. (۱۳۹۶). برآوردگرهای بهبودیافته در برخی مدل های رگرسیونی خطی جریمه شده، رساله دکتری، دانشگاه صنعتی شاهرود.
- [5] Ahmad, I., Leelahanon, S. and Li, Q. (2005). Efficient estimation of a semiparametric partially linear varying coefficient model, *The Annals of Statistics*, 39, 305-332.
- [6] Akdeniz, E., Akdeniz, F. and Roozbeh, M. (2018). A new difference-based weighted mixed Liu estimator in partially linear models, *Statistics*, 52(6), 1309-1327.
- [7] Arashi, M. and Roozbeh, M. (2016). Some improved estimation strategies in high-dimensional semiparametric regression models with application to riboflavin production data, *Statistical Papers*, DOI: 10.1007/s00362-016-0843-y.
- [8] Barut, E., Fan, J. and Verhasselt, A. (2016). Conditional sure independence screening, *Journal of the American Statistical Association*, 111(515), 1266-1277.
- [9] Breheny, P. and Huang, J. (2019). *High-Dimensional Regression Modeling: Methodology, Applications and Software*, Chapman and Hall/CRC.

- [10] Breheny, P. and Huang, J. (2015). Group descent algorithms for nonconvex penalized linear and logistic regression models with grouped predictors, *Statistics and Computing*, 25(2), 173-187.
- [11] Breheny, P. and Huang, J.(2011). Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection, *Annals of Applied Statistics*, 5, 232-253.
- [12] Breiman, L. (1996). Heuristics of instability and stabilization in model selection, *The Annals of Statistics*, 24(6), 2350–2383.
- [13] Bushel, P., Wolfinger, R. D. and Gibson, G. (2007). Simultaneous Clustering of Gene Expression Data with Clinical Chemistry and Pathological Evaluations Reveals Phenotypic Prototypes, *BMC Systems Biology*, Doi:10.1186/1752-0509-1-15.
- [14] Bühlmann, P. and Van De Geer, S. (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*, Springer Science & Business Media.
- [15] Bühlmann, P., Kalisch, M. and Maathuis, M. H. (2010). Variable selection in high-dimensional linear models: partially faithful distributions and the PC-simple algorithm, *Biometrika*, 97(2), 261-278.
- [16] Chatterjee, S. and Hadi, A. S. (2012). *Regression Analysis by Example, 5th Edition*, John Wiley.
- [17] Chen, X. (2018). Model-free conditional feature screening for ultra-high dimensional right censored data, *Journal of Statistical Computation and Simulation*, 88(12), 2425-2446.
- [18] Chen, R., Liang, H. and Wang, J. (2011). Determination of linear components in additive models, *Journal of Nonparametric Statistics*, 23(2), 367-383.
- [19] Chen, J., Li, D., Linton, O. and Lu, Z. (2018). Semiparametric ultra-high dimensional model averaging of nonlinear dynamic time series, *Journal of the American Statistical Association*, 113(522), 919-932.

- [20] Cleveland, W. S., Grosse, E. and Shyu, W. M. (1991). Local regression models, In *Statistical Models in S* (Chambers, J. M. and Hastie, T. J., eds), 309-376. Wadsworth & Brooks, Pacific Grove.
- [21] De Boor, C. (2001). *A Practical Guide to Splines*, Revised Edition, Springer-Verlag, New York.
- [22] Deng, G. and Liang, H. (2010). Model averaging for semiparametric additive partial linear models, *Science China Mathematics*, 53(5), 1363-1376.
- [23] Desboulets, L. (2018). A Review on Variable Selection in Regression Analysis, *Econometrics*, 6(4), 45.
- [24] Dierckx, P. (1995). *Curve and Surface Fitting with Splines*, Oxford University Press.
- [25] Ding, J., Tarokh, V. and Yang, Y. (2018). Model selection techniques: an overview, *IEEE Signal Processing Magazine*, 35(6), 16-34.
- [26] Drapper, N. R. and Smith, H. (1981). *Applied Regression Analysis, 2nd Edition*, John Wiley.
- [27] Du, J., Li, G. and Peng, H. (2015). Variable selection for semiparametric partially linear covariate-adjusted regression models, *Communication in Statistics- Theory and Methods*, 44, 2809-2826.
- [28] Dvoretzky, A., Kiefer, J. and Wolfowitz, J. (1956). Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator, *Annals of Mathematical Statistics*, 27, 642-669.
- [29] Dziak, J., Li, R. and Collins, L. (2005). *Critical Review and Comparison of Variable Selection Procedures for Linear Regression*, State College, PA: Pennsylvania State University.
- [30] Efron B., Hastie T., Johnstone I. and Tibshirani, R. (2004). Least angle regression, *The Annals of Statistics*, 32, 407-451.
- [31] Engle R. F., Granger C. W. J., Rice J. and Weiss A. (1986 ). Semiparametric estimates of the relation between weather and electricity sales, *Journal of the American Statistical Association*, 81, 310-320.

- 
- [32] Fan, J., Feng, Y. and Wu, Y. (2010). High-dimensional variable selection for Cox's proportional hazards model, *In Borrowing Strength: Theory Powering Applications—A Festschrift for Lawrence D. Brown* (pp. 70-86). Institute of Mathematical Statistics.
- [33] Fan J., Feng Y. and Song R. (2011). Nonparametric independence screening in sparse ultrahigh-dimensional additive models, *Journal of the American Statistical Association*, 106(494), 544-557.
- [34] Fan, J., Han, F. and Liu, H. (2014). Challenges of big data analysis, *National Science Review*, 1(2), 293-314.
- [35] Fan, J. and Huang, T. (2005). Profile likelihood inferences on semiparametric varying coefficient partially linear models, *Bernoulli*, 11, 1031-1057.
- [36] Fan J. and Li R., (2001). Variable selection via nonconcave penalized likelihood and its oracle properties, *Journal of the American Statistical Association*, 96, 1348-1360.
- [37] Fan, J. and Li, R. (2006). Statistical challenges with high dimensionality: Feature selection in knowledge discovery, In *Proceedings of the International Congress of Mathematicians*, 3, 595-622.
- [38] Fan, J. and Lv, J. (2008). Sure independence screening for ultrahigh dimensional feature space, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(5), 849-911.
- [39] Fan, J., Ma, Y. and Dai, W. (2014). Nonparametric independence screening in sparse ultra-high-dimensional varying coefficient models, *Journal of the American Statistical Association*, 109(507), 1270-1284.
- [40] Fan, J., Samworth, R. and Wu, Y. (2009). Ultrahigh dimensional feature selection: beyond the linear model, *Journal of Machine Learning Research*, 10(Sep), 2013-2038.
- [41] Fan J. and Song R. (2010). Sure independence screening in generalized linear models with NP-dimensionality, *The Annals of Statistics*, 6, 3567–3604.



- [42] Frank, L. E. and Friedman, J. H. (1993). A statistical view of some chemometrics regression tools, *Technometrics*, 35(2), 109-135.
- [43] Fu, W. J. (1998). Penalized regressions: the bridge versus the lasso, *Journal of Computational and Graphical Statistics*, 7(3), 397-416.
- [44] Geenens, G. (2011). Curse of dimensionality and related issues in nonparametric functional regression, *Statistics Surveys*, 5, 30-43.
- [45] Giraud, C. (2014). *Introduction to high-dimensional statistics*, Chapman and Hall/CRC.
- [46] Green, P. J. and Silverman, B. W. (1993). *Nonparametric Regression and Generalized Linear Models: a Roughness Penalty Approach*, CRC Press.
- [47] Gujarati, D. N. (2004). *Basic Econometrics (4th Edition)*, Tata McGraw- Hill.
- [48] Gunst, R. F. and Webster, J. T. (1975). Regression analysis and problems of multicollinearity, *Communications in Statistics-Theory and Methods*, 4(3), 277-292.
- [49] Guo, J., Tang, M., Tian, M. and Zhu, K. (2013). Variable selection in high-dimensional partially linear additive models for composite quantile regression, *Computational Statistics and Data Analysis*, 65, 56-67.
- [50] Hall, P. and Miller, H. (2009). Using generalized correlation to effect variable selection in very high dimensional problems, *Journal of Computational and Graphical Statistics*, 18(3), 533-550.
- [51] Härdle, W. (1990). *Applied Non-parametric Regression*, Economic Society Monographs, vol, 19. Cambridge University Press.
- [52] Härdle, W., Hall, P. and Ichimura, H. (1993). Optimal smoothing in single-index models, *The Annals of Statistics*, 21(1), 157-178.
- [53] Härdle, W., Liang, H. and Gao, J. T. (2000). *Partially Linear Models*, Heidelberg: Springer Physica- Verlag.
- [54] Härdle, W. K., Muller, M., Sperlich, S. and Werwatz, A. (2004). *Nonparametric and Semiparametric Models*, Springer.

- 
- [55] Härdle, W. and Stoker, T. M. (1989). Investigating smooth multiple regression by the method of average derivatives, *Journal of the American Statistical Association*, 84(408), 986-995.
- [56] Hastie, T. and Tibshirani, R. (1990). *Generalized Additive Models*, Chapman and Hall.
- [57] Hastie, T. and Tibshirani, R. (1993). *Varying-coefficient models*, *Journal of the Royal Statistical Society, Series B (Methodological)*, 55(4), 757-796.
- [58] Hastie, T., R. Tibshirani and J. Friedman (2009). *The Elements of Statistical Learning: Data Mining, Inference and Prediction (2nd Edition)*, Springer.
- [59] Heinze, G., Wallisch, C. and Dunkler, D. (2018). Variable selection—a review and recommendations for the practicing statistician, *Biometrical Journal*, 60(3), 431-449.
- [60] Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: biased estimation for nonorthogonal problems, *Technometrics*, 12, 55–67.
- [61] Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: applications nonorthogonal problems, *Technometrics*, 12, 69–82.
- [62] Hu, T., Xia, Y., (2012). Adaptive semi-varying coefficient model selection, *Statistica Sinica*, 22, 575-599.
- [63] Huang, J., Wei, F. and Ma, S. (2012a). Semiparametric regression pursuit, *Statistica Sinica*, 22, 1403-1426.
- [64] Huang, J., Breheny, P. and Ma, S. (2012b). A selective review of group selection in high-dimensional models, *Statistical Science*, 27(4), 481–499.
- [65] Huang, J., Horowitz, J. L. and Wei, F., (2010). Variable selection in nonparametric additive models, *The Annals of Statistics*, 38, 2282-2313.
- [66] Huang, J., Wu, C., Zhou, L. (2002). Varying-coefficient models and basis function approximation for the analysis of repeated measurements, *Biometrika*, 89, 111–128.

- 
- [67] Huang, Q. and Zhu, Y. (2016). Model-free sure screening via maximum correlation, *Journal of Multivariate Analysis*, 148, 89-106.
- [68] Huber, P. (1981). *Robust Estimation*, New York: John Wiley.
- [69] Hunter, D. R. and Li, R. (2005). Variable selection using MM algorithms, *The Annals of Statistics*, 33(4), 1617-1642.
- [70] James, G., Witten, D., Hastie, T. and Tibshirani, R. (2013). *An Introduction to Statistical Learning (Vol. 6)*, New York: springer.
- [71] Jia, J. and Yu, B. (2010). On model selection consistency of elastic net when  $p \gg n$ , *Statistica Sinica*, 20, 595–611.
- [72] Johnstone, I. M. and Titterington, D. M. (2009). Statistical challenges of high-dimensional data, *Philosophical Transactions of the Royal Society A*, 367, 4237–4253
- [73] Kai, B., Li, R. and Zou, H. (2011). New efficient estimation and variable selection method for semiparametric varying-coefficient partially linear models, *The Annals of Statistics*, 39, 305-332.
- [74] Kazemi, M., Shahsavani, D. and Arashi, M. (2019a). A sure independence screening procedure for ultra-high dimensional partially linear additive models, *Journal of Applied Statistics*, 46(8), 1385-1403.
- [75] Kazemi, M., Shahsavani, D. and Arashi, M. (2018). Variable selection and structure identification for ultrahigh-dimensional partially linear additive models with application to cardiomyopathy microarray data, *Statistics, Optimization and Information Computing*, DOI: 10.19139/soic.v6i3.576.
- [76] Kazemi, M., Shahsavani, D., Arashi, M. and Rodrigues, P. C. (2019b). Estimation in partial linear model with spline modal function, *Communications in Statistics - Simulation and Computation* (Revised).
- [77] Kong, X. B., Liu, Z., Yao, Y. and Zhou, W. (2017). Sure screening by ranking the canonical correlations, *Test*, 26(1), 46-70.

- [78] Leng, C. (2009). A simple approach for varying-coefficient model selection. *Journal of Statistical Planning and Inference*, 139, 2138-2146.
- [79] Leng, C., Lin, Y. and Wahba, G. (2006). A note on the lasso and related procedures in model selection, *Statistica Sinica*, 16, 1273–1284.
- [80] Li, Q. and Lin, N. (2010). The Bayesian elastic net, *Bayesian Analysis*, 5(1), 151-170.
- [81] Li, Y., Li, G., Lian, H. and Tong, T. (2017). Profile forward regression screening for ultra-high dimensional semiparametric varying coefficient partially linear models, *Journal of Multivariate Analysis*, 155, 133-150.
- [82] Li, X., Ma, X. and Zhang, J. (2018). Conditional quantile correlation screening procedure for ultrahigh-dimensional varying coefficient models, *Journal of Statistical Planning and Inference*, 197, 69-92.
- [83] Li, J., Ray, S. and Lindsay, B. (2007). A nonparametric statistical approach to clustering via mode identification, *Journal of Machine Learning Research*, 8, 1687–1723.
- [84] Li, R., Zhong, W. and Zhu, L. (2012). Feature screening via distance correlation learning, *Journal of the American Statistical Association*, 107(499), 1129-1139.
- [85] Lian, H. (2012a). Variable selection in high-dimensional partly linear additive models, *Journal of Nonparametric Statistics*, 24, 825-839.
- [86] Lian, H. (2012b). Shrinkage estimation for identification of linear components in additive models, *Statistics and Probability Letters*, 82, 225-231.
- [87] Lian, H. (2012c). Semiparametric estimation of additive quantile regression models by two-fold penalty, *Journal of Business and Economic Statistics*, 30(3), 337-350.
- [88] Lian, H., Chen, X. and Yang, J. Y. (2012). Identification of partially linear structure in additive models with an application to gene expression prediction from sequences, *Biometrics*, 68(2), 437-445.

- [89] Lian, H, Du, P, Li, Y., and Liang, H. (2014). Partially linear structure identification in generalized additive models with NP-dimensionality, *Computational Statistics and Data Analysis*, 80, 197–208.
- [90] Lian, H., Liang, H. and Ruppert, D. (2015). Separation of covariates into nonparametric and parametric parts in high-dimensional partially linear additive models, *Statistica Sinica*, 591-607.
- [91] Liu, J. (2016). Feature screening and variable selection for partially linear models with ultrahigh-dimensional longitudinal data, *Neurocomputing*, 195, 202-210.
- [92] Liu, J., Li, R. and Wu, R. (2014). Feature selection for varying coefficient models with ultrahigh-dimensional covariates, *Journal of the American Statistical Association*, 109(505), 266-274.
- [93] Liu, X., Wang, L. and Liang, H. (2011). Estimation and variable selection for semi-parametric additive partial linear models, *Statistica Sinica*, 21(3), 1225-1248.
- [94] Lv, J., Yang, H. and Guo, C. (2017). Variable selection in partially linear additive models for modal regression, *Communications in Statistics-Simulation and Computation*, 46(7), 5646-5665.
- [95] Meier, L., Van de Geer, S. and Bühlmann, P. (2009). High-dimensional additive modeling, *The Annals of Statistics*, 37(6B), 3779-3821.
- [96] Montgomery, D. C., Peck, E. A. and Vining, G. G. (2012). *Introduction to Linear Regression Analysis* (Vol. 821), John Wiley.
- [97] Nadaraya, E. A. (1964). On Estimating Regression, *Theory of Probability and Its Applications*, 9(1), 141–142.
- [98] Noh, H., Keilegom, V. (2012). Efficient model selection in semivarying coefficient models, *Electronic Journal of Statistics*, 6, 2519-2534.
- [99] Opsomer, J. D. and Ruppert, D. (1999). A root-n consistent backfitting estimator for semiparametric additive modeling, *Journal of Computational and Graphical Statistics*, 8(4), 715-732.

- 
- [100] Park, T. and Casella, G. (2008). The bayesian lasso, *Journal of the American Statistical Association*, 103(482), 681-686.
- [101] Parzen, E. (1962). On estimation of a probability density function and mode, *The Annals of Mathematical Statistics*, 33, 1065–1076
- [102] Qiao, X. (2014). Variable selection using  $L_q$  penalties, *Wiley Interdisciplinary Reviews: Computational Statistics*, 6(3), 177-184.
- [103] Roozbeh, M. (2015). Shrinkage ridge estimators in semiparametric regression models, *Journal of Multivariate Analysis*, 136, 56-74.
- [104] Roozbeh, M. (2016). Robust ridge estimator in restricted semiparametric regression models, *Journal of Multivariate Analysis*, 147, 127-144.
- [105] Ruppert, D., Wand, M. P. and Carroll, R. J. (2003). *Semiparametric Regression*, Cambridge University Press, Cambridge.
- [106] Schmalensee, R. and Stoker, T.M. (1999). Household gasoline demand in the United States, *Econometrica*, 67, 645-662.
- [107] Schumaker, L. (1981). *Splines Function: Basic Theory*, New York: John Wiley.
- [108] Shen, S., Cui, J., Mei, C., Wang, C. (2014). Estimation and inference of semi-varying coefficient models with heteroscedastic errors, *Journal of Multivariate Analysis*, 124, 70-93.
- [109] Song, R., Yi, F. and Zou, H. (2014). On varying-coefficient independence screening for high-dimensional varying-coefficient models, *Statistica Sinica*, 24(4), 1735.
- [110] Song, A., Ma, T., Lv, S. and Lin, C. (2018). A model-free variable selection method for reducing the number of redundant variables, *Statistics*, 52(6), 1212-1248.
- [111] Stone, C. (1982). Optimal global rates of convergence for nonparametric regression, *The Annals of Statistics*, 10, 1040-1053.
- [112] Sun, J., Lin, L. (2014). Local rank estimation and related test for varying-coefficient partially linear models, *Journal of Nonparametric Statistics*, 26, 187-206.

- [113] Szekely, G. J. and Rizzo, M. L. (2009). Brownian Distance Covariance, *The Annals of Applied Statistics*, 3, 1233–1303.
- [114] Székely, G. J., Rizzo, M. L. and Bakirov, N. K. (2007). Measuring and testing dependence by correlation of distances, *The Annals of Statistics*, 35(6), 2769-2794.
- [115] Tang, Y., Wang, H., Zhu, Z., Song, X. (2012). A unified variable selection approach for varying coefficient models, *Statistica Sinica*, 22, 601-628.
- [116] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society, Series B (Methodological)*, 267-288.
- [117] Wang, X. and Leng, C. (2016). High dimensional ordinary least squares projection for screening variables, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(3), 589-611.
- [118] Wang, L., Li, H., Huang, J. Z. (2008). Variable selection in nonparametric varying coefficient models for analysis of repeated measurements, *Journal of the American Statistical Association*, 103, 1556-1569.
- [119] Wang, K., Lin, L. (2016). Robust structure identification and variable selection in partial linear varying coefficient models, *Journal of Statistical Planning and Inference*, 174, 153-168.
- [120] Wang, M. and Song, L. (2013). Identification for semiparametric varying coefficient partially linear models, *Statistics and Probability Letters*, 83(5), 1311-1320.
- [121] Wang, H., Xia, Y. (2009). Shrinkage estimation of the varying coefficient model, *Journal of the American Statistical Association*, 104, 747-757.
- [122] Wang, Z., Xue, L., Li, G., Lu, F. (2018). Spline estimator for ultra-high dimensional partially linear varying coefficient models, *Annals of the Institute of Statistical Mathematics*, 1-21.
- [123] Watson, G.S. (1964). Smooth regression analysis, *Sankhya A*, 26, 359–372.
- [124] Wei, C. H. and Liu, C. (2012). Statistical inference on semi-parametric partial linear additive models, *Journal of Nonparametric Statistics*, 24(4), 809-823.

- [125] Willems, J. P., Saunders, J. T., Hunt, D. E. and Schorling, J. B. (1997). Prevalence of coronary heart disease risk factors among rural blacks: a community-based study, *Southern Medical Journal*, 90(8), 814-820.
- [126] Wu, T. T. and Lange, K. (2008). Coordinate descent algorithms for lasso penalized regression, *The Annals of Applied Statistics*, 2(1), 224-244.
- [127] Xia, X., Yang, H. and Li, J. (2016). Feature screening for generalized varying coefficient models with application to dichotomous responses, *Computational Statistics and Data Analysis*, 102, 85-97.
- [128] Xia, Y., Zhang, W., Tong, H. (2004). Efficient estimation for semivarying-coefficient models, *Biometrika*, 91, 661-681.
- [129] Xue, L. (2009). Consistent variable selection in additive models, *Statistica Sinica*, 1281-1296.
- [130] Xue, L., Qu, A. (2012). Variable selection in high-dimensional varying coefficient models with global optimality, *Journal of Machine Learning Research*, 13, 1973-1998
- [131] Yao, W., Li, L. (2014). A new regression model: modal linear regression, *Scandinavian Journal of Statistics*, 41(3), 656-671.
- [132] Yao, W., Lindsay, B., Li, R. (2012). Local modal regression, *Journal of Nonparametric Statistics*, 24, 647-663.
- [133] Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1), 49-67.
- [134] Yuan, M. and Lin, Y. (2007). Model selection and estimation in the Gaussian graphical model, *Biometrika*, 94(1), 19-35.
- [135] Yuan, M. and Lin, Y. (2007). On the non-negative garrotte estimator, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2), 143-161.



- [136] Zhang, C. H. (2010). Nearly unbiased variable selection under minimax concave penalty, *The Annals of Statistics*, 38(2), 894-942.
- [137] Zhang, H. H., Cheng, G. and Liu, Y. (2011). Linear or nonlinear? Automatic structure discovery for partially linear models, *Journal of the American Statistical Association*, 106(495), 1099-1112.
- [138] Zhang, Y. and Li, R. (2011). Iterative conditional maximization algorithm for non-concave penalized likelihood, In *Nonparametric Statistics and Mixture Models: A Festschrift in Honor of Thomas P Hettmansperger*, 336-351.
- [139] Zhang, J., Liu, Y., Wu, Y. (2017). Correlation rank screening for ultrahigh-dimensional survival data, *Computational Statistics and Data Analysis*, 108, 121-132.
- [140] Zhang, R., Zhao, W., Liu, J. (2013). Robust estimation and variable selection for semiparametric partially linear varying coefficient model based on modal regression, *Journal of Nonparametric Statistics*, 25, 523-544.
- [141] Zhao, P., Xue, L. (2009). Variable selection in semiparametric regression analysis for longitudinal data, *Annals of the Institute of Statistical Mathematics*, 64, 213-231.
- [142] Zhao, P., Xue, L. (2010). Variable selection for semiparametric varying coefficient partially linear errors-in-variables models, *Journal of Multivariate Analysis*, 101, 1872-1883.
- [143] Zhao, P. and Yu, B. (2006). On model selection consistency of Lasso, *Journal of Machine Learning Research*, 7(Nov), 2541-2563.
- [144] Zhao, W., Zhang, R., Liu, J., Lv, Y. (2014). Robust and efficient variable selection for semiparametric partially linear varying coefficient model based on modal regression, *Annals of the Institute of Statistical Mathematics*, 66, 165-191.
- [145] Zhao, W., Zhang, R., Lv, Y., Liu, J. (2013). Variable selection of the quantile varying coefficient regression models, *Journal of the Korean Statistical Society*, 42(3), 343-358.

- 
- [146] Zhou, Y. and Liang, H. (2009). Statistical inference for semiparametric varying-coefficient partially linear models with error-prone linear covariates, *The Annals of Statistics*, 37(1), 427-458.
- [147] Zhu, L.P., Li, L., Li, R., Zhu, L.X., (2011). Model-free feature screening for ultra-high dimensional data, *Journal of the American Statistical Association*, 106, 1464-1475.
- [148] Zou, H. (2006). The adaptive lasso and its oracle properties, *Journal of the American Statistical Association*, 101(476), 1418-1429.
- [149] Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2), 301-320.
- [150] Zou, H. and Li, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models, *The Annals of Statistics*, 36(4), 1509-1533.





# واژه‌نامه فارسی به انگلیسی

Interaction	اثر متقابل
Evaluation	ارزیابی
Bias	اریبی
Multivariate adaptive regression splines	اسپلاین رگرسیونی تطبیقی چند متغیره
Cubic spline	اسپلاین مکعبی
Cross validation	اعتبارسنجی متقابل
Flexibility	انعطاف‌پذیری
Residual	باقیمانده
Ill-Condition	بدشرطیدگی
Shrinkage estimation	برآورد انقباضی
Oracle estimator	برآوردگر پیشگو
Optimal	بهینه
Overfitting	بیش‌برازش
Penalty parameter	پارامتر جریمه
Basis	پایه
Assumption	پذیره
Bandwidth	پهنای باند
Prediction	پیشگویی
Continuous	پیوسته
Density function	تابع چگالی
Differentiable function	تابع مشتق‌پذیر
Cholesky decomposition	تجزیه چولسکی
Structure identification	تشخیص ساختار
Orthogonal projection	تصویرسازی متعامد
Interpretability	تفسیرپذیری

Locally quadratic approximation	تقریب درجه دو موضعی
Unimodal	تک‌مدی
Heavy-tailed distribution	توزیع دم‌سنگین
Asymptotic distribution	توزیع مجانبی
Sparsity	تنکی
Piecewise polynomial	چندجمله‌ای تکه‌ای
Test dataset	داده‌های آزمایشی
Training dataset	داده‌های آموزشی
Outlier	داده پرت
High-dimensional data	داده‌های با بعد بالا
Ultrahigh-dimensional data	داده‌های با بعد بسیار بالا
Human capital earnings	درآمد سرمایه انسانی
Full rank	رتبه کامل
Bridge regression	رگرسیون بریج
Penalized regression	رگرسیون جریمه‌شده
Quantile regression	رگرسیون چندکی
Ridge regression	رگرسیون ریج
Lasso regression	رگرسیون لاسو
Local regression	رگرسیون موضعی
Principal component regression	رگرسیون مؤلفه‌های اصلی
Modal regression	رگرسیون نما
Surface	رویه
Microarray	ریزآرایه
Consistency	سازگاری
Identifiability	شناساپذیری
Coefficient of determination	ضریب تعیین
Variance inflation factor	عامل تورم واریانس
Screening	غربالگری
Efficiency	کارایی
Stepwise	گام به گام
Internal knot	گره داخلی
Diagonal matrix	ماتریس قطری

Gram-Schmidt orthonormalization	متعامدسازی گرام-اشمیت
Response variable	متغیر پاسخ
Explanatory Variable	متغیر توضیحی
Categorical variable	متغیر طبقه‌ای
Dummy variable	متغیر ظاهری
Important variable	متغیر مهم
Symmetric	متقارن
Asymptotic	مجانبی
Convex	محدب
Parametric model	مدل پارامتری
Single index model	مدل تک شاخص
Additive model	مدل جمعی
Partial linear model	مدل خطی-جزئی
Varying coefficient model	مدل ضریب متغیر
Non-parametric model	مدل ناپارامتری
Semiparametric model	مدل نیمه پارامتری
Curse of dimensionality	مشقت بعد
Genetic association studies	مطالعات پیوند ژنتیکی
Threshold	مقدار آستانه
Initial value	مقدار اولیه
True value	مقدار درست
Concave	مقعر
Minimax concave	مقعر مینیماکس
Unbiased	نااریب
Convergence rate	نرخ همگرایی
Robust	نیرومند
Profile	نیمرخ
Sure screening property	ویژگی غربالگری مطمئن
Kernel	هسته
Distance correlation	همبستگی فاصله‌ای
Collinearity	همخطی
Exact collinearity	همخطی کامل

Partial collinearity .....	همخطی ناقص
Smooth.....	هموار
Kernel smoothing .....	هموارساز هسته



# واژه‌نامه انگلیسی به فارسی

Assumption	پذیره
Asymptotic	مجانبی
Asymptotic distribution	توزیع مجانبی
Additive model	مدل جمعی
Bandwidth	پهنای باند
Basis	پایه
Bias	اریبی
Bridge regression	رگرسیون بریج
Categorical variable	متغیر طبقه‌ای
Cholesky decomposition	تجزیه چولسکی
Collinearity	همخطی
Concave	مقعر
Consistency	سازگاری
Continuous	پیوسته
Convex	محدب
Curse of dimensionality	مشقت بعد
Convergence rate	نرخ همگرایی
Cross validation	اعتبارسنجی متقابل
Coefficient of determination	ضریب تعیین
Cubic spline	اسپلاین مکعبی
Density function	تابع چگالی
Diagonal matrix	ماتریس قطری
Differentiable function	تابع مشتق‌پذیر
Distance correlation	همبستگی فاصله‌ای
Dummy variable	متغیر ظاهری

Efficiency	کارایی
Exact collinearity	همخطی کامل
Explanatory Variable	متغیر توضیحی
Evaluation	ارزیابی
Flexibility	انعطاف‌پذیری
Full rank	رتبه کامل
Genetic association studies	مطالعات پیوند ژنتیکی
Gram-Schmidt orthonormalization	متعامدسازی گرام-اشمیت
Heavy-tailed distribution	توزیع دم‌سنگین
High-dimensional data	داده‌های با بعد بالا
Human capital earnings	درآمد سرمایه انسانی
Ill-Condition	بدشرطیدگی
Important variable	متغیر مهم
Identifiability	شناساپذیری
Interaction	اثر متقابل
Internal knot	گره داخلی
Interpretability	تفسیرپذیری
Initial value	مقدار اولیه
Kernel	هسته
Kernel smoothing	هموارسازی هسته
Lasso regression	رگرسیون لاسو
Locally quadratic approximation	تقریب درجه دو موضعی
Local regression	رگرسیون موضعی
Microarray	ریزآرایه
Minimax concave	مقعر مینیماکس
Multivariate adaptive regression splines	اسپلاین رگرسیونی تطبیقی چند متغیره
Modal regression	رگرسیون نما
Non-parametric model	مدل ناپارامتری
Orthogonal projection	تصویرسازی متعامد
Overfitting	بیش‌برازش
Outlier	داده پرت
Oracle estimator	برآوردگر پیشگو

Optimal.....	بهینه.....
Parametric model.....	مدل پارامتری.....
Partial collinearity.....	همخطی ناقص.....
Partial least squares regression.....	رگرسیون کمترین توان‌های دوم جزئی.....
Partial linear model.....	مدل خطی- جزئی.....
Penalty parameter.....	پارامتر جریمه.....
Penalized regression.....	رگرسیون جریمه‌شده.....
Piecewise polynomial.....	چندجمله‌ای تکه‌ای.....
Prediction.....	پیشگویی.....
Principal component regression.....	رگرسیون مؤلفه‌های اصلی.....
Profile.....	نیمرخ.....
Quantile regression.....	رگرسیون چندکی.....
Ridge regression.....	رگرسیون ریج.....
Residual.....	باقیمانده.....
Robust.....	نیرومند.....
Response variable.....	متغیر پاسخ.....
Symmetric.....	متقارن.....
Shrinkage estimation.....	برآورد انقباضی.....
Semiparametric model.....	مدل نیمه‌پارامتری.....
Stepwise.....	گام به گام.....
Structure identification.....	تشخیص ساختار.....
Sure screening property.....	ویژگی غربالگری مطمئن.....
Screening.....	غربالگری.....
Smooth.....	هموار.....
Sparsity.....	تنکی.....
Surface.....	رویه.....
Single index model.....	مدل تک‌شاخص.....
Test dataset.....	داده‌های آزمایشی.....
Threshold.....	مقدار آستانه.....
Training dataset.....	داده‌های آموزشی.....
True value.....	مقدار درست.....
Unbiased.....	نالریب.....

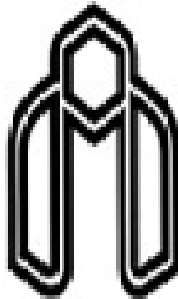
Unimodal .....	تک‌مدی
Ultrahigh-dimensional data .....	داده‌های با بعد بسیار بالا
Varying coefficient model .....	مدل ضریب متغیر
Variance inflation factor .....	عامل تورم واریانس



## **Abstract**

Semiparametric models are often used in real data analysis for their flexibility and parsimony. Two well-known examples of semiparametric models are the partially linear additive model and the partially linear varying coefficient model. Statistical inference of these models is restricted to a condition that the parametric and nonparametric parts are known in advance. However, such prior information is usually unavailable, especially when the number of covariates is large. Therefore, it is of great interest to develop some efficient methods to distinguish parametric components from nonparametric ones. In this dissertation, we introduce a two-step procedure, in the context of ultra-high dimensional additive models, which aims to reduce the size of covariates vector and distinguish linear and nonlinear effects among nonzero components. Also, we propose a robust method for simultaneously variable selection and parametric component identification in varying coefficient models based on modal regression, which is robust with respect to non-normal errors and outliers in the response. The performance of the two proposed methods is examined by simulation studies and real data analysis. Results of numerical studies demonstrate the superiority in comparison with the existing methods.

**Keywords:** Dimensionality reduction, Modal regression, Partially linear additive model, Partially linear varying coefficient model, Structure identification, Sure independence screening, Variable selection.



**Shahrood University of Technology**

**Faculty of Mathematical Sciences**

**Ph.D. Thesis in Statistics**

# **Variable Selection and Structure Identification in Semiparametric Models**

**By: Mohammad Kazemi**

**Supervisors**

**Dr. Davood Shamsavani**

**Dr. Mohammad Arashi**

**February 2019**