

انتخاب متغیر و تشخیص ساختار در بعد بالا برای مدل‌های جمعی خطی-جزیی

محمد کاظمی، داود شاهسونی، محمد آرشی

گروه آمار، دانشگاه صنعتی شاهرود

تاریخ دریافت: ۱۳۹۶/۱۰/۰۷ تاریخ آخرین بازنگری: ۱۳۹۷/۰۴/۲۹

چکیده: در این مقاله یک روش دو مرحله‌ای برای انتخاب متغیر و تشخیص مؤلفه‌های خطی و غیرخطی در مدل‌های جمعی با بعد بالا معرفی می‌شود. در مرحله اول، از یک روش غربالگری برای کاهش بعد فضای متغیرها 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 دم سنگین یا دارای مقادیر فرین باشد، عملکرد خوبی را از خود نشان می‌دهد. در مرحله دوم، از روشی مبتنی بر دو تابع تاوان برای انتخاب همزمان مؤلفه‌های غیرصفر و خطی استفاده می‌شود. کارایی این روش دو مرحله‌ای با مطالعه شبیه‌سازی و تحلیل یک مجموعه داده واقعی بررسی شده است.

واژه‌های کلیدی: انتخاب متغیر، تشخیص ساختار، غربالگری، کاهش بعد، مدل جمعی خطی-جزیی.

۱ مقدمه

مدل جمعی ناپارامتری را به صورت

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

در نظر بگیرید، که در آن متغیر پاسخ، $Y_i = (X_{i1}, \dots, X_{ip})^T$ بردار متغیرهای توضیحی و f_j ها توابع یک متغیره هموار نامعلوم هستند. منظور از تابع هموار، تابعی است که بی‌نهایت بار مشتق پذیر باشد.

آدرس الکترونیکی نویسنده مسئول مقاله: محمد کاظمی، m.kazemie64@yahoo.com

کد موضوع بندی ریاضی (۲۰۱۰): 62G08, 62G05, 62J07

همچنین $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$ بردار خطای تصادفی با میانگین صفر و $E(\varepsilon\varepsilon^T) = \sigma^2 I_n$ 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، اپسومر و راپرت (۱۹۹۹) مدل جمعی خطی-جزیی را به صورت

$$Y_i = \sum_{j \in S_1} \beta_j X_{ij} + \sum_{j \in S_2} f_j(X_{ij}) + \varepsilon_i, \quad (2)$$

معرفی کردند، که در آن S_1 و S_2 دو زیر مجموعه مکمل و دو به دو ناسازگار از $\{1, \dots, p\}$ است. این مدل از دو ویژگی مثبت تفسیرپذیری مدل‌های خطی و انعطاف‌پذیری مدل‌های ناپارامتری استفاده می‌کند. از نظر آماری، مؤلفه‌های پارامتری دارای همگرایی سریع‌تر از مؤلفه‌های ناپارامتری هستند. در نتیجه در نظر گرفتن مؤلفه‌های خطی به صورت غیرخطی کارایی برآورد را کاهش می‌دهد. بنابراین، استفاده از مدل‌های جمعی خطی-جزیی در برخی از مسائل، نسبت به مدل‌های جمعی ناپارامتری، مناسب‌تر است. مسئله برآورد و انتخاب متغیر در مدل جمعی خطی-جزیی تاکنون توسط محققین زیادی مورد مطالعه قرار گرفته است. از جمله می‌توان به لیو و همکاران (۲۰۱۱)، لیان (۲۰۱۲a)، گو و همکاران (۲۰۱۳)، دو و همکاران (۲۰۱۵) و لیو و همکاران (۲۰۱۶) اشاره کرد.

مهمترین فرض در استفاده از مدل (۲) این است که اجزای خطی و غیرخطی مدل از پیش تعیین شده‌اند. اما در عمل چنین اطلاعات پیشین به ندرت در دسترس است، به ویژه وقتی که تعداد متغیرهای توضیحی زیاد باشد، یا به عبارتی مدلی با بعد بالا باشد. بنابراین، علاوه بر تشخیص مؤلفه‌های غیرصفر، شناسایی مؤلفه‌های خطی و غیرخطی نیز بسیار حائز اهمیت است.

ژانگ و همکاران (۲۰۱۱) از دو تابع تاوان برای شناسایی همزمان مؤلفه‌های خطی و غیرخطی در مدل جمعی خطی-جزیی استفاده کردند، اما نتوانستند ویژگی سازگاری در انتخاب را برای روش ارائه شده ثابت کنند. لذا، هوآنگ و همکاران (۲۰۱۲) یک روش موسوم به «تعقیب رگرسیون نیمه پارامتری»^۱ را برای تشخیص مؤلفه‌های خطی و غیرخطی با استفاده از تابع تاوان محدب مینیماکس^۲ (MCP) ارائه دادند. لیان (۲۰۱۲b) روشی برای تشخیص مؤلفه‌های خطی با استفاده از تابع تاوان SCAD^۳ معرفی کرد. سپس لیان و همکاران (۲۰۱۲c) با به کار بردن دو تابع تاوان SCAD توانستند مؤلفه‌های خطی و

¹ Semiparametric regression pursuit

² Minimax Concave Penalty

³ Smoothly Clipped Absolute Deviation

غیرصفر را به طور همزمان تشخیص دهند. همچنین لیان و همکاران (۲۰۱۵) با استفاده از دو تابع تاوان $LASSO^4$ گروهی تطبیقی به شناسایی مؤلفه‌های خطی و غیرصفر در یک مدل جمعی خطی-جزیی با بعد بالا^۵ پرداختند. در داده‌های با بعد بالا، تعداد متغیرهای توضیحی، p ، می‌تواند بزرگ‌تر از حجم نمونه، n ، باشد ($p > n$). در این نوع داده‌ها، معمولاً فقط تعداد اندکی از متغیرهای توضیحی واقعاً با متغیر پاسخ مرتبط هستند. به دلیل وجود تعداد زیادی از متغیرها در مدل‌های با بعد بالا، تفسیر این مدل‌ها بسیار مشکل است. لذا مسئله انتخاب متغیر نقش بسیار مهمی را در مدل‌سازی آماری با بعد بالا ایفا می‌کند.

موضوع بعد بالای داده‌ها، چالشی اجتناب ناپذیر در مسئله مدل‌سازی است که مرهون ارتقای فناوری تولید داده‌های بزرگ در علوم مختلف، از جمله ژنتیک و پزشکی و همچنین افزایش توان ذخیره‌سازی کامپیوترها است. فراتر از آن، داده‌های با «بعد بسیار بالا»^۶ فرصتی پدید آورده‌اند که روش‌های نوین انتخاب متغیر آماری و یادگیری ماشین بتوانند هنگامی که p بسیار بزرگتر از n است ($p \gg n$)، نقش خود را بطور مؤثر ایفا کنند. در این نوع داده‌ها p یک تابع نمایی از n است، به عبارت دیگر $\log(p) = O(n^\alpha)$ ، که در آن $\alpha > 0$. در این خصوص، مطالعات انجام شده بر اساس توابع تاوان که در بالا بدان اشاره شد، به دلایل هزینه محاسباتی، دقت آماری و ناپایداری الگوریتمی از کارایی کافی برخوردار نیستند و 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 راستا، قالب کلی راه حل موجود مبتنی بر رهیافت توابع تاوان، یک روش دو مرحله‌ای است که در مرحله اول ابتدا بعد مدل توسط یک روش غربالگری مستقل کارآ کاهش یافته و در مرحله دوم از روشی مبتنی بر دو تابع تاوان برای شناسایی مؤلفه‌های غیرصفر و خطی در زیر مدل غربال شده استفاده می‌شود. البته لازم به ذکر است که رهیافت استفاده از تابع تاوان، تنها راه موجود برای غربالگری و انتخاب متغیر نیست. به عنوان مثال، چانگ (۲۰۱۰) در قالب یک فن یادگیری ماشین به نام SLiM، رده‌بندی داده‌های بیان ژنی با بعد بالا را مورد بررسی قرار داد. همچنین غربالگری متغیرها در بعد بسیار بالا توسط وانگ (۲۰۰۹) و از منظر رگرسیون پیش رو مورد مطالعه قرار گرفت. فراتی و هال (۲۰۱۵) نیز به نقش آفرینی روش ناپارامتری رگرسیون موضعی در غربالگری و انتخاب متغیرها پرداختند. مسئله انتخاب متغیر در داده‌های با بعد بالا، از دیدگاه بیزی و توسط فرآیند پواسون – دیریکله نیز مورد بحث قرار گرفت (گوها و بلادان دیوزاپانی، ۲۰۱۶).

⁴Least Absolute Shrinkage and Selection Operator

⁵High dimensional

⁶Ultrahigh dimensional

در رابطه با مرحله اول روش دو مرحله‌ای مذکور، فن و لیو (۲۰۰۸) روش غربالگری مستقل مطمئن^۷ (SIS) را با استفاده از همبستگی پیرسن برای مدل‌های خطی با بعد بسیار بالا معرفی کردند. پس از آن، روش SIS به مدل‌های آماری مختلف مانند مدل خطی تعمیم‌یافته (فن و همکاران ۲۰۰۹، فن و سانگ ۲۰۱۰) و مدل جمعی ناپارامتری (فن و همکاران ۲۰۱۱، کاظمی و همکاران ۲۰۱۷) تعمیم داده شد. این روش‌ها مدل مبنا بوده و در صورت نادرست بودن مدل آماری مفروض، ممکن است دارای عملکرد مناسبی نباشند. بنابراین برای اجتناب از تشخیص نادرست ساختار مدل، معرفی یک روش غربالگری آزاد - مدل، که به نوع مدل بستگی نداشته باشد، ضروری به نظر می‌رسید. لذا، ژو و همکاران (۲۰۱۱) یک روش غربالگری آزاد- مدل به نام غربالگری و رتبه‌بندی مستقل مطمئن^۸ (SIRS) را معرفی کردند. این روش را می‌توان برای بسیاری از مدل‌های پارامتری و نیمه‌پارامتری مانند مدل خطی، مدل خطی تعمیم یافته، مدل شاخص و بسیاری از مدل‌های متداول دیگر به‌کار برد.

با توجه به اهمیت روش‌های غربالگری آزاد-مدل، یک روش دیگر موسوم به غربالگری مستقل مطمئن براساس همبستگی فاصله‌ای^۹ (DC-SIS) توسط لی و همکاران (۲۰۱۲) براساس همبستگی فاصله‌ای بین متغیرهای توضیحی و متغیر پاسخ ارائه شد. آنها نشان دادند که این روش دارای ویژگی غربالگری مطمئن^{۱۰} است، یعنی با احتمال نزدیک به یک متغیرهای مهم را برای ورود به مدل انتخاب می‌کند. روش DC-SIS دارای چندین مزیت است: الف. برای پیاده سازی 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 این مقاله ضمن رعایت الگوی کلی و متعارف دو مرحله‌ای ذکر شده، ابتدا از یک روش غربالگری مستقل مطمئن بر اساس همبستگی فاصله‌ای بین متغیرهای توضیحی و تابع توزیع حاشیه‌ای متغیر پاسخ برای کاهش بعد مدل استفاده می‌شود. این روش آزاد - مدل بوده و انتظار می‌رود وقتی که متغیر پاسخ دم سنگین، چوله یا دارای مقادیر فرین باشد، عملکرد خوبی داشته باشد. این روش غربالگری تعمیمی از روش ناستوار لی و همکاران (۲۰۱۲) است که به جای Y از $F(Y)$ ، تابع توزیع حاشیه‌ای متغیر پاسخ، استفاده می‌شود. سپس دو تابع تاوان SCAD گروهی را برای تشخیص همزمان متغیرهای مهم و مؤلفه‌های خطی در زیر مدل غربال شده در مرحله قبل به‌کار برده می‌شود.

⁷ Sure independence screening

⁸ Sure independence ranking and screening

⁹ Sure independence screening based on the distance correlation

¹⁰ Sure screening property

در بخش ۲ تعمیمی از روش لی و همکاران (۲۰۱۲) برای کاهش بعد معرفی می شود. در بخش ۳ روش انتخاب متغیر و تشخیص ساختار با استفاده از دو تابع تاوان با جزئیات ارائه می شود. در بخش ۴ در مطالعه‌ای شبیه‌سازی به بررسی عملکرد روش معرفی شده پرداخته می شود. در بخش ۵ تحلیل یک مجموعه داده واقعی ارائه می شود. بحث و نتیجه‌گیری کلی بخش پایانی مقاله حاضر است.

۲ غربالگری مستقل مطمئن نیرومند بر اساس همبستگی فاصله‌ای

در این مقاله مدل (۲) بررسی می شود، اما از آنجا که در مسائل واقعی مؤلفه‌های خطی و غیرخطی مشخص نیستند، ابتدا فرض می شود که مدل به صورت جمعی ناپارامتری (۱) است، سپس روشی ارائه خواهد شد تا مؤلفه‌های خطی در مدل (۱) تعیین شوند و صورت کلی رابطه (۲) بدست آید. برای این منظور، ابتدا با استفاده از یک روش غربالگری نیرومند به کاهش بعد فضای متغیرها به یک بعد کوچکتر از n پرداخته می شود. این روش غربالگری بر اساس همبستگی فاصله‌ای بین متغیرهای توضیحی و تابع توزیع حاشیه‌ای متغیر پاسخ تعریف می شود. استفاده از تابع توزیع متغیر پاسخ به جای خود متغیر پاسخ باعث می شود که روش ارائه شده نیرومند باشد. دلیل استفاده از روش غربالگری این است که روش‌های تاوان برای تشخیص اجزای خطی و غیرخطی در مدل جمعی ناپارامتری با بعد بسیار بالا مستقیماً قابل استفاده نیستند. لازم است ابتدا با یک روش غربالگری بعد مدل را کاهش داده، سپس از توابع تاوان برای تشخیص ساختار مدل استفاده شود.

سزکلی و همکاران (۲۰۰۷) همبستگی فاصله‌ای را به عنوان یک معیار وابستگی بین دو بردار تصادفی معرفی کردند. همبستگی فاصله‌ای بین بردارهای تصادفی $U \in \mathbb{R}^q$ و $V \in \mathbb{R}^r$ ، با گشتاورهای مرتبه اول متناهی، یک عدد نامنفی به صورت

$$dcorr(U, V) = \frac{dcov(U, V)}{\sqrt{dcov(U, U)dcov(V, V)}}$$

تعریف می شود، که در آن $dcov(U, V)$ کوواریانس فاصله‌ای U و V است و ثابت کردند

$$dcov^2(U, V) = S_1 + S_2 - 2S_3$$

که در آن

$$\begin{aligned} S_1 &= E\{\|U - \tilde{U}\|_q \|V - \tilde{V}\|_r\}, \\ S_2 &= E\{\|U - \tilde{U}\|_q\} E\{\|V - \tilde{V}\|_r\}, \\ S_3 &= E\{E(\|U - \tilde{U}\|_q | U) E(\|V - \tilde{V}\|_r | V)\}, \end{aligned}$$

و $\|\cdot\|$ نرم اقلیدسی و (\tilde{U}, \tilde{V}) بردارهای تصادفی مستقل از (U, V) و هم‌توزیع با آنها هستند. سزکلی و همکاران (۲۰۰۷) نشان دادند که U و V مستقل‌اند اگر و تنها اگر $dcorr(U, V) = 0$. همچنین $dcorr(U, V)$ تابعی اکیدا صعودی از قدر مطلق همبستگی پیرسن بین U و V است. به دلیل این دو ویژگی، 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. آنها همچنین نشان دادند که این روش دارای ویژگی غربالگری مطمئن است.

در این مقاله 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

$$\omega_k = dcorr(X_k, F(Y)), \quad k = 1, \dots, p, \quad (3)$$

تعریف می‌شود، که در آن $F(y)$ تابع توزیع حاشیه‌ای Y است. این روش غربالگری نسبت به روش‌های موجود دارای دو مزیت است: الف. با توجه به ویژگی‌های همبستگی فاصله‌ای، X_k و $F(Y)$ مستقل‌اند، اگر و تنها اگر $dcorr(X_k, F(Y)) = 0$. 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) براساس یک نمونه تصادفی برآورد شود. فرض کنید $(X_i, Y_i)_{i=1}^n$ یک نمونه تصادفی از مدل جمعی ناپارامتری (۱) باشد.

مقادیر S_1, S_2, S_3 با روش گشتاوری به صورت

$$\begin{aligned}\hat{S}_{k,1} &= \frac{1}{n^r} \sum_{i=1}^n \sum_{j=1}^n \|X_{ik} - X_{jk}\|_q \|F_n(Y_i) - F_n(Y_j)\|_r, \\ \hat{S}_{k,2} &= \frac{1}{n^r} \sum_{i=1}^n \sum_{j=1}^n \|X_{ik} - X_{jk}\|_q \frac{1}{n^r} \sum_{i=1}^n \sum_{j=1}^n \|F_n(Y_i) - F_n(Y_j)\|_r, \\ \hat{S}_{k,3} &= \frac{1}{n^r} \sum_{i=1}^n \sum_{j=1}^n \sum_{l=1}^n \|X_{ik} - X_{lk}\|_q \|F_n(Y_j) - F_n(Y_l)\|_r,\end{aligned}$$

برآورد می‌شود، که در آن $F_n(y) = \frac{1}{n} \sum_{i=1}^n I(Y_i \leq y)$ تابع توزیع تجربی Y است. بنابراین، برآورد $dcov^2(X_k, F(Y))$ به صورت

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

است. همچنین کوواریانس‌های فاصله‌ای نمونه‌ای $dcov(X_k, X_k)$ و $dcov(F(Y), F(Y))$ را نیز می‌توان به طور مشابه برآورد نمود. در نتیجه همبستگی فاصله‌ای نمونه‌ای بین X_k و $F(Y)$ برابر است با

$$\widehat{dcorr}(X_k, F(Y)) = \frac{\widehat{dcov}(X_k, F(Y))}{\sqrt{\widehat{dcov}(X_k, X_k)} \sqrt{\widehat{dcov}(F(Y), F(Y))}}.$$

بنابراین $\hat{\omega}_k = \widehat{dcorr}(X_k, F(Y))$ برای کاهش بعد فضای متغیرها، مجموعه‌ای از متغیرهای توضیحی را که دارای مقادیر بزرگ $\hat{\omega}_k$ هستند، به عنوان مجموعه متغیرهای مهم انتخاب می‌شوند. فرض کنید \hat{M} مجموعه اندیس متغیرهای مهم باشد، در نتیجه

$$\hat{M}_{\nu_n} = \{1 \leq k \leq p : \hat{\omega}_k \geq \nu_n\}, \quad (4)$$

که در آن ν_n یک عدد مثبت از پیش تعیین شده است و در زیربخش ۱.۲ به نحوه انتخاب آن پرداخته خواهد شد. این روش بعد فضای متغیرها را از p به یک فضای بسیار کوچکتر با اندازه $d = |\hat{M}_{\nu_n}|$ کاهش می‌دهد. این روش پیشنهادی، غربالگری مستقل مطمئن نیرومند براساس همبستگی فاصله‌ای نامیده می‌شود. در ادامه برای سادگی، از نمایش اختصاری RDC-SIS استفاده می‌شود.

۱.۲ قانون آستانه

استفاده از روش غربالگری RDC-SIS مستلزم تعیین یک مقدار آستانه^{۱۱} معقول ν_n است که در عمل، تعیین مقدار آن معمولاً مشکل است. در راستای یافتن راهکاری مناسب برای این مسئله، یک روش جایگزین، انتخاب d متغیر با بیشترین مطلوبیت حاشیه‌ای یا بیشترین همبستگی است. انتخاب مقدار d نقش مهمی را در مرحله غربالگری ایفا می‌کند. فن و لیو (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. این مقادیر ممکن است در برخی شرایط خوب عمل کنند اما دارای دو عیب عمده‌اند:

الف. هنوز بطور واضح مقدار دقیق d مشخص نیست. برای یک مجموعه داده واقعی، دقیقاً مشخص نیست کدام یک از مقادیر d_1 ، d_2 ، d_3 ، یا حتی یک مقدار بزرگتر باید استفاده شود.

ب. فرمول $d = [n/\log(n)]$ تنها به اندازه نمونه، n ، وابسته است و تعداد متغیرهای توضیحی را نادیده می‌گیرد. در این راستا، ژائو و لی (۲۰۱۲) روشی را برای انتخاب d در مدل کاکس معرفی کردند، اما روش آنها صرفاً برای روش‌های غربالگری مدل-مینا قابل استفاده است.

ژو و همکاران (۲۰۱۱) با افزودن q متغیر کمکی به مجموعه داده‌ها، یک روش دیگر برای تعیین d در روش غربالگری SIRS پیشنهاد دادند. مشکل عمده این روش، انتخاب تعداد متغیرهای کمکی، q ، است. آنها بطور تجربی مقدار $q = p$ را انتخاب کردند و در مطالعه‌ای شبیه‌سازی نشان دادند که این مقدار q مناسب است.

در عمل باید تعیین شود کدام یک از دو روش فوق برای انتخاب d مناسب‌تر است. ژو و همکاران (۲۰۱۱) نشان دادند وقتی مدل واقعی بسیار تنک است، یا به عبارتی دیگر تعداد متغیرهای توضیحی مهم بسیار اندک است، روش فن و لیو (۲۰۰۸) به روش ژو و همکاران (۲۰۱۱) برتری دارد، اما وقتی تعداد متغیرهای با اهمیت زیاد است، روش ژو و همکاران (۲۰۱۱) نسبت به روش فن و لیو (۲۰۰۸) عملکرد بهتری دارد.

۲.۲ غربالگری مستقل مطمئن تکراری

روش‌های غربالگری مستقل با مطلوبیت حاشیه‌ای، تنها از اطلاعات حاشیه‌ای متغیرها به جای مدل کامل استفاده می‌کنند. لذا دو مسئله مهم ممکن است عملکرد این روش‌ها را با مشکل مواجه کند:

الف. در این روش‌ها، برخی از متغیرهای بی اهمیت که همبستگی بالایی با متغیرهای با اهمیت دارند،

¹¹Threshold

نسبت به سایر متغیرهای با اهمیت که همبستگی ضعیفی با متغیر پاسخ دارند، ارجحیت دارند. ب. متغیری که به صورت حاشیه‌ای با متغیر پاسخ ناهمبسته اما به صورت توأم و از طریق سایر متغیرها با متغیر پاسخ همبسته است، توسط این روش‌ها انتخاب نمی‌شود. فن و لیو (۲۰۰۸) نشان دادند که در صورت وجود دو مشکل فوق، با به کار بردن SIS ممکن است برخی متغیرهای مهم از دست داده شوند. آنها برای رفع مشکل و افزایش کارایی انتخاب متغیر، روش SIS را به صورت مکرر برای مدل‌های خطی به کار گرفتند و این روش را ISIS نامگذاری کردند (فن و لیو، ۲۰۰۸). تاثیر روش‌های SIS و ISIS به فرض خطی بودن مدل بستگی دارد، 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. ارائه یک روش تکراری برای افزایش کارایی RDC-SIS با یک مثال شروع می‌شود. بدین منظور، مدل

$$Y = \delta X_1 + \delta X_2 + \delta X_3 - 15\sqrt{\rho}X_4 + \varepsilon \quad (5)$$

را در نظر بگیرید که هر کدام از متغیرهای توضیحی آن از توزیع $N(0, 1)$ تولید می‌شوند. در اینجا فرض می‌شود که به استثنای X_4 ، ضریب همبستگی هر متغیر با سایر متغیرها یکسان و برابر با $\rho \neq 0$ است و X_4 دارای همبستگی $\sqrt{\rho}$ با $p - 1$ متغیر دیگر است. در این مثال، متغیر X_4 بطور توأم مهم اما بطور 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 تاثیر 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_k روی (X_1, X_2, X_3) به صورت خطی برای $k = 4, \dots, p$ است. باقیمانده‌های بدست آمده از این رگرسیون‌های خطی با (X_1, X_2, X_3) ناهمبسته هستند. سپس روش RDC-SIS برای این مانده‌ها (به جای متغیر X_k) و Y به کار برده می‌شود. مانده‌های 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 می‌شود که M_1 مجموعه اندیس متغیرهای انتخاب شده با اندازه $d_1 < d$ و یک مقدار از پیش تعیین شده است. در اینجا برای

سهولت از $d = 2 \lceil n / \log(n) \rceil$ استفاده می‌شود.

• گام ۲: فرض کنید \mathbf{X}_1 ماتریس طرح متناظر با متغیرهای مجموعه \mathcal{M}_1 و \mathbf{X}_1^c ماتریس طرح متناظر با متغیرهای مجموعه \mathcal{M}_1^c است. ماتریس $\mathbf{X}_{new}^c = \{I_n - \mathbf{X}_1(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T\} \mathbf{X}_1^c$ را محاسبه کنید و سپس روش RDC-SIS را برای Y و تمام ستونهای ماتریس \mathbf{X}_{new}^c به کار برید. فرض کنید در این مرحله d_1 متغیر توضیحی انتخاب می‌شوند و مجموعه اندیس متغیرهای انتخاب شده را با \mathcal{M}_2 نشان می‌دهیم. مجموعه \mathcal{M}_1 را با $\mathcal{M}_1 \cup \mathcal{M}_2$ بروز رسانی کنید.

• گام ۳: گام ۲ را $k - 1$ بار تکرار کنید تا تعداد متغیرهای توضیحی انتخاب شده از d تجاوز کند، یعنی $d_1 + \dots + d_k \geq d$. در پایان مجموعه متغیرهای توضیحی انتخاب شده $\mathcal{M}_1 \cup \dots \cup \mathcal{M}_k$ 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^{۱۲} استفاده می‌کند که در مقاله 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 تنظیم‌کننده در نظر گرفته می‌شود و مقدار بهینه آن با مینیمم کردن میانگین مربع خطای پیش‌بینی تعیین می‌گردد.

۳ انتخاب متغیر و تشخیص ساختار

انتخاب یک مقدار بزرگتر برای d ، احتمال انتخاب متغیرهای بی‌اهمیت را برای ورود به مدل افزایش می‌دهد. در این بخش، یک روش مبتنی بر تاوان برای حذف متغیرهای اضافی از مدل و تشخیص مؤلفه‌های خطی و غیرخطی ارائه می‌شود. فرض کنید d متغیر در مرحله غربالگری برای ورود به مدل انتخاب شده‌اند. مدل جمعی ناپارامتری

$$Y = \sum_{j=1}^d f_j(X_j) + \epsilon,$$

¹²Orthogonal projection

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

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

استفاده می‌شود یک روش رایج برای برآورد و انتخاب متغیر بطور همزمان، استفاده از رگرسیون تاوانیده است. در اینجا برای انتخاب متغیر و تشخیص ساختار مدل، از دو تابع تاوان به طور همزمان استفاده می‌شود، یعنی بردار ضرایب $b = (b_1^T, \dots, b_d^T)^T$ که $b_j = (b_{j1}, \dots, b_{jK})^T$ ، با حل مسئله بهینه‌سازی

$$\begin{aligned} \hat{b} = \arg \min_b & \frac{1}{n} \sum_{i=1}^n (Y_i - \mu - \sum_{j=1}^d \sum_{k=1}^K b_{jk} B_{jk}(X_{ij}))^2 \\ & + n \sum_{j=1}^d p_{\lambda_1}(\|b_j\|_{A_j}) + n \sum_{j=1}^d p_{\lambda_2}(\|b_j\|_{D_j}), \end{aligned} \quad (6)$$

برآورد می‌شود که در آن توابع تاوان $p_{\lambda_1}(\cdot)$ و $p_{\lambda_2}(\cdot)$ با پارامترهای تنظیم کننده λ_1 و λ_2 به ترتیب برای شناسایی مؤلفه‌های صفر و خطی به کار می‌روند. علاوه بر این، A_j و D_j دو ماتریس $K \times K$ و بسیار حائز اهمیت است. این ماتریس‌ها باید طوری انتخاب شوند که $\|b_j\|_{A_j} = 0$ باشد، اگر فقط $\sum_k b_{jk} B_{jk}(x) \equiv 0$ و بطور مشابه $\|b_j\|_{D_j} = 0$ باشد، اگر و فقط اگر $\int_0^1 B_{jk}''(x) B_{jk'}''(x) dx \}_{k,k'=1}^K$ انتخاب می‌کنند. انتخاب D_j به صورت فوق از این واقعیت نتیجه می‌شود که مشتق دوم یک تابع خطی برابر صفر است. فرض کنید

$$Z_j = \begin{pmatrix} B_{j1}(X_{1j}) & B_{j2}(X_{1j}) & \cdots & B_{jK}(X_{1j}) \\ \vdots & \vdots & & \vdots \\ B_{j1}(X_{nj}) & B_{j2}(X_{nj}) & \cdots & B_{jK}(X_{nj}) \end{pmatrix}_{n \times K},$$

$Y = (Y_1, \dots, Y_n)$ و $Z = (Z_1, \dots, Z_d)$ باشد. پس رابطه (۶) را می‌توان به صورت ماتریسی

$$\hat{b} = \arg \min_b \frac{1}{n} \|Y - Zb\|^2 + n \sum_{j=1}^d p_{\lambda_1}(\|b_j\|_{A_j}) + n \sum_{j=1}^d p_{\lambda_2}(\|b_j\|_{D_j}). \quad (7)$$

نوشت. به دلیل سهولت، تابع هدف سمت راست رابطه (۷) با $Q(b)$ نشان داده می‌شود.

تاکنون توابع تاوان زیادی توسط افراد مختلف معرفی شده‌است. سوالی که اینجا ممکن است پیش بیاید این است که از چه نوع تابع تاوان باید استفاده کرد. فن و لی (۲۰۰۱) نشان دادند که یک تابع تاوان خوب باید برآوردگری با سه ویژگی مطلوب زیر را نتیجه دهد:

۱. تنکی: برآوردگر نتیجه شده باید به طور خودکار ضرایب برآورد شده‌ای که مقدار کوچکی دارند را برابر صفر قرار دهد تا متغیرهای مناسب را انتخاب کند. این کار پیچیدگی مدل را کاهش می‌دهد.

۲. نارایی: برآوردگر به دست آمده برای ضرایب رگرسیونی که در واقع مقادیر آنها بزرگ است، تقریباً نارایب باشد. این ویژگی آریبی مدل را کاهش می‌دهد.

۳. پیوستگی: برآوردگر نتیجه شده دارای تغییرات پیوسته باشد و ناپایداری در پیش بینی مدل را کاهش دهد.

بنابراین، فن و لی (۲۰۰۱) یک تابع تاوان محدب به نام SCAD را معرفی کردند که دارای هر سه ویژگی مطلوب است. این تابع تاوان دارای مشتق مرتبه اول به صورت

$$p'_{a,\lambda}(x) = \lambda \left\{ I(|x| \leq \lambda) + \frac{(a\lambda - |x|)_+}{(a-1)\lambda} I(|x| > \lambda) \right\}, \quad x \geq 0,$$

است، که در آن $a > 2$ و $p'_{a,\lambda}(0) = 0$. فن و لی (۲۰۰۱) از دیدگاه بیزی مقدار $a = 3/7$ را به عنوان یک مقدار مناسب برای مسائل مختلف پیشنهاد دادند.

برای حل مسئله بهینه‌سازی (۷) از روش تکراری تقریب موضعی درجه دو^{۱۳} (LQA) معرفی شده توسط فن و لی (۲۰۰۱) استفاده می‌کنیم. با استفاده از بسط تیلور و داشتن یک مقدار اولیه b_j^0 ، توابع

¹³Local quadratic approximation

تاوان را می‌توان به صورت

$$p_{\lambda_1}(\|b_j\|_{A_j}) \approx p_{\lambda_1}(\|b_j^{(\circ)}\|_{A_j}) + \frac{1}{\gamma} \frac{p'_{\lambda_1}(\|b_j^{(\circ)}\|_{A_j})}{\|b_j^{(\circ)}\|_{A_j}} \{\|b_j\|_{A_j}^\gamma - \|b_j^{(\circ)}\|_{A_j}^\gamma\},$$

$$p_{\lambda_2}(\|b_j\|_{D_j}) \approx p_{\lambda_2}(\|b_j^{(\circ)}\|_{D_j}) + \frac{1}{\gamma} \frac{p'_{\lambda_2}(\|b_j^{(\circ)}\|_{D_j})}{\|b_j^{(\circ)}\|_{D_j}} \{\|b_j\|_{D_j}^\gamma - \|b_j^{(\circ)}\|_{D_j}^\gamma\}.$$

تقریب زد. با به کار بردن تقریب‌های فوق و حذف جملات ثابت می‌توان رابطه (۷) را به صورت

$$Q(b) = \frac{1}{n} \|Y - Zb\|^\gamma + \frac{1}{\gamma} b^T (\Omega_1 + \Omega_2) b \quad (۸)$$

خلاصه کرد، که در آن Ω_1 و Ω_2 دو ماتریس بلوکی با بعد $dK \times dK$ به صورت

$$\Omega_1 = \text{diag}\left(\frac{p'_{\lambda_1}(\|b_1^{(\circ)}\|_{A_1})}{\|b_1^{(\circ)}\|_{A_1}} A_1, \dots, \frac{p'_{\lambda_1}(\|b_d^{(\circ)}\|_{A_d})}{\|b_d^{(\circ)}\|_{A_d}} A_d\right)$$

$$\Omega_2 = \text{diag}\left(\frac{p'_{\lambda_2}(\|b_1^{(\circ)}\|_{D_1})}{\|b_1^{(\circ)}\|_{D_1}} D_1, \dots, \frac{p'_{\lambda_2}(\|b_d^{(\circ)}\|_{D_d})}{\|b_d^{(\circ)}\|_{D_d}} D_d\right)$$

هستند. توجه شود که رابطه (۸) یک تابع درجه دو برحسب b بوده و در نتیجه دارای جواب صریح

$$\hat{b} = (Z^T Z + n(\Omega_1 + \Omega_2))^{-1} Z^T Y,$$

است، به عنوان برآورد اولیه در تکرار بعدی استفاده می‌شود. این الگوریتم به طور مکرر مسئله بهینه‌سازی (۸) را حل می‌کند. یعنی برآورد $b^{(m)}$ توسط $b^{(m+1)}$ به روز می‌شود که $m = 0, 1, 2, \dots$. به عبارت دیگر، در تکرار $m + 1$ باید مسئله بهینه‌سازی (۸) حل شود، اما به جای استفاده از $b^{(\circ)}$ ، از برآورد $b^{(m)}$ به عنوان مقدار اولیه استفاده می‌شود. برآورد بدست آمده با مینیم کردن (۸) برآورد جدید ($b^{(m+1)}$) است. این الگوریتم تا رسیدن به همگرایی تکرار می‌شود. فن و لی (۲۰۰۱) نشان دادند که الگوریتم LQA پس از تعداد کمی تکرار به همگرایی می‌رسد. در هر تکرار به محض اینکه برخی از مؤلفه‌های $\|b_j\|_{A_j}$ یا $\|b_j\|_{D_j}$ کوچکتر از 10^{-6} شود، مولفه $f_j(x)$ به ترتیب به عنوان تابع صفر یا تابع خطی تشخیص داده می‌شود.

۴ مطالعه شبیه‌سازی

در این بخش، با شبیه‌سازی عملکرد روش دو مرحله‌ای بررسی می‌شود. ابتدا، در مثال ۱ حساسیت روش RDC-SIS نسبت به پارامتر d تحلیل می‌شود. سپس در مثال ۲، عملکرد RDC-SIS را 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 و ۴، 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 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 S و P_j باید نزدیک به یک باشند.

برای پیاده‌سازی RDC-ISIS، مطالعات تجربی نشان می‌دهد 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 مقاله، برای شبیه‌سازی با انتخاب $d_1 = 5$ و $d_2 = p - 5$ الگوریتم فقط یکبار تکرار می‌شود.

در مثال ۵ کارایی روش دو مرحله‌ای را در شناسایی مؤلفه‌های غیرصفر و همچنین تشخیص مؤلفه‌های خطی و غیرخطی در یک مدل جمعی خطی-جزئی با استفاده از دو تابع تاوان بررسی می‌شود. برای انتخاب بهینه پارامترهای λ_1 و λ_2 از معیار BIC به صورت

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

استفاده می‌شود، که در آن مقدار مینیمم‌کننده (۷) برای $(\lambda_1, \lambda_2) = \lambda$ داده شده است. همچنین d_1 تعداد مؤلفه‌های ناپارامتری و d_2 تعداد مؤلفه‌های پارامتری به ازای مقدار داده شده λ است.

مثال ۱: برای تحلیل حساسیت عملکرد RDC-SIS نسبت به مقادیر مختلف d ، داده‌ها از مدل

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

جدول ۱: احتمال تجربی P_j و احتمال S در مثال ۱

S	$\varepsilon \sim t(1)$					S	$\varepsilon \sim N(0, 1)$					d	n	p
	p_+	p_+	p_+	p_+	p_+		p_+	p_+	p_+	p_+				
۰/۴۱	۰/۴۵	۰/۹۲	۱/۰۰	۰/۸۴	۰/۵۳	۰/۶۰	۰/۹۶	۱/۰۰	۰/۹۱	d_1	۱۰۰			
۰/۵۸	۰/۶۲	۰/۹۶	۱/۰۰	۰/۹۴	۰/۷۳	۰/۷۵	۰/۹۸	۱/۰۰	۰/۹۸	d_2				
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	۲۰۰	۱۰۰۰		
۰/۹۷	۰/۹۷	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	d_2				
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	۴۰۰			
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	d_2				
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	۱۰۰			
۰/۵۸	۰/۶۶	۰/۹۶	۱/۰۰	۰/۸۹	۰/۷۱	۰/۷۲	۰/۹۹	۱/۰۰	۰/۹۶	d_2				
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	۲۰۰	۲۰۰۰		
۰/۹۷	۰/۹۷	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	d_2				
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	۴۰۰			
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	d_2				
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	d_3				

تولید شده‌اند، که در آن

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

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

$Cov(X_i, X_j) = 0.8^{|i-j|}$ همبستگی نرمال استاندارد و متغیرهای توضیحی دارای توزیع حاشیه‌ای نرمال استاندارد و همبستگی d هستند. سه مقدار مختلف $d_3 = 3[n/\log(n)]$, $d_2 = 2[n/\log(n)]$, $d_1 = [n/\log(n)]$ 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 گزارش شده است. همان‌طور که ملاحظه می‌شود، با 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 را پیشنهاد می‌شود. اما برای $n = 200, 400$ و هر دو نوع توزیع خطا، هر یک از مقادیر d_1, d_2, d_3 را می‌توان استفاده کرد. در مثال‌های ۲ تا ۴، حجم نمونه $n = 200$ ، تعداد متغیرهای توضیحی را $p = 1000$ و $d = 2[n/\log(n)]$ در نظر گرفته و شبیه‌سازی ۵۰۰ بار تکرار شده است.

جدول ۲: چندک‌های M ، احتمال تجربی P_j و احتمال S در مدل ۱

S	P					M					روش	خطا	c
	۵	۴	۳	۲	۱	%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	$N(0, 1)$
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	(۳۱)
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		

مثال ۲: دو مدل به صورت

مدل ۱: $Y = c\beta^T X + \sigma\varepsilon$ مدل ۲: $Y = X_1 + 2X_2 + 3X_3 + 4X_4 + \varepsilon$
 بگیریید، 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^{۱۴} مقادیر مختلفی برای c در نظر گرفته شده است. مقادیر $c = 1, 2$ که متناظر با $R^2 = 50\%, 80\%$ هستند، انتخاب شده‌اند. در هر دو مدل، بردار 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 $\mu = 0, \sigma = 1, \alpha = 2$ نیز در نظر گرفته شده و نتایج شبیه‌سازی پس از ۵۰۰ تکرار در جداول ۲ و ۳ ارائه شده‌اند.

با توجه به جدول ۲، 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_5, X_4, X_3, X_2, X_1 بسیار خوب عمل می‌کنند و همواره متغیرهای مهم را به درستی شناسایی می‌کنند، اما برای خطای غیر نرمال نتایج کاملاً متفاوت است. برای توزیع خطای تی-استودنت، عملکرد روش‌های DC-SIS و SIS بسیار ضعیف است، در حالیکه روش RDC-SIS با احتمال تجربی تقریباً ۱۰۰٪ متغیرهای مهم را به درستی تشخیص می‌دهد. همچنین روش SIRS نیز در شناسایی متغیرهای مهم خوب عمل می‌کند و عملکرد دو روش RDC-SIS و SIRS

¹⁴Signal-to-noise ratio

جدول ۳: چندک‌های M ، احتمال تجربی P_j و احتمال S در مدل ۲

S	P					M					روش	خطا
	۴	۳	۲	۱		%۹۵	%۷۵	%۵۰	%۲۵	%۵		
۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰		۴	۴	۴	۴	۴	RDC-SIS	$N(0, 1)$
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	

تقریباً یکسان است. با توجه به مقادیر P_j و S در جدول ۲ روش SIS شانس بسیار کمی برای انتخاب متغیرهای مهم دارد. این روش در حالت $c = 1$ با احتمال تجربی ۱۰٪ و در حالت $c = 2$ با احتمال ۲۳٪ همه متغیرهای مهم را انتخاب می‌کند.

جدول ۳ که حاوی نتایج مدل ۲ است، نشان می‌دهد که برای هر سه نوع توزیع خطا، عملکرد روش RDC-SIS بسیار خوب است و نسبت به روش‌های دیگر بهتر عمل می‌کند. برای $j = 1, \dots, 4$ 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 این موارد بهتر است از یک روش تکراری مناسب برای حذف متغیرهای بی‌اهمیت و بازگرداندن متغیرهای مهم به مدل استفاده شود.

مثال ۳: در این مثال، عملکرد روش RDC-ISIS را با روش‌های SIS، ISIS، SIRS، DC-SIS و RDC-SIS در مدل خطی (۵) مقایسه می‌شود. علاوه بر این، یک مدل خطی دیگر، با نسبت سیگنال به نوفه ضعیف‌تر، به صورت

$$Y = 2.5X_1 + 2.5X_2 + 2.5X_3 - 7.5\sqrt{\rho}X_4 + \varepsilon \quad (9)$$

در نظر گرفته شده و مدل (۵) را مدل ۱ و مدل (۹) را مدل ۲ می‌نامیم. در این مثال، $\varepsilon \sim N(0, 1)$

جدول ۴: احتمال تجربی P_j و احتمال S در مثال ۳

ρ	روش	مدل ۱					مدل ۲				
		P					P				
		S	۴	۳	۲	۱	S	۴	۳	۲	۱
۰٫۲	SIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	ISIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	SIRS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	DC-SIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	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	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
۰٫۵	SIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	ISIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	SIRS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	DC-SIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	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	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
۰٫۸	SIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	ISIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	SIRS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰
	DC-SIS	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	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	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰	۱/۰۰

و سه مقدار مختلف $\rho = ۰٫۲, ۰٫۵, ۰٫۸$ برای ضریب همبستگی در نظر گرفته شده و شبیه‌سازی ۵۰۰ بار تکرار شده است. در این مدل، متغیر $X_۴$ بطور توأم مهم اما بطور حاشیه‌ای با Y ناهمبسته است. بنابراین روش‌های غربالگری حاشیه‌ای SIS ، $SIRS$ ، $DC-SIS$ و $RDC-SIS$ به سختی می‌توانند متغیر مهم $X_۴$ را شناسایی کنند. اما همان‌طور که در جدول ۴ ملاحظه می‌شود روش $RDC-ISIS$ در هر دو مدل خطی و به ازای مقادیر مختلف ρ در انتخاب متغیر $X_۴$ بسیار تواناست. ما همچنین $RDC-ISIS$ را با $ISIS$ ، نسخه تکراری $ISIS$ ، مقایسه کرده‌ایم. نتایج نشان می‌دهد که در برخی موارد روش $RDC-ISIS$ در انتخاب $X_۴$ حتی از $ISIS$ نیز بهتر عمل می‌کند. به عنوان مثال، در مدل ۲ به ازای $\rho = ۰٫۸$ روش $RDC-ISIS$ با احتمال تجربی ۹۱٪ همه متغیرهای مهم را به درستی انتخاب می‌کند، در حالی که $ISIS$ 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 حالی که $ISIS$ 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 به‌کار رود.

مثال ۴: در این مثال، مدل رگرسیونی جمعی ناپارامتری به صورت

$$Y = 2f_1(X_1) + \sqrt{6}f_2(X_{1,01}) + 3f_3(X_{2,01}) - 0.6f_4(X_{2,02}) + \varepsilon,$$

است، که در آن توابع پایه‌ای بصورت

$$f_1(x) = \exp(2x/3), \quad f_2(x) = \begin{cases} x + 4 & x < -2 \\ |x| & |x| \leq 2 \\ 4 - x & x > 2, \end{cases}$$

$$f_3(x) = \frac{\sin(3\pi x/4 + 3/2)}{2 - \sin(3\pi x/4 + 3/2)}, \quad f_4(x) = \log(x^2).$$

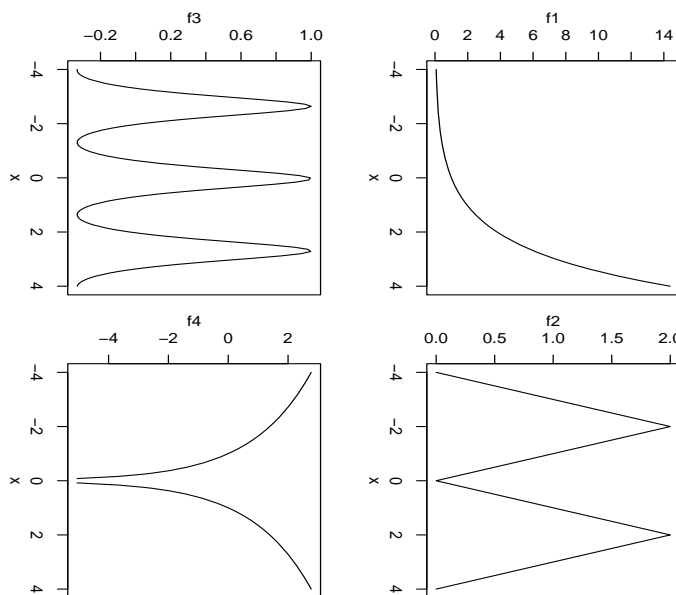
تعریف شده و نمودارهای آنها در شکل ۱ رسم شده‌اند. با توجه به مدل، متغیرهای $(X_1, X_{1,01}, X_{2,01}, X_{2,02})$ متغیرهای مؤثر بر پاسخ هستند. متغیرهای توضیحی از توزیع نرمال چند متغیره با میانگین صفر و همبستگی $Cov(X_i, X_j) = \rho^{|i-j|}$ تولید شده‌اند. سپس متغیر X_k با متغیر جدید ζ_k $X_k = 0.8X_1 + \zeta_k$ جایگزین می‌شود، که در آن $\zeta_k \sim N(0, 1)$ و $k = 2, \dots, 100$ ، یعنی 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 گزارش شده‌اند.

همان‌طور که ملاحظه می‌شود روش‌های غربالگری مستقل حاشیه‌ای SIS ، NIS ، $SIRS$ و DC در شناسایی متغیرهای مهم $X_{1,01}$ ، $X_{2,01}$ و $X_{2,02}$ ناتوان هستند، چون متغیرهای بی‌اهمیت یعنی (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 $X_{1,01}$ ، $X_{2,01}$ و $X_{2,02}$ اولویت دارند. بنابراین یک روش غربالگری تکراری مناسب برای حذف ۹۹ متغیر بی‌اهمیت و بازگرداندن سه متغیر مهم $X_{1,01}$ ، $X_{2,01}$ و $X_{2,02}$ به مدل نیاز است. با توجه به جدول ۵ ملاحظه می‌شود که روش $ISIS$ نمی‌تواند روابط غیرخطی بین متغیرهای توضیحی و متغیر پاسخ را شناسایی کند، در حالی‌که 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.

مثال ۵: در این مثال، ابتدا روش $RDC-SIS$ برای کاهش بعد به کار رفته، سپس با استفاده از دو تابع

جدول ۵: احتمال تجربی P_j و احتمال S در مثال ۴

ρ	روش	$\varepsilon \sim t(2)$					$\varepsilon \sim N(0,1)$				
		P					P				
		S	۴	۳	۲	۱	S	۴	۳	۲	۱
۰.۵	SIS	۰/۰۰	۰/۰۱	۰/۰۰	۰/۰۱	۰/۹۹	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۱/۰۰
	ISIS	۰/۰۰	۰/۰۵	۰/۰۵	۰/۰۶	۰/۹۸	۰/۰۰	۰/۰۵	۰/۰۶	۰/۰۷	۰/۹۸
	NIS	۰/۳۲	۰/۶۱	۰/۷۱	۰/۷۵	۰/۹۸	۰/۵۰	۰/۷۴	۰/۸۳	۰/۸۲	۱/۰۰
	SIRS	۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۱/۰۰	۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۱	۱/۰۰
	DC-SIS	۰/۰۳	۰/۲۷	۰/۳۴	۰/۳۱	۱/۰۰	۰/۰۲	۰/۲۴	۰/۲۹	۰/۲۵	۱/۰۰
	RDC-SIS	۰/۱۱	۰/۳۵	۰/۳۶	۰/۳۰	۱/۰۰	۰/۱۱	۰/۳۳	۰/۴۰	۰/۳۴	۱/۰۰
RDC-ISIS	۰/۸۲	۰/۸۸	۰/۹۹	۰/۹۴	۱/۰۰	۰/۹۶	۰/۹۹	۰/۹۹	۰/۹۸	۱/۰۰	
۰.۸	SIS	۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۰/۹۹	۰/۰۰	۰/۰۰	۰/۰۰	۰/۰۰	۱/۰۰
	ISIS	۰/۰۰	۰/۰۵	۰/۰۵	۰/۰۹	۰/۹۷	۰/۰۰	۰/۰۴	۰/۰۴	۰/۰۷	۰/۹۳
	NIS	۰/۵۹	۰/۸۸	۰/۹۳	۰/۷۲	۰/۹۹	۰/۸۱	۰/۹۶	۰/۹۹	۰/۸۵	۱/۰۰
	SIRS	۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۱/۰۰	۰/۰۰	۰/۰۱	۰/۰۱	۰/۰۰	۱/۰۰
	DC-SIS	۰/۱۴	۰/۶۲	۰/۷۲	۰/۳۰	۱/۰۰	۰/۱۳	۰/۶۱	۰/۷۸	۰/۲۸	۱/۰۰
	RDC-SIS	۰/۳۲	۰/۷۲	۰/۷۹	۰/۳۶	۱/۰۰	۰/۳۱	۰/۷۲	۰/۷۸	۰/۳۵	۱/۰۰
RDC-ISIS	۰/۸۳	۰/۹۵	۰/۹۸	۰/۹۱	۰/۹۸	۰/۹۰	۰/۹۳	۰/۹۸	۱/۰۰	۱/۰۰	



شکل ۱: نمودارهای توابع پایه‌ای برای مدل جمعی ناپارامتری مثال ۴ در فاصله $[-۴, ۴]$.

تاوان SCAD به‌طور همزمان یک مدل جمعی خطی-جزیی برازش داده می‌شود. داده‌ها از مدل

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

تولید شده‌اند، که در آن $f_2(x) = 6x(1-x)$ ، $f_1(x) = 3 \sin(2\pi x)/(2 - \sin(2\pi x))$ ، $f_3(x) = 2x$ ، $f_4(x) = x$ ، $f_5(x) = -x$ ، 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 همانند مثال ۲ تولید 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 $n = 30, 70, 200, 400$ و $p = 1000, 2000$ به تعداد ۲۰۰ بار تکرار شده است. نتایج انتخاب متغیر و تشخیص ساختار مدل در جدول ۶ 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. 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 به عنوان مؤلفه غیرصفر ناپارامتری تشخیص می‌دهد، اما برای خطای غیرنرمال در برخی از تکرارها بعضی از مؤلفه‌های ناپارامتری به عنوان مؤلفه صفر یا خطی تشخیص داده می‌شود. همچنین در برخی از تکرارها، مؤلفه‌های خطی به عنوان مؤلفه‌های صفر یا غیرخطی انتخاب می‌شوند. با توجه به جدول، هنگامی که حجم نمونه بسیار کوچک است، این روش عملکرد چندان مطلوبی ندارد. همانطور که می‌بینید برای $n = 30$ عملکرد این روش در انتخاب متغیر و تشخیص ساختار مدل چندان رضایت بخش نیست. 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.

۵ مثال کاربردی

مجموعه داده این مثال توسط بشل و همکاران (۲۰۰۷) مورد تحلیل و بررسی قرار گرفته و در بسته mixOmics نرم‌افزار R قابل دسترسی است. این داده‌ها مربوط به بیماری مسمومیت کبد بوده و شامل سطوح بیان ۳۱۱۶ ژن و ۹ اندازه‌گیری بالینی برای ۶۴ موش است. بنابراین در این داده‌ها ۳۱۱۶ متغیر توضیحی و ۹ متغیر پاسخ کمی وجود دارد. متغیرهای پاسخ به اختصار عبارتند از: BUN، نیتروژن اوره؛ TP، پروتئین کل؛ ALB، آلبومین؛ ALT، آلانین آمینوترانسفراز؛ SDH، سوربیتول دهیدروژناز؛

جدول ۶: نتایج تشخیص ساختار مدل در مثال ۵

NLT	NL	NNT	NN	N	خطا	n	p
۲/۶۰(۰/۸۷)	۲/۷۶(۱/۲۲)	۲(۰)	۲/۹۱(۱/۱۶)	۵/۶۷	N(۰, ۱)	۳۰	
۲/۶۱(۰/۹۸)	۳/۶۵(۱/۲۶)	۱/۹۵(۰/۲۳)	۳/۰۹(۱/۳۱)	۶/۷۴	t(۲)		
۲/۶۹(۰/۶۳)	۲/۷۴(۰/۸۶)	۲(۰)	۲/۶۸(۱/۰۴)	۵/۴۲	N(۰, ۱)	۷۰	
۲/۶۲(۰/۸۵)	۳/۶۵(۱/۱۵)	۱/۹۵(۰/۱۹)	۲/۹۶(۱/۱۰)	۶/۶۱	t(۲)		۱۰۰۰
۲/۷۴(۰/۳۲)	۲/۹۳(۰/۸۸)	۲(۰)	۲/۳۱(۰/۹۶)	۵/۲۴	N(۰, ۱)	۲۰۰	
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	N(۰, ۱)	۴۰۰	
۲/۷۴(۰/۷۵)	۳/۲۸(۰/۶۶)	۱/۹۹(۰/۰۸)	۲/۴۹(۰/۵۴)	۵/۷۷	t(۲)		
۲/۶۱(۱/۳۶)	۳/۶۴(۱/۳۹)	۲(۰)	۳/۱۱(۱/۳۳)	۶/۷۵	N(۰, ۱)	۳۰	
۲/۴۲(۱/۲۴)	۳/۶۸(۱/۵۴)	۱/۸۳(۱/۰۶)	۳/۷۰(۱/۴۶)	۷/۳۸	t(۲)		
۲/۸۵(۱/۰۵)	۳/۲۳(۱/۰۱)	۲(۰)	۲/۷۱(۱/۰۹)	۵/۹۴	N(۰, ۱)	۷۰	
۲/۷۲(۰/۹۸)	۳/۳۸(۱/۴۸)	۱/۹۰(۰/۵۸)	۳/۲۶(۱/۳۱)	۶/۶۴	t(۲)		۲۰۰۰
۲/۸۴(۰/۶۴)	۳/۱۲(۰/۹۴)	۲(۰)	۲/۴۶(۰/۸۴)	۵/۵۸	N(۰, ۱)	۲۰۰	
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	N(۰, ۱)	۴۰۰	
۲/۸۱(۰/۶۳)	۲/۹۷(۰/۸۵)	۱/۹۹(۰/۰۶)	۲/۵۴(۰/۷۴)	۵/۵۱	t(۲)		

جدول ۷: چهار ژن بسیار مهم انتخاب شده توسط روش‌های مختلف

ژن‌های انتخاب شده				روش
A_۴۲_P۴۹۶۹۵۵۱	A_۴۳_P۲۰۴۳۸	A_۴۳_P۱۴۱۶۳	A_۴۲_P۴۹۶۶۲۲	RDC-SIS
A_۴۲_P۴۹۶۹۵۵۱	A_۴۳_P۱۲۷۲۴	A_۴۳_P۲۰۴۳۸	A_۴۲_P۴۹۶۶۲۲	SIRS
A_۴۳_P۱۱۷۲۴	A_۴۲_P۸۴۰۷۷۶	A_۴۲_P۶۲۰۹۱۵	A_۴۳_P۱۴۱۳۱	DC-SIS
A_۴۳_P۱۴۱۳۱	A_۴۳_P۱۱۷۲۴	A_۴۳_P۱۰۶۰۶	A_۴۲_P۸۲۵۲۹۰	NIS

AST، آسپاراتات آمینوترانسفراز؛ ALP، آلکالین فسفاتاز؛ TBA، کل اسیدهای صفراوی؛ و CHOL، کلسترول. در این بخش ALT به عنوان متغیر پاسخ است و هدف پیدا کردن مؤثرترین ژنها در پیش‌بینی ALT و نوع تاثیر این ژنهاست.

جدول ۷ نشان دهنده ژن‌هایی است که توسط روش‌های غربالگری مختلف در رتبه‌های اول تا چهارم قرار می‌گیرند. به عنوان مثال، RDC-SIS و SIRS رتبه اول را به ژن A_۴۲_P۴۹۶۶۲۲ اختصاص می‌دهند. در مقابل، DC-SIS ژن A_۴۳_P۱۴۱۳۱ و NIS ژن A_۴۲_P۸۲۵۲۹۰ را به عنوان مهمترین ژن انتخاب می‌کنند. برای مقایسه کارایی این روش‌ها، می‌توان یک مدل جمعی با ۴ متغیر به صورت

$$Y = g_{k1}(X_{k1}) + g_{k2}(X_{k2}) + g_{k3}(X_{k3}) + g_{k4}(X_{k4}) + \varepsilon_k, \quad k = 1, 2, 3, 4$$

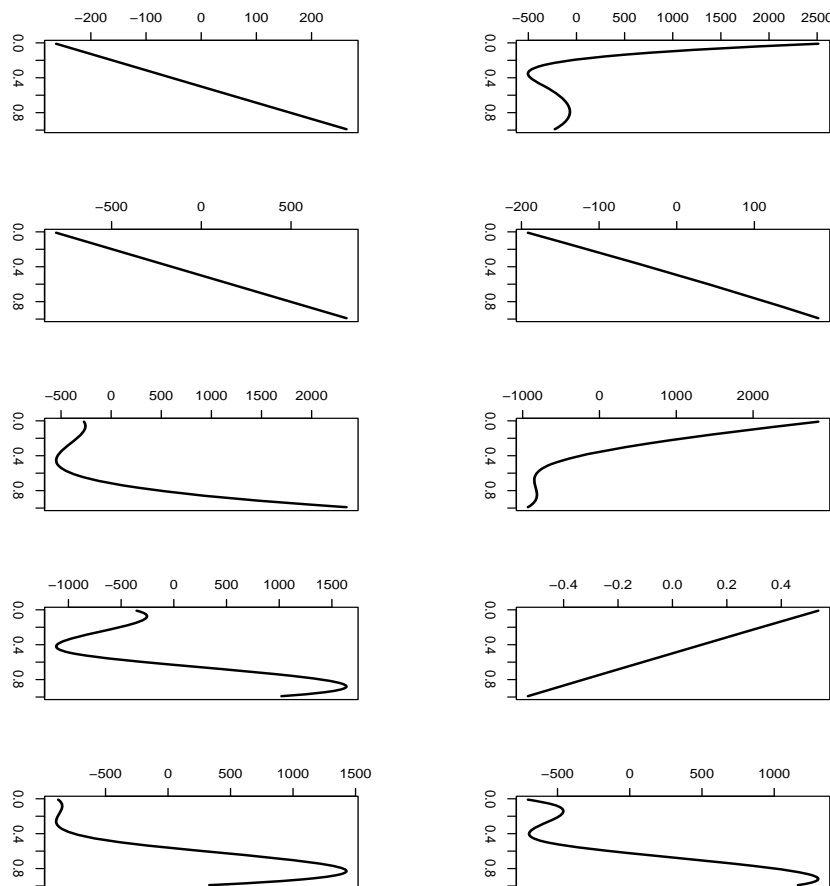
برازش داد، که در آن $X_{k1}, X_{k2}, X_{k3}, X_{k4}$ سطوح بیان ژن‌های انتخاب شده توسط هر کدام از روش‌های غربالگری، $j = 1, \dots, 4$ ، توابع پیوند نامعلوم هستند. برای برازش مدل فوق از

جدول ۸: 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

تابع gam در بسته mgcv نرم‌افزار R استفاده می‌کنیم. تابع gam برای برازش مدل جمعی تعمیم‌یافته استفاده می‌شود. برای مقایسه عملکرد نیکویی برازش از R^2 تعدیل شده مدل‌های جمعی برازش شده فوق استفاده می‌شود. پس از محاسبه R^2 تعدیل شده داریم: $R^2_{RDC-SIS} = 0.94$ ، $R^2_{SIRS} = 0.85$ ، $R^2_{NIS} = 0.89$ و $R^2_{DC-SIS} = 0.89$. با توجه به مقادیر R^2 ، برای این مجموعه داده، روش‌های DC-SIS و NIS عملکرد بسیار خوبی دارند و توانایی RDC-SIS نیز قابل قبول است، اما SIRS دارای عملکرد نسبتاً ضعیفی است. بنابراین RDC-SIS یک روش کارا برای کاهش بعد داده‌های با بعد بالاست. برای انتخاب مدل نهایی، در ابتدا نسخه تکراری RDC-SIS، یعنی RDC-ISIS با $d = 2[n/\log(n)] = 28$ و $p_1 = d/2$ به کار رفته و بعد داده‌ها را به $d = 28$ کاهش داده شده است. سپس از دو تابع تاوان SCAD برای انتخاب متغیر و تشخیص ساختار مدل استفاده شده است. 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. بعد از به کارگیری 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 ۸ گزارش شده است. نمودار توابع اثرات ژن‌های انتخاب شده را نیز در شکل ۲ می‌توان ملاحظه نمود.

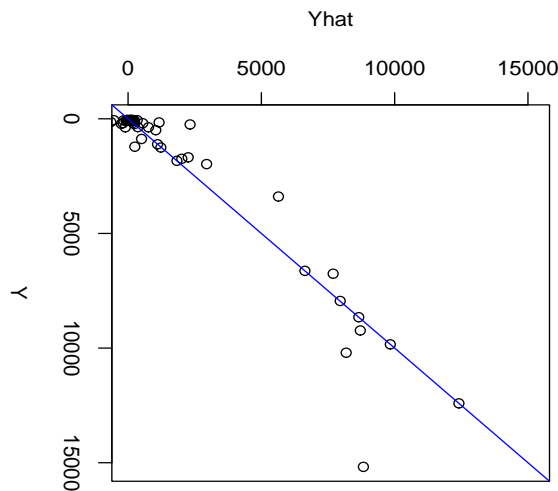
برای ارزیابی روش دو مرحله‌ای ارائه شده، از اعتبارسنجی متقابل استفاده می‌کنیم. ابتدا یک مشاهده را به عنوان داده آزمون و سایر مشاهدات را به عنوان داده‌های آموزشی در نظر بگیریم. پس از برازش مدل با استفاده از داده‌های آموزشی، مقدار متغیر پاسخ به ازای داده آزمون پیش بینی می‌شود. این روند برای تمام مشاهدات تکرار شده و نمودار مقادیر مشاهده شده متغیر پاسخ در مقابل مقادیر پیش بینی شده در شکل ۳ رسم شده است. با توجه به نمودار، روش معرفی شده در پیش بینی مقادیر ALT عملکرد خوبی دارد.



شکل ۲: توابع رگرسیونی برازش شده برای ژن‌های مهم

بحث و نتیجه‌گیری

در این مقاله، یک روش دو مرحله‌ای برای انتخاب متغیر و تشخیص ساختار در مدل‌های جمعی خطی-جزیی با بعد بالا ارائه شد. در مرحله اول، یک روش غربالگری مستقل مطمئن برای کاهش بعد مدل استفاده شد که در آن متغیرهای توضیحی براساس میزان همبستگی فاصله‌ای آنها با تابع توزیع حاشیه‌ای متغیر پاسخ رتبه‌بندی می‌شوند. کارایی این روش غربالگری در مطالعه‌ای 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 بالا،



شکل ۳: نمودار مقادیر مشاهده شده متغیر پاسخ در مقابل مقادیر پیش بینی شده

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 تاوان 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 مقاله مدل جمعی خطی-جزیی مد نظر بوده است اما روش ارائه شده برای غربالگری و تشخیص ساختار را می‌توان به مدل‌های نیمه‌پارامتری دیگر مانند مدل‌های با ضریب متغیر خطی-جزیی و مدل‌های جمعی تعمیم یافته خطی-جزیی تعمیم داد. همچنین لازم به ذکر است که در این مقاله متغیرهای توضیحی و پاسخ صرفاً از نوع کمی در نظر گرفته شده است. موضوع غربالگری در حالتی که این متغیرها از نوع کیفی و چند سطحی است، می‌تواند موضوعی برای تحقیقات آینده در نظر گرفته شود. البته تحت یک روش غربالگری متفاوت، مسئله غربالگری متغیرهای کیفی توسط هوآنگ و همکاران (۲۰۱۴) مورد بررسی قرار گرفته است.

در خصوص تقلیل بعد بالای فضای متغیرهای توضیحی به بعد مرتبه d توسط روش غربالگری بکار رفته در این مقاله می‌توان گفت که انتخاب d ، همانند انتخاب پارامتر تاوان 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 تاوان در تشخیص ساختار مدل، مسئله دیگری است که بایستی به آن توجه شود. اگر چه ما از تابع تاوان SCAD استفاده کرده‌ایم، اما می‌توان از ترکیب روش غربالگری با سایر توابع تاوان نظیر LASSO، MCP یا LASSO تطبیقی نیز استفاده کرد.

تقدیر و تشکر

نویسندگان مقاله از داوران محترم که نظرات ارزشمند ایشان باعث بهبود مطالب ارائه شده در این مقاله گردید، کمال تشکر و قدردانی را دارند.

مراجع

- Bushel, P., Wolfinger, R. D., Gibson, G. (2007), Simultaneous Clustering of Gene Expression Data with Clinical Chemistry and Pathological Evaluations Reveals Phenotypic Prototypes, *BMC Systems Biology*, **1**. doi:10.1186/1752-0509-1-15.
- Cheng, Q. (2010), A Sparse Learning Machine for High-dimensional Data with Application to Microarray Gene Analysis, *IEEE/ACM transactions on computational biology and bioinformatics*, **7**, 636-646.
- 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.
- Fan, J., and Li, R., (2001), Variable Selection via Nonconcave Penalized Likelihood and It's Oracle Properties, *Journal of the American Statistical Association*, **96**, 1348-1360.
- Fan, J., and Lv, J. (2008), Sure Independence Screening for Ultrahigh Dimensional Feature Space, *Journal of the Royal Statistical Society: Series B*, **70**, 849-911.
- Fan, J., Samworth, R. J., and Wu, Y. (2009), Ultrahigh Dimensional Feature Selection: Beyond the Linear Model, *Journal of Machine Learning Research*, **10**, 1829-1853.
- Fan, J., Feng, Y., and Song, R. (2011), Nonparametric Independence Screening in Sparse Ultrahigh-Dimensional Additive Models, *Journal of the American Statistical Association*, **106**, 544-557.

- Fan, J., and Song, R. (2010), Sure Independence Screening in Generalized Linear Models with NP-dimensionality, *Annals of Statistics*, **6**, 3567-3604.
- Ferraty, F., and Hall, P. (2015), An Algorithm for Nonlinear, Nonparametric Model Choice and Prediction, *Journal of Computational and Graphical Statistics*, **24**, 695-714.
- Guha, S., and Baladandayuthapani, V. (2016), A Nonparametric Bayesian Technique for High-dimensional Regression. *Electronic Journal of Statistics*, **10**, 3374-3424.
- 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 & Data Analysis*, **65**, 56-67.
- 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**, 533-550.
- Huang, J., Wei, F., and Ma, S. (2012), Semiparametric Regression Pursuit, *Statistica Sinica*, **22**, 1403-1426.
- Huang, D., Li, R., and Wang, H. (2014), Feature Screening for Ultrahigh Dimensional Categorical Data with Applications, *Journal of Business and Economic Statistics*, **32**, 237-244.
- Kazemi, M., Shahsavani, D., and Arashi, M. (2017), A Sure Independence Screening Procedure for Ultra-high Dimensional Partially Linear Additive Models, arXiv preprint arXiv:1708.08604.
- Li, R.Z., Zhong, W., and Zhu, L.P. (2012), Feature Screening via Distance Correlation Learning, *Journal of the American Statistical Association*, **107**, 1129-1139.
- Lian, H. (2012), Variable Selection in High-dimensional Partly Linear Additive Models, *Journal of Nonparametric Statistics*, **24**, 825-839.
- Lian, H. (2012), Shrinkage Estimation for Identification of Linear Components in Additive Models, *Statistics & Probability Letters*, **82**, 225-231.
- Lian, H., Chen, X., & Yang, JY. (2012), Identification of Partially Linear Structure in Additive Models with an Application to Gene Expression Prediction from Sequences. *Biometrics* **68**, 437-445.
- 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*, **25**, 591-607.

- Liu, X., Wang, L., and Liang, H., (2011), Estimation and Variable Selection for Semiparametric Additive Partial Linear Models, *Statistica Sinica*, **21**, 1225-1248.
- Lu, J., Yang, H., and Guo, C. (2016), Variable Selection in Partially Linear Additive Models for Modal Regression, *Communication in Statistics-Simulation and Computation*, DOI: 10.1080/03610918.2016.1171346.
- Opsomer, J. D., and Ruppert, D. (1999), A Root-n Consistent Backfitting Estimator for Semiparametric Additive Modeling. *Journal of Computational and Graphical Statistics*, **8**, 715-732.
- Segal, M. R., Dahlquist, K. D., and Conklin, B. R. (2003). Regression Approach for Microarray Data Analysis, *Journal of Computational Biology*, **10**, 961-980.
- Zhao, S. D., and Li, Y. (2012), Principled Sure Independence Screening for Cox Models with Ultra-High-Dimensional Covariates, *Journal of Multivariate Analysis*, **105**, 397-411.
- Szekely, G. J., Rizzo, M. L., and Bakirov, N. K. (2007), Measuring and Testing Dependence by Correlation of Distances, *Annals of Statistics*, **35**, 2769-2794.
- Tibshirani, R. (1996), Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society: Series B*, **58**, 267-288.
- Wang, H. (2009), Forward Regression for Ultra-high Dimensional Variable Screening, *Journal of the American Statistical Association*, **104**, 1512-1524.
- Zhang, C. H. (2010), Nearly Unbiased Variable Selection under the Minimax Concave Penalty, *Annals of Statistics*, **38**, 894-942.
- 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**, 1099-1112.
- Zhu, L.P., Li, L., Li, R., and Zhu, L.X., (2011), Model-free Feature Screening for Ultrahigh Dimensional Data, *Journal of the American Statistical Association*, **106**, 1464-1475.