

مدل‌بندی داده‌های تابعی با رویکرد رگرسیون مؤلفه اصلی بر اساس معیار اعتبار سنجی متقابل تعمیم‌یافته

مهدی روزبه^۱، آرتا روحی^۲، فاطمه جهادی^۳

تاریخ دریافت: ۱۴۰۰/۰۴/۱۲

تاریخ پذیرش: ۱۴۰۲/۰۲/۰۸

چکیده:

تحلیل داده‌های تابعی برای توسعه رویکردهای آماری در داده‌هایی مورد استفاده قرار می‌گیرد که دارای ماهیت تابعی و پیوسته هستند و چون این توابع به فضاها با بعد بی‌نهایت تعلق دارند، استفاده از روش‌های متداول در آمار کلاسیک برای تحلیل آن‌ها، با چالش روبرو است. مشهورترین تکنیک تحلیل داده‌های آماری، رویکرد مؤلفه‌های اصلی تابعی می‌باشد که ابزاری مهم برای کاهش بعد است، در این مقاله با استفاده از روش رگرسیون مؤلفه اصلی تابعی بر اساس جریمه مشتق دوم، ریچ و لاسو به تحلیل داده‌های تابعی آب‌وهوای کانادا و داده‌های تابعی طیف‌سنج پرداخته خواهد شد. بدین منظور برای تعیین مقدار بهینه پارامتر جریمه در روش‌های مورد استفاده از اعتبار سنجی متقابل تعمیم‌یافته، که معیاری معتبر و کارآمد است، استفاده می‌گردد.

واژه‌های کلیدی: اعتبار سنجی متقابل تعمیم‌یافته، رگرسیون تابعی، رگرسیون مؤلفه اصلی، تحلیل داده‌های تابعی.

۱ مقدمه

۱.۱ روش مؤلفه‌های اصلی

روش مؤلفه‌های اصلی، از روش‌های کاربردی با ایده‌ی کاهش ابعاد و حفظ بیشترین اطلاعات ممکن از متغیرهای توضیحی است. در این روش به دنبال ترکیب خطی از متغیرهای توضیحی هستیم که واریانس را به بیشترین مقدار خود برساند. در این روش از مقادیر ویژه و بردارهای ویژه ماتریس واریانس-کوواریانس یا ماتریس همبستگی استفاده می‌شود [۱۲]. از دیدگاه هندسی، این روش تبدیل خطی متعامدی است که داده‌ها را از یک دستگاه مختصات به دستگاه مختصات جدید می‌برد به نحوی که بزرگ‌ترین واریانس داده‌ها بر روی محور مختصات اول، دومین بزرگ‌ترین واریانس بر روی محور مختصات دوم قرار می‌گیرد و به همین ترتیب ادامه می‌یابد. در نتیجه با این کار متغیرها با کاهش بعد به مؤلفه‌هایی تبدیل می‌شوند که مؤلفه‌های اصلی نامیده می‌شوند. مؤلفه‌های اصلی ضمن ناهمبستگی به نحوی سازمان‌دهی می‌شوند که تعداد کمی از مؤلفه‌ها بتوانند درصد قابل توجهی از تغییرات متغیرهای توضیحی اولیه را توجیه کنند. انتخاب تعداد مناسب مؤلفه‌ها مورد توجه است و روش‌های مختلفی برای تعیین تعداد مناسب این مؤلفه‌ها پیشنهاد شده که در زیر به برخی از آن‌ها اشاره می‌شود. یکی از روش‌های

در داده‌های تابعی نخستین گام تبدیل داده‌های گسسته به پیوسته است که بدین منظور از روش‌های هموارسازی مانند پایه‌هایی فوریه برای داده‌های دارای دوره تناوب و اسپلاین بهره گرفته خواهد شد و سپس به برآورد ضرایب رگرسیونی در آن پرداخته می‌شود. در ادامه چون در عمل داده‌ها از قسمتی عددی و قسمتی تابعی تشکیل شده‌اند، تعدادی از مدل‌های رگرسیون تابعی بیان می‌شوند که به طور هم‌زمان بتوان پیش‌بینی‌های تابعی و غیرتابعی متعدد را کنترل و با جریمه مناسب، عامل‌های مهم خطا در پیش‌بینی را شناسایی کرد. به دلیل پرکاربرد بودن، روش‌های مختلفی ارائه شده است که هسته اصلی این روش‌ها، کاهش بعدهایی از پیش‌بینی‌کننده‌ها یا تنظیم منظم تابع ضریب در قالب یک جریمه ناهموار است که برای اطلاعات بیشتر می‌توان به هوروث و کوکوزکا [۵]، هزینگ و او بانک [۶] و در اواخر به کوکوزکا و ریمهر [۷]، فیرو بانده و همکاران [۸] و رییس و همکاران [۱۱] مراجعه کرد. در این مقاله علاوه بر رگرسیون مؤلفه اصلی تابعی از جریمه‌های کمینه کردن طول منحنی مشتق دوم، جریمه‌ی لاسو و ریچ استفاده شده است.

^۱ هیئت علمی گروه آمار، دانشگاه سمنان، سمنان، ایران (نویسنده مسئول): mahdi.roozbeh@semnan.ac.ir

^۲ دانش‌آموخته کارشناسی ارشد آمار، دانشگاه سمنان، سمنان، ایران.

^۳ دانش‌آموخته کارشناسی ارشد آمار، دانشگاه سمنان، سمنان، ایران.

۲.۱ روش رگرسیون لاسو

رویکرد رگرسیون لاسو^۵ نخستین بار توسط لئو برایمن^۶ با نام یقه فلزی (برای خفه کردن محکوم)^۷ بیان شد [۱۳]. نخستین بار مسئله بهینه‌سازی نامنفی او به صورت زیر مطرح شد:

$$\sum_{i=1}^n \left(Y_i - \sum_{j=1}^p c_j \beta_j X_{ji} \right)^2 \quad c_j \geq 0, \quad \sum_j c_j \leq s \quad (1)$$

که در آن برآوردگر اولیه $\hat{\beta}_j$ با روش کمترین توان‌های دوم و s در آن پارامتر جریمه است که با کاهش آن گروتی محدود می‌شود. این روش در بین محققان به نام جریمه مرگ^۸ نیز معروف است. در این روش برخی از متغیرها حذف شده و بقیه آن‌ها منقبض می‌شود. از مزایای این روش این است که خطای آن در مقایسه با خطای رگرسیون بهترین مجموعه کمتر است و فقط هنگامی که ضرایب کوچک غیر صفر دارد، باید در استفاده از این روش احتیاط لازم انجام شود. با توجه به این مسئله، لاسو تابع هدفی ارائه می‌دهد که در آن از برآوردگرهای کمترین توان‌های دوم استفاده نمی‌کند که به صورت زیر نوشته می‌شود:

$$\min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n \left(Y_i - \sum_{j=1}^p X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\},$$

$$i = 1, \dots, n, \quad j = 1, \dots, p$$

از مزایای اثبات شده این روش این است که برآوردگرهای پایدار و پیوسته تولید می‌کند [۱۳]. همچنین این روش برخی از ضرایب را منقبض و بقیه را صفر می‌کند. از معایب روش برآورد لاسو می‌توان به اریبی و عملکرد ضعیف آن در حضور متغیرهای توضیحی همبسته اشاره کرد زیرا که بین متغیرهای همبسته به ویژه متغیرهایی که به صورت گروهی همبسته اند یکی از آن‌ها را انتخاب می‌کند و متغیر انتخاب شده لزومی ندارد که بهترین متغیر باشد. لاسو برای یک مدل رگرسیونی خطی با p متغیر توضیحی و n مشاهده، حداکثر n متغیر را انتخاب می‌کند، بنابراین هنگامی که متغیرهای بیشتری (بیشتر از n) در مدل معنی‌دار باشند، در این روش شانسی برای انتخاب ندارند [۲].

تعیین تعداد مؤلفه‌های اصلی انتخاب شماری از مؤلفه‌هاست که بتوانند درصد قابل توجهی از واریانس (تغییرات کل) را توجیه کنند. برای مثال اگر ۱۰۰ متغیر توضیحی داشته باشیم و با اتخاذ ۲۰ مؤلفه ۷۰ درصد تغییرات واریانس کل و با اتخاذ ۴۰ مؤلفه ۷۵ درصد تغییرات واریانس کل توجیه شود، اضافه کردن ۲۰ مؤلفه دیگر برای دست یافتن به ۵ درصد تغییرات بیشتر به صرفه نخواهد بود. روش دیگر، برگزیدن تعداد مؤلفه‌هایی است که واریانس آن‌ها بزرگ‌تر یا مساوی متوسط واریانس کل است. توجه شود که اگر از ماتریس همبستگی استفاده می‌کنیم مؤلفه‌هایی که واریانس آن‌ها بزرگ‌تر یا مساوی ۱ است، برگزیده می‌شوند. از روش‌های شهودی نیز برای این امر استفاده از نمودار بازو^۴ می‌باشد. در این نمودار مقادیر ویژه هر مؤلفه (λ_i) در برابر i رسم می‌شود. با عمود کردن ناحیه‌ای از نمودار که شیب آن به طور ناگهانی کم می‌شود، تعداد مؤلفه‌ها تعیین می‌گردد. البته می‌توان لگاریتم (λ_i) در برابر i رسم کرد. استفاده از لگاریتم به این علت است که تابع لگاریتم مقیاس را کاهش می‌دهد. البته استفاده از این روش یا روش نمودار بازو تفاوت چندانی ندارد، گرچه نمودار بازو جهش‌های ناگهانی را بهتر نشان می‌دهد. توجه به این نکته ضروری است که ممکن است هر کدام از روش‌های مطرح شده پاسخ‌های متفاوتی ارائه دهند که محقق می‌تواند بنا بر کاربرد خود روش مورد نظر خود را به کارگیرد. در استفاده از این روش باید به پایا بودن مقیاس متغیرهای توضیحی اولیه توجه نمود چراکه مؤلفه‌ی اصلی در مواجهه با متغیر دارای واریانس بالا بسیار وابسته آن می‌شود که باعث خطای زیاد در کار می‌شود و در این حالت استفاده از ماتریس همبستگی به علت پایایی می‌تواند جلوی این مشکل را بگیرد. بعد از انتخاب تعداد متغیرها نوبت به برازش مدل رگرسیون مؤلفه‌های اصلی می‌شود. مدل رگرسیونی دیگر با مشکل هم خطی روبرو نیست و می‌توان از روش کمترین توان‌های دوم معمولی برای برآورد ضرایب آن استفاده کرد. از کاربردهای روش مؤلفه‌های اصلی می‌توان به تبدیل متغیرهای همبسته به متغیرهای ناهمبسته، یافتن ترکیبات خطی با واریانس نسبی بزرگ یا کوچک و کاهش در حجم داده‌ها نام برد.

⁴Scree Diagram

⁵Lasso

⁶Leo Breiman

⁷Garrote

⁸Death Penalty

⁹Andrey Nikolayevich Tikhonov

۲ رگرسیون ریب

پس عناصر قطری آن بزرگتر یا مساوی صفر هستند. حال در عبارت (۲) با جایگزینی ماتریس زیر بجای ماتریس $X^T X$ می توان نوشت:

$$X^T X + kI = VDV^T + V kIV^T = V(D + kI)V^T$$

در عبارت فوق $VV^T = I$ است. پس ماتریس مربوطه وارون پذیر است. برای k های مختلف داریم:

۱- اگر $k = 0$ باشد، برآوردگر ریب همان برآوردگر کمترین توان های دوم معمولی می شود.

۲- اگر $k = \infty$ آنگاه $\beta = 0$ است.

۳- اگر k بین صفر و بی نهایت باشد، مسئله در حال تعادل بین دو روش برازش کردن یک مدل خطی روی متغیرها و کوچک کردن ضرایب است.

۳ برآورد مدل تابعی

در ابتدا به برآورد یک منحنی یا خط ساده برای برازش روی داده ها نیاز است و چون داده ها به صورت گسسته می باشند، باید داده های گسسته را به داده های پیوسته تبدیل کرد. اگر داده ها بدون خطا باشند، از روش درون یابی استفاده می شود و اگر دارای خطا باشند، از روش هموارسازی می توان بهره گرفت. چون اکثراً داده ها دارای خطا هستند، پس باید از روش هموارسازی استفاده کرد.

۴ تبدیل داده های گسسته به پیوسته

با توجه به فرم قرارگیری داده ها، رگرسیون خطی به صورت زیر تعریف می شود:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \epsilon_i \quad (3)$$

و اگر به فرم منحنی باشند، به صورت زیر نوشته می شود:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3 + \dots + \epsilon_i \quad (4)$$

به طور کلی اگر متغیر وابسته به صورت یک ترکیب خطی از توابع پایه و ضرایب مربوطه باشد، آنکه مدل رگرسیونی به فرم زیر قابل نمایش است:

$$Y_i = \sum_{j=1}^k c_j \phi_j(t_i) + \epsilon_i = f(t_i) + \epsilon_i \quad (5)$$

به طوری که در معادله بالا $\phi_j(t)$ تابع پایه و وابسته به نوع داده ها می باشد و c_j ها ضرایب می باشند. در ادامه مثال هایی از توابع پایه آورده خواهند شد.

محققان در مسائل رگرسیونی باهم خطی متغیرها روبرو می شوند که معمولاً در مدل هایی با تعداد پارامتر بالا اتفاق می افتد. آندری نیکولایویچ تیخونوف^۹ که تحقیقات مهم او در توپولوژی، تجزیه و تحلیل عملکردی، فیزیک، ریاضیات شناخته شده است، منظم سازی تیخونوف را به عنوان راه حلی برای رویارویی با مسائل بدشرطی معرفی نمود. رگرسیون خطی می تواند نااریب باشد اما واریانس بالایی داشته باشد که ممکن است بعضی از اربیبی ها را از دست بدهد تا واریانس کمتری داشته باشد. در داده های با بعد بالا چون ماتریس طرح وارون پذیر نیست، روش کمترین توان های دوم رگرسیون خطی نمی تواند مفید باشد. رگرسیون خطی زمانی که تعداد متغیرها بیشتر از داده ها باشد تعریف نمی شود. در این شرایط با استفاده از روش ریب می توان معکوس ماتریس را به دست آورد، بدین صورت که با افزودن مقدار λ به درایه های قطر اصلی ماتریس $X^T X$ این ماتریس معکوس پذیر خواهد شد [۱۴]. اگر نوك كوه نقطه صفر در نظر گرفته شود، با افزایش انقباض متغیرها، ضرایب کوچک و کوچک تر شده که همانند کاهش عرض کوه و افزایش ارتفاع آن می باشد و ضرایب هیچ گاه دقیقاً مقدار صفر را نمی پذیرند. برآوردگر ریب دیگر نااریب نیست اما دارای واریانس کمتری می باشد و مقدار برآوردگر پایا است یعنی تغییرات جزئی در داده ها، بر برآوردگر تأثیری نخواهد داشت [۳]. برآورد β به روش ریب عبارت است از:

$$\hat{\beta} = (X^T X + kI)^{-1} X^T Y, \quad k \geq 0 \quad (2)$$

که به k پارامتر ریب می گوئیم و یافتن این پارامتر بسیار مهم می باشد. برای انتخاب k باید به گونه ای عمل کنیم که کاهش در واریانس برآوردگر اربیب بیش از افزایش مربع اربیبی باشد [۱۰]، در نتیجه MSE آن کمتر از واریانس برآوردگر نااریب خواهد بود.

طبق قضیه تجزیه مقادیر منفرد ماتریس ها که بر اساس آن یک ماتریس به دو ماتریس متعامد و یک ماتریس قطری تجزیه می شود، می توان گفت روش ریب همواره منجر به یافتن برآورد مدل رگرسیونی خواهد شد، به بیان دیگر می توان نوشت:

$$X^T X = VDV^T$$

که در آن D ماتریس قطری با عناصر قطری مقادیر ویژه ی ماتریس $X^T X$ بوده و V و V^T ماتریس های متعامد هستند، ماتریس $X^T X$ نیمه معین مثبت است زیرا:

$$\forall Z \neq 0, Z^T X^T X Z = (XZ)^T (XZ) = \|XZ\| \geq 0$$

۱۰۴ مثال‌هایی از توابع پایه

مرسوم‌ترین توابع پایه عبارت‌اند از:

$$1. \text{ ثابت }^{10}: \phi(t) = 1$$

$$2. \text{ توانی }^{11}: \phi(t) = t^{\lambda_1}, t^{\lambda_2}, t^{\lambda_3}, \dots$$

$$3. \text{ نمایی }^{12}: \phi(t) = e^{\lambda_1 t}, e^{\lambda_2 t}, e^{\lambda_3 t}, \dots$$

۴۰۴ برآورد ضرایب تابعی

تابع رگرسیونی به صورت

$$Y_i = f(t_i) + \epsilon_i$$

را در نظر می‌گیریم که در آن خطاها مستقل و دارای توزیع نرمال با میانگین صفر و واریانس σ^2 و $i = 1, \dots, n$ می‌باشند. برای برآورد $f(t_i)$ بر اساس توابع پایه داریم:

$$\hat{f}(t_i) = \sum_{j=1}^p c_j \phi_j(t_i)$$

برای برآورد ضرایب تابعی باید مجموع مربع خطاها را به صورت زیر حداقل کرد:

$$H(c) = \sum_{i=1}^n (Y_i - f(t_i))^2 = \sum_{i=1}^n (Y_i - \sum_{j=1}^p c_j \phi_j(t_i))^2 \quad (6)$$

همچنین می‌توان معادله‌ی فوق را به صورت زیر نوشت:

$$H(c) = (Y - \Phi c)^T (Y - \Phi c) \quad (7)$$

اکنون با مشتق‌گیری و برابر صفر قرار دادن آن داریم:

$$\hat{c} = (\Phi^T \Phi)^{-1} \Phi^T Y \quad (8)$$

حال با توجه به برآورد بالا مقدار $\hat{f}(t)$ برابر است با:

$$\hat{f}(t) = \hat{c}^T \Phi$$

و همچنین \hat{Y} به صورت:

$$\hat{Y} = \underbrace{\Phi(\Phi^T \Phi)^{-1} \Phi^T}_S Y = SY$$

به دست می‌آید که به S ماتریس هموارساز می‌گویند. در عبارت فوق مقدار \hat{Y} ضریبی از مقدار Y است، یعنی با هموارکردن داده‌ها توسط S مقدار \hat{Y} به دست می‌آید. انتخاب تعداد توابع پایه بسیار اهمیت دارد، اگر تعداد این توابع کم در نظر گرفته شود منجر به اریبی زیاد و واریانس کم شده و زیاد در نظر گرفتن تعداد این توابع منجر به اریبی کم و واریانس زیاد می‌شود.

۲۰۴ تقریب توابع پایه با استفاده از پایه‌های فوریه

تابع‌های تقریب‌زن پایه‌های فوریه اکثراً برای داده‌هایی کاربرد دارد که نوسانی بوده و دارای دوره تناوب باشند، همانند داده‌های مربوط به آب‌وهوا. پایه‌های فوریه به صورت زیر تعریف می‌شود [۱۵]:

$$\{1, \sin(\omega t), \cos(\omega t), \sin(2\omega t), \cos(2\omega t), \dots, \sin(m\omega t), \cos(m\omega t)\}$$

در تابع پایه بالا ω دوره‌ی نوسان نامیده شده و برابر است با:

$$\omega = \frac{2\pi}{p}$$

که p دوره تکرارشونده است. برای مثال دوره تکرارشونده برای داده‌های آب‌وهوا برابر با $p = 365$ روز است.

۳۰۴ تقریب توابع پایه با استفاده از پایه‌های اسپلاین

پایه‌های اسپلاین^{۱۳} یا بی اسپلاین^{۱۴} هر دو نقطه‌ی مجاور را با یک تابع برآورد می‌کند. اسپلاین‌ها همانند تابع‌های چندجمله‌ای هستند که ابتدا داده‌های گسسته را به چند قسمت مساوی تقسیم کرده و سپس به دنبال بهترین منحنی برای برازش به هر قسمت می‌باشند که اگر درجه آن صفر باشد با یک خط عمودی و افقی به برآورد می‌پردازد و اگر درجه‌ی آن یک باشد به صورت خطی و درجه‌های بالاتر را به صورت منحنی برآورد می‌کند. در ضمن قسمت‌هایی از منحنی که در محل پیوستن به هم هستند را می‌توان هموار کرد، نقاطی را که در قسمت اتصال قرار می‌گیرند گره نامیده می‌شود.

برای انتخاب گره‌ها، اگر تعداد آن‌ها خیلی زیاد باشد اریبی بسیار کم و واریانس بسیار زیاد می‌شود و باعث ناهمواری نمودار می‌شود [۱۵].

¹⁰Constant

¹¹Power

¹²Exponential

¹³Spline

¹⁴B-Spline

در معادله‌ی فوق R را به‌عنوان ماتریس جریمه می‌شناسیم. با حداقل کردن مربع خطاها و اندازه طول مشتق دوم، مقدار برآورد c به‌صورت زیر خواهد بود:

$$\hat{c} = [\Phi^T \Phi + \lambda R]^{-1} \Phi^T Y$$

و مقدار برآورد Y برابر است با:

$$\hat{Y} = \Phi \underbrace{[\Phi^T \Phi + \lambda R]^{-1} \Phi^T}_{S} Y = SY$$

برای انتخاب پارامتر هموارساز می‌توان از مینیم سازی معیارهای متداول زیر که برآوردهای مجانباً ناریب از تابع ریسک با زیان توان دوم هستند، استفاده کرد [۱۶]:

۱. اعتبار سنجی متقابل معمولی^{۱۵}:

$$OCV(\lambda) = \frac{1}{n} \sum_{i=1}^n \frac{(Y_i - \hat{Y}_i)^2}{(1 - S_{ii})^2} \quad (10)$$

که در آن S_{ii} مؤلفه i ام روی قطر اصلی ماتریس S است.

۲. اعتبار سنجی متقابل تعمیم‌یافته^{۱۶}:

$$GCV(\lambda) = \frac{1}{n} \sum_{i=1}^n \frac{(Y_i - \hat{Y}_i)^2}{[n^{-1} \text{trace}(I_n - S)]^2}$$

که در آن I_n ماتریس همانی از مرتبه n است. برای دیدن جزئیات بیشتر این معیار می‌توان به [۹، ۴، ۱] مراجعه کرد.

۶ تحلیل داده‌های واقعی

۱۰۶ داده‌های طیف‌سنجی

در این قسمت کلیه تحلیل‌ها در نرم‌افزار R انجام شده است و برای انجام آن نیاز به فراخوانی کتابخانه‌های fda ، $fda.usc$ ، $goffda$ و $refund$ می‌باشد.

در این بخش به بررسی داده‌هایی با متغیر مستقل تابعی و متغیر پاسخ عددی با بعد پایین و همچنین انجام رگرسیون مؤلفه‌های اصلی با جریمه‌های متفاوت و مقایسه آن‌ها پرداخته می‌شود. در ابتدا داده‌های $tecator$ از کتابخانه $fda.usc$ را فراخوانی می‌کنیم. این داده‌ها شامل ۱۰۰ طیف جذب مادون‌قرمز نزدیک (NIT)^{۱۷} است که برای پیش‌بینی مقدار رطوبت، چربی و پروتئین برای تیکه‌های گوشت خرد شده استفاده

۵ جریمه کمینه کردن طول مشتق دوم

با داشتن متغیر تابعی $X(t)$ و با تعریف مشتق‌های توابع می‌توان منحنی‌هایی هموار داشت. در ابتدا داریم:

$$DX(t) = \frac{d}{dt} X(t)$$

و می‌توان مشتقات مراتب بالاتر را به‌صورت زیر تعریف کرد:

$$D^k X(t) = \frac{d^k}{dt^k} X(t), \dots, D^k X(t) = \frac{d^k}{dt^k} X(t)$$

$D^k X(t)$ را می‌توان طبق تعریف مشتق دوم به‌صورت زیر نوشت:

$$D^2 X(t_i) = \frac{Y_{i+1} - Y_{i-1} - 2Y_i}{(\Delta t)^2}$$

اکنون طول منحنی را با توجه به فرمول زیر می‌توان نوشت:

$$j_2(X(t)) = \int (D^2 X(t))^2 dt \quad (9)$$

اگر طول منحنی کم باشد بدین معنی است که به نوسان‌های داده‌ها کمتر توجه شده و اگر اندازه منحنی زیاد باشد، پس باید جریمه افزایش یابد، به بیانی دیگر باید مشتق دوم ضرایب تابعی را هموار کرد تا منحنی این ضرایب هموار شود. با در نظر گرفتن مربع خطای جریمه‌شده به فرم زیر داریم:

$$PENSSSE_\lambda(X(t)) = (Y - X(t))(Y - X(t))^T + \lambda j_2(X(t))$$

به λ پارامتر هموارساز گفته می‌شود که کنترل بین مقادیر برازش شده و همواری منحنی بر داده‌ها می‌باشد. اگر مقدار λ افزایش یابد، هموارکننده بسیار جریمه می‌شود و $X(t)$ تبدیل به خط می‌شود و اگر λ کاهش یابد، جریمه هم کاهش می‌یابد و اجازه می‌دهد $X(t)$ روی داده‌ها برازش شود. باید مربع خطای جریمه‌شده را حداقل کرد و این امر در پایه‌های بی‌اسپلین زمانی می‌افتد که مقدار درجه آن برابر ۴ باشد، زیرا دارای مشتق دوم بوده و می‌توان مشتق دوم را هموار کرد و منحنی‌های هموار مناسب‌تری ایجاد کرد. هنگامی که درجه بی‌اسپلین برابر ۴ باشد به آن اسپلین هموار مکعب گفته می‌شود. اکنون با تعریف $X(t) = \Phi^T c$ و با توجه به اینکه $X(t)$ یک ترکیب خطی می‌باشد داریم:

$$\int [D^m X(t)]^2 dt = c^T \underbrace{\int D^m \Phi D^m \Phi^T dt}_R c = c^T R c$$

¹⁵Ordinary Cross Validation

¹⁶Generalized Cross Validation

¹⁷Near Infrared Transmission

سمت راست برآورد ضرایب تابعی برای مدل رگرسیون مؤلفه اصلی با جریمه لاسو می‌باشد.

اکنون داده‌های آزمون را در مدل‌ها بررسی و سپس به انتخاب یکی از سه مدل یعنی مدل بدون جریمه، مدل با جریمه مشتق دوم و مدل با جریمه ریج می‌پردازیم. می‌توان در شکل ۵ نمودارهای رسم شده مقادیر برازش شده در مقابل مقادیر واقعی برای داده‌های آزمون را مشاهده کرد. نتایج به‌دست‌آمده در جدول ۱ گویای آن است که با اختلاف خیلی کم مدل رگرسیون مؤلفه اصلی با جریمه لاسو مناسب‌تر نسبت به سایر جریمه‌ها عمل کرده است.

۲.۶ داده‌های آب‌وهوای کانادا

در این نوع داده دمای هوا و میزان بارندگی ۳۵ شهر کانادا در ۳۶۵ روز از سال مشخص است. حال باید به دنبال رگرسیون مؤلفه اصلی بین دمای این ۳۵ شهر با میزان بارندگی آن‌ها بود. با فراخوانی این داده‌ها از کتابخانه FDA می‌توان مدل تابعی زیر را در نظر گرفت:

$$Y_i = \beta_0(t) + \sum_{i=1}^{25} \beta_i(t)X_i(t) \quad (13)$$

به‌طوری‌که در مدل بالا Y_i میزان بارندگی برای ۳۵ ایستگاه آب‌وهوایی می‌باشد و $X_i(t)$ میزان دمای ۳۵ ایستگاه در طول ۳۶۵ روز می‌باشد و $\beta(t)$ ضرایب تابعی مدل می‌باشند. با استفاده از توابع پایه فوریه متغیرهای پیش‌بین را می‌توان تبدیل به منحنی کرد:

$$X_i(t) = \sum_{i=1}^n c_i \phi_i(t)$$

در ابتدا به تعریف ۳۵ تابع پایه برای تبدیل داده‌های گسسته دمای هوا به پیوسته توسط توابع پایه فوریه پرداخته می‌شود و دلیل استفاده از توابع فوریه تناوبی بودن داده‌ها می‌باشد که منحنی آن را می‌توان در نمودار ۶ مشاهده کرد. اکنون حالت‌های مختلف رگرسیون مؤلفه اصلی بدون جریمه، با جریمه مشتق دوم و جریمه لاسو را بررسی می‌شود. مدل رگرسیون مؤلفه اصلی بر اساس سه مؤلفه اول بدون در نظر گرفتن جریمه به‌صورت زیر می‌باشد:

$$Y_i = \alpha_1(t)PC_1(t) + \alpha_2(t)PC_2(t) + \alpha_3(t)PC_3(t) \quad (14)$$

در معادله فوق Y_i میزان بارندگی برای ۳۵ ایستگاه آب‌وهوایی می‌باشد. پس از برازش، مقدار همبستگی بین مقادیر واقعی و مقادیر برازش شده عدد ۰.۳۶ می‌باشد که حاکی از آن است که روش بدون جریمه مناسب عمل نکرده است. اکنون به دنبال برآوردی مناسب‌تر با استفاده از اعتبارسنجی متقابل برای انتخاب λ و انتخاب تعداد مؤلفه‌های

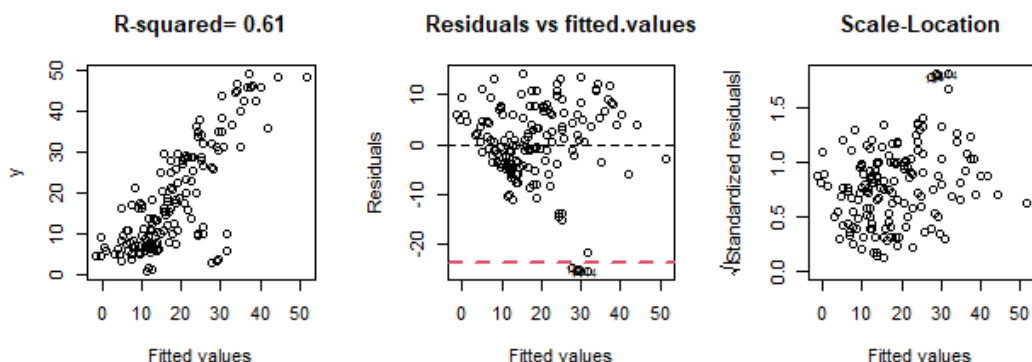
می‌شود. این داده‌ها بر اساس آنالیز غذا و خوراک Tecator Infratec می‌باشد که در محدوده طول موج ۱۰۵۰-۸۵۰ نانومتر کار می‌کند و توسط اصل انتقال مادون‌قرمز نزدیک (NIT) ثبت می‌شود. هر نمونه حاوی گوشت خالص ریز خردشده با رطوبت، چربی و پروتئین متفاوت است. هدف پیش‌بینی محتوای چربی می‌باشد. این داده‌ها برای هر نمونه گوشت از ۱۰۰ کانال طیف تشکیل شده است که به اندازه‌گیری جذب و محتویات رطوبت، چربی و پروتئین خواهیم پرداخت. مدل تابعی برای داده‌های بالا به فرم زیر نوشته می‌شود:

$$Fat = \beta_0(t) + \sum_{i=1}^{100} \beta_i(t)X_i(t) \quad (11)$$

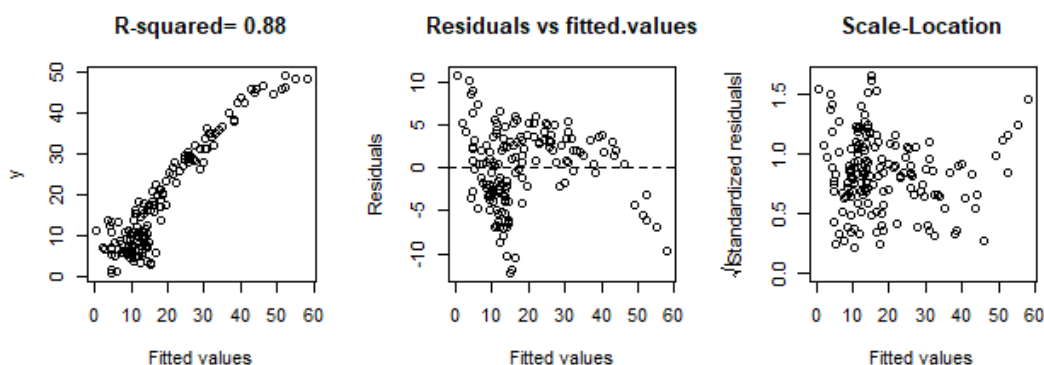
در مدل فوق Fat میزان چربی موجود در هر قسمت گوشت خالص و $X_i(t)$ اطلاعات ثبت‌شده توسط امواج مادون‌قرمز با طول موج ۸۵۰-۱۰۵۰ می‌باشد و t برابر ۱۰۰ طیف کانال می‌باشد. با تقسیم داده‌ها و اختصاص دادن هشتاد درصد آن‌ها به داده‌های آموزش و بیست درصد داده‌ها به داده‌های آزمون به بررسی مدل روی داده‌های آموزش می‌پردازیم. اکنون با توجه به عددی بودن متغیر پاسخ و تابعی بودن متغیر مستقل می‌توان با استفاده از رگرسیون با سه مؤلفه اصلی مدل زیر را در نظر گرفت:

$$Fat = \alpha_1(t)PC_1(t) + \alpha_2(t)PC_2(t) + \alpha_3(t)PC_3(t) \quad (12)$$

به‌طوری‌که در مدل فوق PC_i ، $i = 1, 2, 3$ مؤلفه‌های اصلی برای داده‌های مذکور می‌باشند. بدون هیچ‌گونه جریمه و بر اساس رگرسیون مؤلفه اصلی، مقدار مجذور همبستگی بین مقادیر برازش شده و مقادیر واقعی ۰.۶۱ می‌باشد باقیمانده و باقیمانده استاندارد با مقادیر برازش شده را می‌توان در شکل ۱ مشاهده کرد. با استفاده از اعتبارسنجی متقابل تعداد ۳ مؤلفه اصلی نقش اساسی دارند. حال جریمه ریج و جریمه لاسو و جریمه مشتق دوم را در رگرسیون مؤلفه اصلی با انتخاب متغیرهای اول و سوم و چهارم و مقدار λ بهینه که توسط اعتبارسنجی متقابل به‌دست‌آمده است، با استفاده از جریمه مشتق دوم مقدار همبستگی بین مقادیر برازش شده و مقادیر واقعی عدد ۰.۸۸ می‌باشد و نمایش نمودار باقیمانده و باقیمانده استاندارد با مقادیر برازش شده را می‌توان در شکل‌های ۲ و ۳ برای مدل‌های جریمه‌ای پیشنهادی مشاهده کرد. اکنون به بررسی برآورد ضرایب تابعی در مدل‌های فوق پرداخته می‌شود. با توجه به شکل ۴ نمودار بالا سمت چپ برآورد ضرایب تابعی مدل بدون جریمه می‌باشد، نمودار بالا سمت راست، نمودار برآورد ضرایب تابعی با جریمه ریج که با دخالت ۳ متغیر اصلی و نمودار پایین سمت چپ برآورد ضرایب تابعی با جریمه مشتق دوم، نمودار پایین



شکل ۱: مدل رگرسیون مؤلفه اصلی بدون جریمه، نمودار سمت چپ: پراکنش مقادیر واقعی در مقابل مقادیر برازش شده، نمودار وسط: باقیمانده‌ها در مقابل مقادیر برازش شده و نمودار سمت راست: ریشه قدرمطلق باقیمانده‌ها در مقابل مقادیر برازش شده.



شکل ۲: مدل رگرسیون مؤلفه اصلی با جریمه مشتق دوم، نمودار سمت چپ: پراکنش مقادیر واقعی در مقابل مقادیر برازش شده، نمودار وسط: باقیمانده‌ها در مقابل مقادیر برازش شده و نمودار سمت راست: ریشه قدرمطلق باقیمانده‌ها در مقابل مقادیر برازش شده در مدل جریمه‌شده.

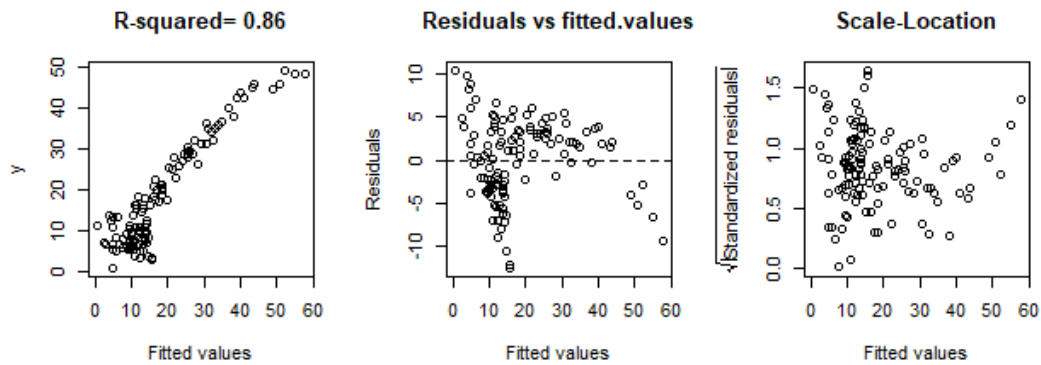
واقعی و مقادیر برازش شده را در تصویر ۷ مشاهده کرد. نمودار برآورد ضرایب تابعی را برای دو رگرسیون مؤلفه اصلی با جریمه مشتق دوم و جریمه لاسو می‌توان در تصویر ۸ مشاهده کرد. ضرایب مؤلفه‌های اصلی برای سه مدل فوق به صورت مشخص شده در جدول ۳ می‌باشد. با توجه به آزمون مجذور همبستگی به دست آمده در جدول ۲ مشخص است که مدل رگرسیون با جریمه مشتق دوم عملکرد مناسبی داشته است.

مدل‌سازی پرداخته شد که با توجه به نتایج به دست آمده رگرسیون مؤلفه اصلی با جریمه لاسو و ریب با میزان همبستگی بین مقادیر مشاهده شده و مقادیر برازش شده 0.9491003 مدلی مناسب‌تر بود. در داده‌های آب‌وهوای کانادا دو جریمه ریب و مشتق دوم انجام شده است و نتایج با

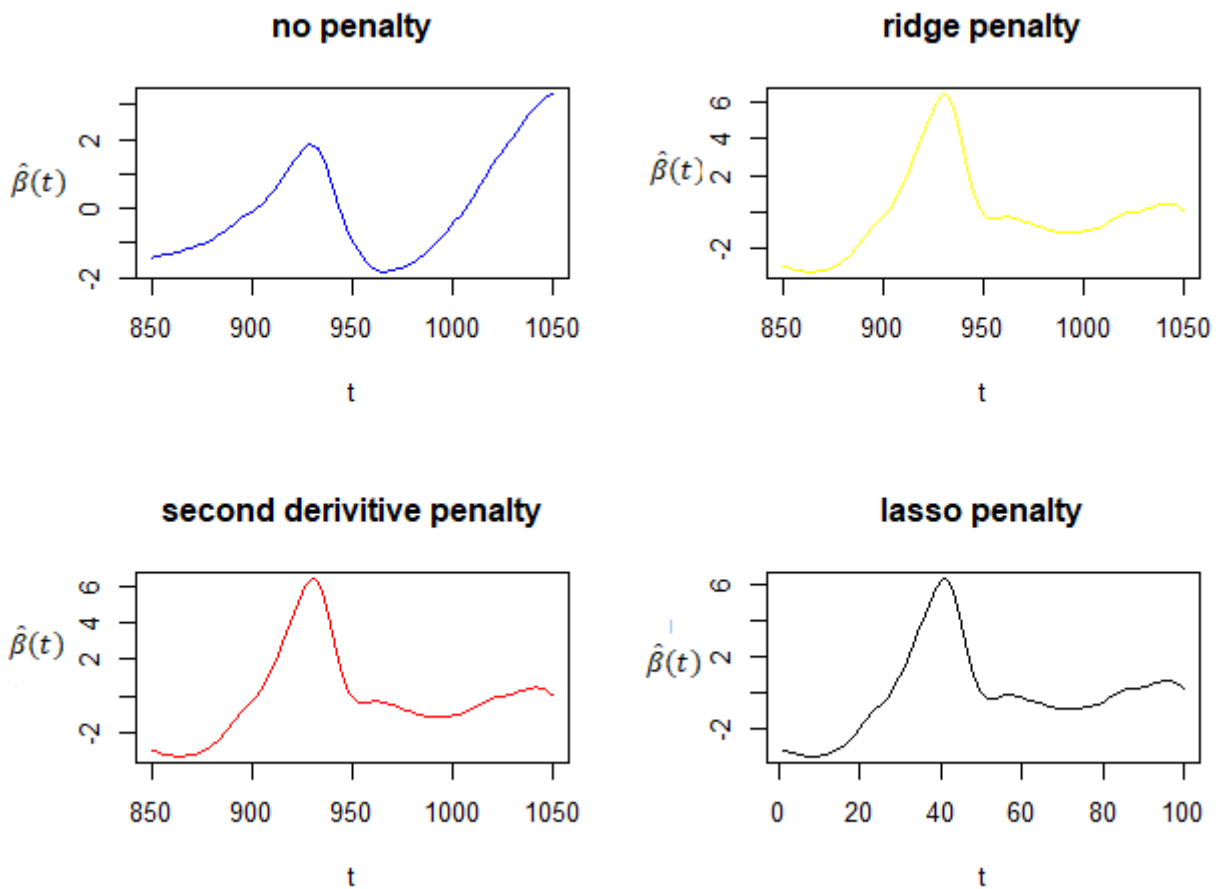
اصلی و جریمه کردن ضرایب تابعی می‌باشیم. با توجه به جریمه مشتق دوم صورت گرفته و عدد همبستگی 0.86 مشخص است که مدل فوق مناسب‌تر شده است. اکنون با جریمه لاسو هم این مدل‌سازی انجام شد که مقدار λ را با استفاده از اعتبارسنجی متقابل به دست آورده و با توجه به عدد همبستگی بین مقادیر برازش شده و مقادیر واقعی می‌توان گفت این روش نسبتاً خوب بوده و می‌توان نمودار رسم شده بین داده‌های

۷ بحث و نتیجه‌گیری

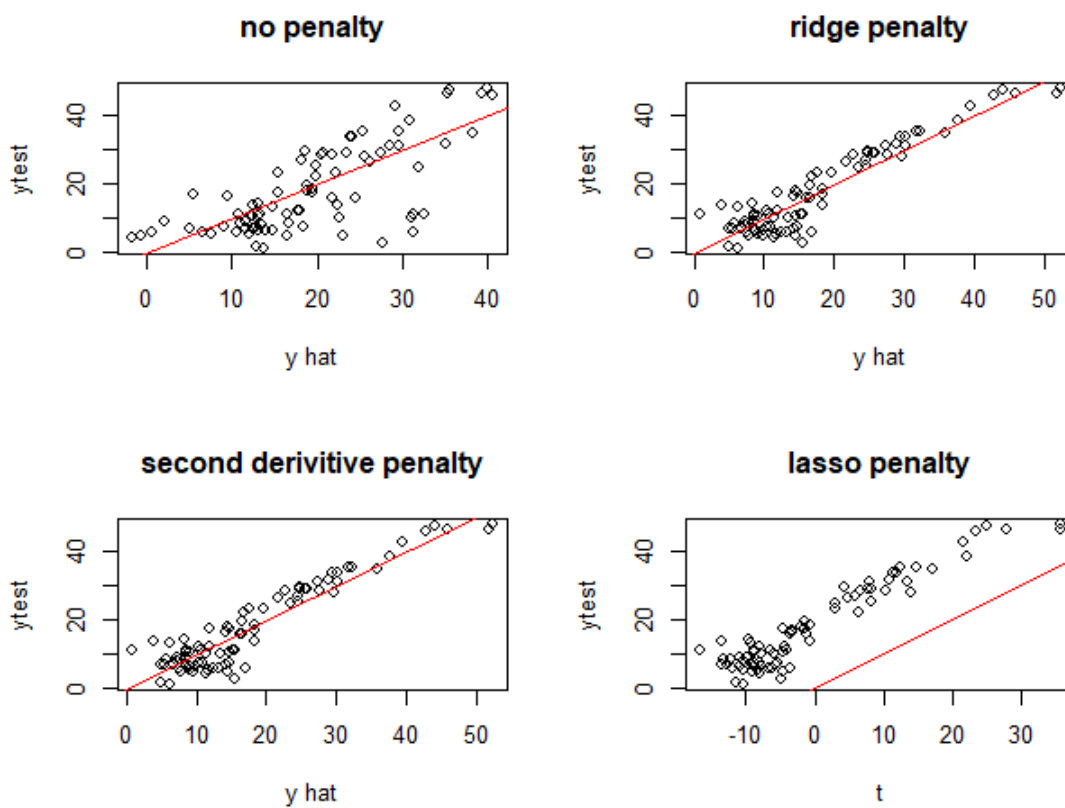
با بررسی داده‌های tecator با استفاده از رگرسیون مؤلفه اصلی و با جریمه‌های کمینه کردن طول منحنی مشتق دوم، جریمه لاسو و ریب به



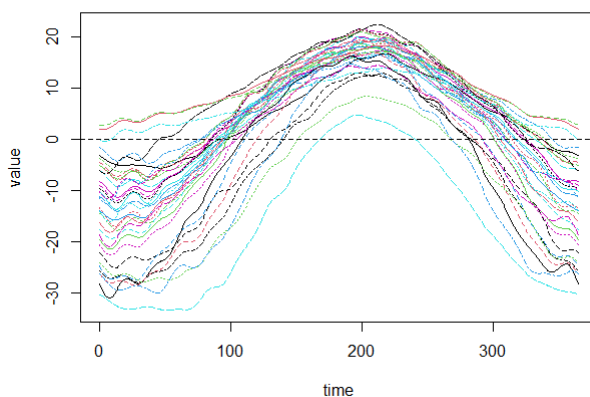
شکل ۳: مدل رگرسیون مؤلفه اصلی با جریمه ریج، نمودار سمت چپ: پراکنش مقادیر واقعی در مقابل مقادیر برازش شده، نمودار وسط: باقیمانده‌ها در مقابل مقادیر برازش شده و نمودار سمت راست: ریشه قدرمطلق باقیمانده‌ها در مقابل مقادیر برازش شده در مدل جریمه‌شده.



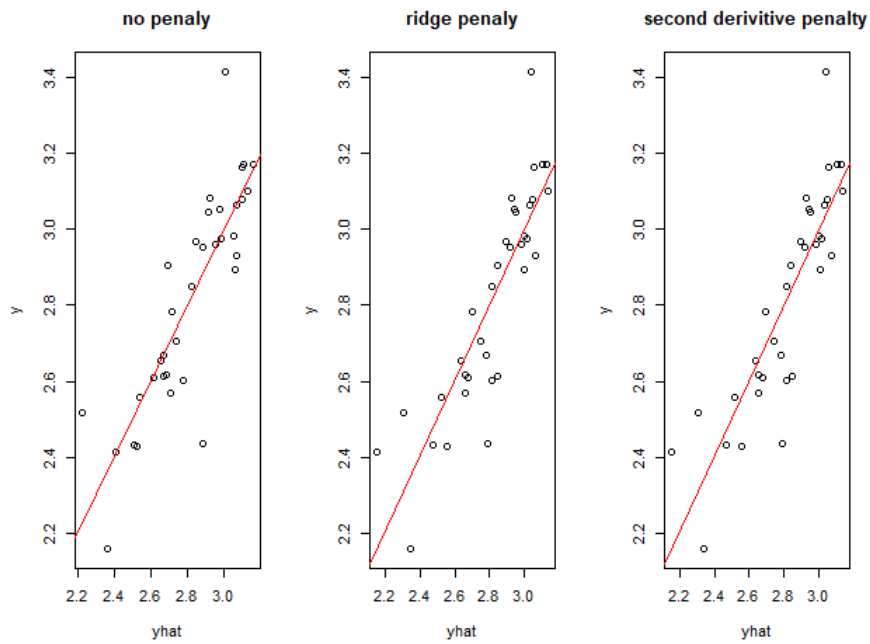
شکل ۴: نمودار برآورد ضرایب تابعی با جریمه‌های متفاوت



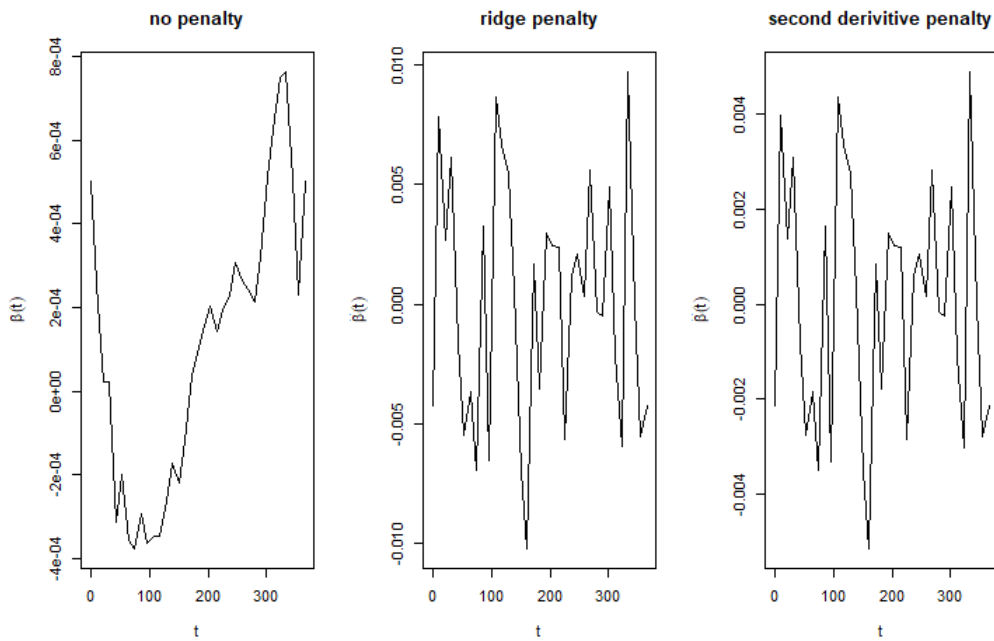
شکل ۵: نمودار مقادیر برازش شده در مقابل مقادیر واقعی برای داده‌های آزمون که نمودار بالا سمت چپ بدون جریمه، نمودار بالا سمت راست با جریمه ریج، نمودار پایین سمت چپ با جریمه مشتق دوم و نمودار پایین سمت راست با جریمه لاسو می‌باشد.



شکل ۶: نمایش منحنی دمای هوای ۳۵ ایستگاه آب‌وهوای کانادا



شکل ۷: مقادیر برازش شده و مقادیر واقعی در مدل بدون جریمه و با رگرسیون مؤلفه اصلی با جریمه لاسو و با جریمه مشتق دوم



شکل ۸: برآورد ضرایب تابعی در سه مدل متفاوت

جدول ۱: آزمون مجذور همبستگی بین مقادیر واقعی و مقادیر برازش شده در داده‌های طیف‌سنجی

مقدور همبستگی	نوع مدل
۰/۷۳۸۴۵۳۷	مدل رگرسیون مؤلفه اصلی بدون جریمه
۰/۹۳۹۱۷۱۳	مدل رگرسیون مؤلفه اصلی با جریمه ریج
۰/۹۳۹۰۹۷۱	مدل رگرسیون مؤلفه اصلی با جریمه مشتق دوم
۰/۹۴۹۱۰۰۳	مدل رگرسیون مؤلفه اصلی با جریمه لاسو

جدول ۲: آزمون مجذور همبستگی بین مقادیر واقعی و مقادیر برازش شده در داده‌های آب‌وهوای کانادا

مقدور همبستگی	نوع مدل
۰/۸۵۳۱۸۲	مدل رگرسیون مؤلفه اصلی بدون جریمه
۰/۸۶۸۳۹۴۲	مدل رگرسیون مؤلفه اصلی با جریمه ریج
۰/۸۶۸۴۰۵۱	مدل رگرسیون مؤلفه اصلی با جریمه مشتق دوم

جدول ۳: در داده‌های آب‌وهوای کانادا ضرایب مؤلفه‌های اصلی

ضرایب	مؤلفه اصلی بدون جریمه	مؤلفه اصلی با جریمه مشتق دوم	مؤلفه اصلی با جریمه ریج
مؤلفه اول	-۰/۰۰۱۹	۰/۰۵۲۵	۰/۰۴۸۳۵
مؤلفه دوم	-۰/۰۰۱۸	-۰/۰۰۵۳	-۰/۰۰۴۹۳
مؤلفه سوم	-۰/۰۰۶۲	۰/۰۰۸۰	۰/۰۰۷۳۷
عرض از مبدأ	۲/۸۱۴۸	۲/۸۱۴۸	۲/۸۱۴۸

تقدیر و تشکر

توجه به میزان همبستگی مقادیر برازش شده مدل رگرسیون مؤلفه اصلی با جریمه مشتق دوم در مقابل مقادیر مشاهده‌شده رضایت‌بخش بودند.

نویسندگان مقاله کمال قدردانی و تشکر را از پیشنهادهای ارزنده داوران، سردبیر و ویراستار محترم مجله که باعث ارائه بهتر و افزایش سطح کیفی مقاله شده است، دارند.

مراجع

- [۱] روزبه، م. و امینی، م. (۱۳۹۸)، برآوردگر استوار مرزبندی شده تعمیم‌یافته محتمل در مدل رگرسیون نیمه‌پارامتری، *مجله علوم آماری ایران* ۱۳(۲)، ۳۴۱-۳۶۰.
- [۲] معنوی، م. و روزبه، م. (۱۳۹۹)، روش‌های تحلیل رگرسیونی برای داده‌های بعد بالا، *اندیشه آماری*، ۲۵(۱)، ۶۹-۹۰.
- [۳] روزبه، م.، ملک جعفریان، م. و معنوی، م. (۱۳۹۹)، کاربرد رگرسیون ستیغی کمترین توان‌های دوم پیراسته محدودشده تصادفی در مدل‌سازی مصرف آب، *اندیشه آماری*، ۲۶(۲)، ۱۹-۹.
- [4] Amini, M. and Roozbeh, M. (2015), Optimal Partial Ridge Estimation in Restricted Semiparametric Regression Models, *Journal of Multivariate Analysis*, **136**, 26-40.
- [5] Horváth, L., and Kokoszka, P. (2012). *Inference for functional data with applications*, Springer Science and Business Media, New York.
- [6] Hsing, T., and Eubank, R. (2015). *Theoretical foundations of functional data analysis, with an introduction to linear operators*, John Wiley and Sons, Hoboken.
- [7] Dette, H., Kokot, K., and Aue, A. (2017). Functional data analysis in the Banach space of continuous functions, *arXiv preprint arXiv*, 1710.07781.
- [8] Febrero Band, M., Galeano, P., and Gonzalez Manteiga, W. (2017). Functional principal component regression and functional partial least-squares regression: An overview and a comparative study, *International Statistical Review* , **85(1)** , 61-83.
- [9] Roozbeh, M. (2018), Optimal QR-Based Estimation in Partially Linear Regression Models with Correlated Errors Using GCV Criterion, *Computational Statistics & Data Analysis*, **117**, 45-61.
- [10] Roozbeh, M. and Arashi, M. (2016), New ridge regression estimator in semiparametric regression models, *Communications in Statistics-Simulation and Computation*, **45**, 3683-3715.
- [11] Reiss, P. T., Goldsmith, J., Shang, H. L. and Ogden, R. T. (2017). Methods for scalar-on-function regression, *International Statistical Review*, **85(2)** , 228-249.
- [12] Jolliffe, I.T. (2002). *Principal Component Analysis*, Springer series in statistics, Aberdeen.
- [13] Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society. Series B (Methodological)*, **58(1)**, 267-288.
- [14] Hoerl, A.E. and Kennard, R.W. (1975). Ridge regression some simulation, *Communication in Statistics*, **4**, 4105-4123.
- [15] Ramsay, J. O. and Silverman, B. W. (2005), *Functional Data Analysis*, Springer-Verlag, New York.
- [16] Wasserman, L. (2005), *All of Nonparametric Statistics*, Springer-Verlag, New York.

Modelling of functional data using principal component regression approach based on the generalized cross validation criterion

Mahdi Roozbeh¹ Arta Rouhi² Fatemeh Jahadi³

Abstract:

Functional data analysis is used to develop statistical approaches to the data sets that are functional and continuous essentially, and because these functions belong to the spaces with infinite dimensional, using conventional methods in classical statistics for analyzing such data sets is challenging. The most popular technique for statistical data analysis is the functional principal components approach, which is an important tool for dimensional reduction. In this research, using the method of functional principal component regression based on the second derivative penalty, ridge and lasso, the analysis of Canadian climate and spectrometric data sets is proceed. To do this, to obtain the optimum values of the penalized parameter in proposed methods, the generalized cross validation, which is a valid and efficient criterion, is applied.

Keywords: Functional Data Analysis, Functional Regression, Generalized Cross Validation, Principal Component Regression.

¹Faculty of mathematics, Semnan university, Semnan, Iran.

²Master's degree graduate, statistics and Computer science, Semnan university, Semnan, Iran.

³Master's degree graduate, statistics and Computer science, Semnan university, Semnan, Iran.