

الحمد لله  
الذي هدانا لهذا  
الذي كنا لنهتدي لولا  
أن هدانا الله





دانشکده علوم ریاضی

رشته آمار، گرایش آمار ریاضی

پایان نامه کارشناسی ارشد

# اسپلاین‌های جریمه‌ای و کاربردهای نوین آن

نگارنده: اکرم قائمی‌زاده

استاد راهنما

دکتر نگار اقبال

استاد مشاور

دکتر حسین باغیشنی

تیر ۱۳۹۷





شماره:

تاریخ:

باسمه تعالی



مدیریت تحصیلات تکمیلی

فرم شماره (۳) صورتجلسه نهایی دفاع از پایان نامه دوره کارشناسی ارشد

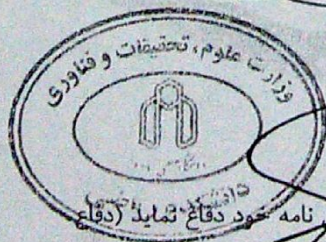
با نام و یاد خداوند متعال، ارزیابی جلسه دفاع از پایان نامه کارشناسی ارشد خانم اکرم قائمی زاده با شماره دانشجویی ۹۴۱۳۹۰۴ رشته آمار گرایش آمار ریاضی تحت عنوان اسپلین های جریمه ای و

کاربردهای نوین آن

که در تاریخ ۹۷/۰۴/۲۵ با حضور هیأت محترم داوران در دانشگاه صنعتی شاهرود برگزار گردید به شرح ذیل اعلام می گردد:

|                                |  |
|--------------------------------|--|
| <input type="checkbox"/> مردود | <input checked="" type="checkbox"/> قبول (با درجه: <u>خیلی خوب</u> ) |
| <input type="checkbox"/> عملی  | <input checked="" type="checkbox"/> نظری                             |

| عضو هیأت داوران           | نام و نام خانوادگی  | مرتبه علمی | امضاء |
|---------------------------|---------------------|------------|-------|
| ۱- استاد راهنمای اول      | دکتر نگار اقبال     | استادیار   |       |
| ۳- استاد مشاور            | دکتر حسین باغیشنی   | استادیار   |       |
| ۴- نماینده تحصیلات تکمیلی | دکتر عبدالله آل هوز | استادیار   |       |
| ۵- استاد ممتحن اول        | دکتر داود شاهسونی   | دانشیار    |       |
| ۶- استاد ممتحن دوم        | دکتر محمدرضا ربیعی  | استادیار   |       |



نام و نام خانوادگی رئیس دانشکده: دکتر ابراهیم هاشمی

تاریخ و امضاء و مهر دانشکده: ۹۷/۰۴/۲۵

تبصره: در صورتی که کسی مردود شود حداکثر یکبار دیگر (در مدت مجاز تحصیل) می تواند بار دیگر نام خود دفاع نماید (دفاع مجدد نباید زودتر از ۴ ماه برگزار شود).





تقدیم بہ ہمسرم

کہ در سایہ ہمکاری و ہمدلی او بہ این منظور نائل شدم.

تقدیم بہ پسرم

امید بخش جانم کہ آسایش او آرامش من است.

## سپاس‌گزاری...

سپاس بی‌کران پروردگار یکتا را که هستی مان بخشید و به طریق علم و دانش رهنمونمان شد و به همنشینی رهروان علم و دانش مفتخرمان نمود و خوشه‌چینی از علم و معرفت را روزیمان ساخت. اما از آنجایی که تجلیل از معلم، سپاس از انسانی است که هدف و غایت آفرینش را تامین می‌کند و سلامت امانت‌هایی را که به دستش سپرده‌اند، تضمین؛ بر حسب وظیفه و از باب ”من لم یشکر المنعم من المخلوقین لم یشکر الله عزّ و جلّ“:

از اساتید با کمالات و شایسته و دلسوز؛ سرکار خانم دکتر نگار اقبال و جناب آقای دکتر حسین باغیشنی که در کمال سعه صدر، با حسن خلق و فروتنی، از هیچ کمکی در این عرصه بر من دریغ ننمودند و زحمت راهنمایی این پایان‌نامه را بر عهده گرفتند نهایت تشکر و امتنان را دارم و برای آن بزرگواران آرزوی سلامتی و کامیابی از درگاه خداوند تبارک و تعالی می‌نمایم.

هم‌چنین بر خود لازم می‌دانم از اساتید گرانقدر، جناب آقای دکتر داود شاهسونی و جناب آقای دکتر محمدرضا ربیعی که زحمت داوری این پایان‌نامه را متقبل شدند، تشکر و قدردانی نمایم.

و وظیفه می‌دانم از همسر عزیزم که همواره در طول تحصیل متحمل زحماتم بود و تکیه‌گاه من در مواجهه با مشکلات و وجودش مایه دلگرمی من می‌باشد صمیمانه سپاس‌گزاری کنم.

در پایان از دوستان عزیز و بسیار مهربانم، سرکار خانم مریم علی‌بیگی و خانم محبوبه محبی به پاس محبت‌های بی‌دریغشان که در تمامی لحظات رفیق راهم بودند کمال تشکر و قدردانی را دارم.

اکرم قائمی‌زاده

تیر ۱۳۹۷



## تعهد نامه

اینجانب اکرم قائمی زاده دانشجوی کارشناسی ارشد رشته آمار دانشکده علوم ریاضی دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان **اسپلاین های جریمه ای و کاربردهای نوین آن**، تحت راهنمایی **نگار اقبال** متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است.
- در استفاده از نتایج پژوهش های دیگر پژوهش گران، به مرجع مورد استفاده استناد شده است.
- مطالب این پایان نامه، تا کنون توسط خود، یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا 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 به دست آوردن نتایج اصلی پایان نامه تاثیرگذار بوده اند، در مقالات مستخرج از پایان نامه رعایت می گردد.
- در تمام مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آن ها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در تمام مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته (یا استفاده شده است)، اصل رازداری و اصول اخلاق انسانی رعایت شده است.

اکرم قائمی زاده

تیر ۱۳۹۷

## مالکیت نتایج و حق نشر

- تمام حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه های رایانه ای، نرم افزارها و تجهیزات ساخته شده) متعلق به دانشگاه صنعتی شاهرود می باشد. این مطلب باید به نحو مقتضی، در تولیدات علمی مربوطه ذکر شود.
- استفاده از اطلاعات و نتایج موجود در این پایان نامه بدون ذکر منبع مجاز نمی باشد.



## چکیده

تحلیل رگرسیونی در حقیقت پرکاربردترین روش در بین تکنیک‌های آماری و یک ابزار آماری برای تشخیص رابطه یک یا چند متغیر تبیینی با متغیر پاسخ است. زمانی که برای تحلیل یک مجموعه داده مدل رگرسیونی خطی کارآمد نباشد، می‌توان از روش‌های ناپارامتری رگرسیون استفاده نمود. اسپلاین‌ها به عنوان یکی از ابزارهای درونیایی، در راستای رفع این مشکل، ابزاری ارزشمند و پرتوان به شمار می‌روند و از جمله روش‌های ناپارامتری برای مدل‌بندی رگرسیون ناپارامتری هستند. اسپلاین‌ها منحنی‌هایی را می‌سازند که شامل چندجمله‌ای‌هایی با درجه یکسان روی زیربازه‌هایی از یک بازه مشخص هستند و با شرایط پیوستگی تعریف‌شده‌ای به هم می‌پیوندند و از گره‌های مشترک بین دو زیربازه عبور می‌کنند. در طی سال‌هایی که از معرفی اسپلاین‌ها می‌گذرد، مبانی نظری آن‌ها توسعه یافته‌اند و نسخه‌های تعمیم‌یافته متفاوتی از آن‌ها معرفی شده‌اند. اسپلاین‌های جریمه‌ای ( $P$ -اسپلاین‌ها) از جمله ابزار پرکاربرد برای مدل‌سازی ناپارامتری و مسایل هموارسازی محسوب می‌شوند. ظهور  $P$ -اسپلاین‌ها با تحولات زیادی در مسایل هموارسازی همراه و به یک زمینه پویا و فعال تحقیقاتی تبدیل شده است. با توجه به کاربرد وسیع و اهمیت شناخت اسپلاین‌ها، در این پایان‌نامه به معرفی اسپلاین‌ها و روند رشد و توسعه آن‌ها می‌پردازیم و برخی از ویژگی‌های آن‌ها را مطرح می‌کنیم.

کلمات کلیدی: اسپلاین جریمه شده، هموارسازی، مدل‌های خطی تعمیم‌یافته، مدل‌های جمعی تعمیم‌یافته، معیار اطلاع آکائیک.



## لیست مقالات مستخرج از پایان نامه

۱. قائمی زاده ا. و اقبال ن، (۱۳۹۶)، ”سیر تحول اسپلین ها“، دومین کنفرانس ملی محاسبات نرم، گیلان – رودسر، ایران، ۸۴۷-۸۳۹.





## پیش‌گفتار

یکی از قدیمی‌ترین و پرکاربردترین روش‌های آماری، رگرسیون خطی می‌باشد. در مدل‌های کلاسیک آماری میانگین متغیر پاسخ به صورت تابع خطی از متغیرهای مستقل بیان می‌شود. مجموعه داده‌های زیادی وجود دارند که مدل‌های خطی، برازش خوبی برای آن‌ها ندارند. اسپلاین‌ها جایگزین مناسبی برای برازش به داده‌ها و توابع چندجمله‌ای تکه‌ای هستند که محدود به اتصال به یکدیگر در نقاط گره می‌باشند. با اتصال اولین تکه به تکه بعدی یک گره ایجاد می‌شود. بدین صورت که هر قسم از شکل تابع را در محدوده‌های موردنظر تکه تکه کرده و در هر تکه به هموارسازی تابع مشخص شده می‌پردازد. به طور معمول توابع چندجمله‌ای درجه سوم و اسپلاین‌های جریمه شده زیر مجموعه‌ی اسپلاین هموار شده است که مبتنی بر مدل رگرسیونی و برآورد ناپارامتریک می‌باشند. مدل رگرسیون ناپارامتریک اختصاص دارد به مدلی که در آن نتوان رابطه‌ی بین متغیر وابسته و متغیرهای مستقل را به فرم تابعی مشخص و معلوم بیان نمود. اسپلاین جریمه شده مدلی غیرخطی است که مقدار اثرات خطی را با در نظر گرفتن اثرات غیر خطی متغیرها بررسی می‌کند. زیرا ممکن است در صورت در نظر نگرفتن اثرات غیر خطی منجر به از دست دادن اطلاعات شود که بسته به اهداف مطالعه ممکن است تاثیرگذار باشد در نتیجه با حضور اثرات غیر خطی، به بررسی اثرات خطی می‌پردازد. وسعت کاربرد اسپلاین‌های جریمه‌ای بسیار گسترده است به عنوان مثال می‌توان به ساختار رگرسیون جمعی خطی و جمعی تعمیم‌یافته، مدل با اثرات متغیر(وابسته به زمان، مکان یا هر دو)، مدل‌هایی با متغیرهای تبیینی تابعی، برآورد توابع دو یا چند متغیره و مدل‌های رگرسیونی چندکی و گشتاوری اشاره کرد.

در این پایان‌نامه با توجه به ویژگی‌های نظری و کاربردی خوب اسپلاین‌ها، تمرکز بر معرفی و استفاده از آن‌ها در مدل‌های مختلف آماری است. در فصل اول تعاریفی از انواع مدل‌های خطی از قبیل مدل‌های خطی تعمیم‌یافته، مدل‌های جمعی و مدل‌های جمعی تعمیم‌یافته برای مکان، مقیاس و شکل (GAMLSS) آورده شده است. در فصل دوم به تفصیل اسپلاین‌ها و انواع آن می‌پردازیم. در فصل سوم انتخاب پارامتر هموارسازی و بهترین مدل برازشی را براساس معیارهای CV, BIC, AIC و Cp معرفی می‌کنیم و در نهایت فصل چهارم، مجموعه داده‌هایی را با استفاده از این مدل‌ها مورد تحلیل و بررسی قرار می‌دهیم.



# فهرست مطالب

ق فهرست تصاویر

ش فهرست جداول

|    |   |    |
|----|---|----|
| ۱  | مفاهیم اولیه و مورد نیاز  | ۱  |
| ۱  | ۱.۱ مقدمه   | ۱  |
| ۱  | ۲.۱ مدل‌های خطی   | ۱  |
| ۳  | ۳.۱ مدل‌های خطی تعمیم‌یافته                                     | ۳  |
| ۴  | ۱.۳.۱ ساختار مدل‌های خطی تعمیم‌یافته                            | ۴  |
| ۵  | ۴.۱ مدل جمعی تعمیم‌یافته  | ۵  |
| ۶  | ۱.۴.۱ مدل‌های جمعی تعمیم‌یافته برای مکان، مقیاس و شکل           | ۶  |
| ۹  | ۲ سیر تحول اسپلین‌ها و انواع آن‌ها                              | ۹  |
| ۹  | ۱.۲ توابع پایه  | ۹  |
| ۱۰ | ۲.۲ تابع توانی بریده‌شده  | ۱۰ |
| ۱۱ | ۳.۲ اسپلین‌ها   | ۱۱ |
| ۱۲ | ۴.۲ اسپلین مکعبی (CS)   | ۱۲ |
| ۱۳ | ۱.۴.۲ درونیابی اسپلین مکعبی                                     | ۱۳ |
| ۱۳ | ۲.۴.۲ خواص اسپلین‌های مکعبی                                     | ۱۳ |
| ۱۶ | ۳.۴.۲ اسپلین مکعبی طبیعی (NCS)                                  | ۱۶ |
| ۱۷ | ۵.۲ B-اسپلین‌ها   | ۱۷ |
| ۲۱ | ۶.۲ هموارسازی اسپلینی   | ۲۱ |
| ۲۳ | ۷.۲ اسپلین‌های جریمه‌ای (p-اسپلین)                              | ۲۳ |
|    | ۱.۷.۲ صورت ماتریسی و برآورد ضرایب اسپلین جریمه‌ای با استفاده از |    |
| ۲۴ | کمترین توان‌های دوم   | ۲۴ |
| ۲۶ | ۲.۷.۲ خواص اسپلین جریمه‌ای                                      | ۲۶ |

|    |  |      |
|----|--|------|
| ۲۷ | اسپلین حاصل ضرب تانسور . . . . .   | ۸.۲  |
| ۳۰ | اسپلین صفحه نازک . . . . .   | ۹.۲  |
| ۳۲ | مزیای اسپلین صفحه نازک . . . . .   | ۱۰.۲ |
| ۳۵ | <b>۳ برخی از مباحث نظری اسپلین‌ها</b>  |      |
| ۳۵ | ۱.۳ نحوه‌ی انتخاب پارامتر هموارسازی . . . . .  |      |
| ۳۶ | ۱.۱.۳ معیار اطلاع آکائیک . . . . .   |      |
| ۳۷ | ۲.۱.۳ اعتبارسنجی متقابل . . . . .  |      |
| ۳۸ | ۳.۱.۳ اعتبارسنجی متقابل تعمیم‌یافته . . . . .  |      |
| ۳۹ | ۲.۳ روش‌های انتخاب تعداد و موقعیت گره . . . . .                                      |      |
| ۴۱ | ۳.۳ انتخاب مدل . . . . .   |      |
| ۴۲ | ۱.۳.۳ معیار اطلاع بیزی (شوارتز) . . . . .  |      |
| ۴۲ | ۲.۳.۳ آماره $C_p$ مالو . . . . .   |      |
| ۴۷ | <b>۴ برخی از کاربردهای اسپلین‌ها</b>   |      |
| ۴۸ | ۱.۴ مدل‌های جمعی . . . . .   |      |
| ۴۸ | ۱.۱.۴ کاربرد اول: داده‌های اجاره مونیخ . . . . .                                     |      |
| ۵۵ | ۲.۱.۴ کاربرد دوم: داده‌های گونه‌های ماهی . . . . .                                   |      |
| ۵۹ | ۳.۱.۴ کاربرد سوم: داده‌های مدت بستری در بیمارستان . . . . .                          |      |
| ۶۱ | ۲.۴ اسپلین دو بعدی . . . . .   |      |
| ۶۱ | ۱.۲.۴ داده‌های فیلم . . . . .  |      |
| ۶۶ | ۳.۴ مدل میدان‌های تصادفی گاوسی . . . . .   |      |
| ۶۷ | ۱.۳.۴ داده‌های جرم و جنایت . . . . .   |      |
| ۶۸ | ۴.۴ نتیجه‌گیری . . . . .   |      |
| ۶۹ | <b>آ</b>   |      |
| ۶۹ | ۱.آ ضرب کرونکر . . . . .   |      |
| ۷۰ | ۲.آ مشتق تابع B-اسپلین . . . . .   |      |
| ۷۱ | ۳.آ جدول توزیع‌های پیوسته و گسسته موجود در بسته‌ی <code>gamlss.dist</code> . . . . . |      |
| ۷۲ | ۴.آ دستورات نرم‌افزار R . . . . .  |      |
| ۸۹ | <b>مراجع</b>   |      |



# فهرست تصاویر

|    |   |      |
|----|---|------|
| ۱۱ | توابع توانی بریده شده مرتبه یک  | ۱.۲  |
| ۱۷ | نمایش اسپلاین های مکعبی و مکعبی طبیعی   | ۲.۲  |
| ۱۹ | نمایش توابع پایه B-اسپلاینی مرتبه یک  | ۳.۲  |
| ۲۰ | نمایش توابع پایه B-اسپلاین  | ۴.۲  |
| ۲۲ | هموارسازی   | ۵.۲  |
| ۲۲ | فلوچارت روش های هموارساز  | ۶.۲  |
| ۲۶ | هموارسازی تابع رگرسیونی برای مقادیر مختلف $\lambda$                           | ۷.۲  |
| ۲۸ | حاصل ضرب دو تابع پایه ی حاشیه ای برای توابع هموار $fx$ و $fz$                 | ۸.۲  |
| ۳۱ | نمایش اسپلاین صفحه نازک   | ۹.۲  |
| ۳۲ | نمایش اسپلاین صفحه نازک هموارتر   | ۱۰.۲ |
| ۳۳ | نمایش اسپلاین صفحه نازک   | ۱۱.۲ |
|    | پارامتر هموارسازی در مقابل معیار انتخاب برای تابع رگرسیونی P-اسپلاینی از      | ۱.۳  |
| ۳۹ | مرتبه $m = 4$   |      |
| ۴۰ | برازش P-اسپلاین مرتبه $k = 4$ با تعداد گره کم                                 | ۲.۳  |
| ۴۰ | برازش P-اسپلاین مرتبه $k = 4$ با تعداد گره زیاد                               | ۳.۳  |
|    | برازش اسپلاین مکعبی طبیعی با چهار درجه آزادی و قرار گرفتن موقعیت گره ها       | ۴.۳  |
| ۴۱ | در چندک داده ها   |      |
| ۴۳ | برازش چند جمله ای از مرتبه یک تا مرتبه هفت برای مجموعه داده های شبیه سازی شده | ۵.۳  |
| ۴۴ | مقایسه اثربخشی معیار اطلاع اکائیک و بیزی در انتخاب بهترین مدل                 | ۶.۳  |
| ۴۴ | انتخاب بهترین مدل با استفاده از معیار اعتبارسنجی متقابل                       | ۷.۳  |
| ۴۹ | نمودار متغیر پاسخ (میزان اجاره R) در مقابل هریک از متغیرهای توضیحی            | ۱.۴  |
| ۵۰ | نمودار مانده های برازش مدل خطی  | ۲.۴  |
| ۵۰ | نمودار مانده های برازش مدل گاما   | ۳.۴  |
| ۵۱ | نمودار مقادیر برازشی برای مدل (GAM)   | ۴.۴  |
| ۵۲ | نمودار ماریپیچ برای مدل (GAM)   | ۵.۴  |

|    |  |      |
|----|--|------|
| ۵۴ | نمودار مارپیچ برای مدل گامای دو پارامتری                                     | ۶.۴  |
| ۵۵ | نمودار مارپیچ برای مدل سه پارامتری   | ۷.۴  |
| ۵۶ | نمودار گونه‌های ماهی   | ۸.۴  |
| ۵۸ | نمودار برازش مدل سیچل به داده‌های گونه‌های ماهی                              | ۹.۴  |
| ۵۸ | نمودار مارپیچ مدل برازشی سیچل به داده‌های گونه‌های ماهی                      | ۱۰.۴ |
|    | نمودار نرخ نامناسب متغیر پاسخ در مقابل متغیرهای توضیحی در داده‌های مدت بستری | ۱۱.۴ |
| ۵۹ | توابع برازشی برای $\mu$ در مدل M4  | ۱۲.۴ |
| ۶۱ | ترسیم نمودار سه بعدی داده‌های فیلم   | ۱۳.۴ |
| ۶۲ | نمودار مارگون برای مدل M6 با برازش رویه برای پارامتر $\mu$                   | ۱۴.۴ |
| ۶۳ | نمودار رویه‌ی برازشی برای مدل M6   | ۱۵.۴ |
| ۶۴ | نمودار مارگون برای مدل M7 با برازش رویه برای پارامترهای $\mu$ و $\sigma$     | ۱۶.۴ |
| ۶۵ | نمودار مارگون قاب سمت چپ مدل M8 و قاب سمت راست مدل M9                        | ۱۷.۴ |
| ۶۶ | نمودار مارگون برای مدل M9 با دو متغیر توضیحی                                 | ۱۸.۴ |
| ۶۸ | برازش مقادیر برای داده‌های جرم جنایت با استفاده از تابع GMRF                 | ۱۹.۴ |

# فهرست جداول

|    |  |      |
|----|--|------|
| ۵  | جدول توزیع‌های خانواده نمایی   | ۱.۱  |
| ۴۵ | جدول مقادیر معیارهای $AIC$ , $CV$ , $BIC$ , $C_p$ برای انتخاب بهترین مدل         | ۱.۳  |
| ۵۲ | تفسیر الگوهای مختلف در نمودار ماریچ  | ۱.۴  |
| ۵۳ | گزینش بهترین مدل با استفاده معیارهای $AIC$                                       | ۲.۴  |
| ۵۳ | گزینش بهترین مدل با استفاده معیارهای $AIC$                                       | ۳.۴  |
| ۵۴ | گزینش بهترین مدل با استفاده معیارهای $AIC$                                       | ۴.۴  |
| ۵۶ | مقادیر به‌دست آمده با استفاده از $AIC$   | ۵.۴  |
| ۵۷ | مقادیر به‌دست آمده با استفاده از معیار $AIC$                                     | ۶.۴  |
| ۵۷ | مقادیر به دست آمده با استفاده از معیار $AIC$                                     | ۷.۴  |
| ۶۰ | گزینش بهترین مدل با استفاده معیارهای $AIC$ و $BIC$                               | ۸.۴  |
| ۶۲ | گزینش بهترین مدل با استفاده معیارهای $AIC$ و $BIC$                               | ۹.۴  |
| ۶۳ | گزینش بهترین مدل با استفاده معیارهای $AIC$                                       | ۱۰.۴ |
| ۶۵ | گزینش بهترین مدل با استفاده معیارهای $AIC$                                       | ۱۱.۴ |
| ۶۷ | مدل برازشی به داده‌های columb با استفاده از تابع $gmrf$                          | ۱۲.۴ |
|    | خلاصه‌ای از توزیع‌های پیوسته موجود در بسته‌ی $gamlss.dist$ با تابع پیوند پیش فرض | ۱.آ  |
| ۷۱ |  |      |
|    | خلاصه‌ای از توزیع‌های گسسته موجود در بسته‌ی $gamlss.dist$ با تابع پیوند پیش فرض  | ۲.آ  |
| ۷۲ |  |      |



# فصل ۱

## مفاهیم اولیه و مورد نیاز

### ۱.۱ مقدمه

در این فصل مفاهیم اساسی و اولیه که در فصل‌های آینده به آن نیاز داریم را معرفی می‌کنیم. ابتدا مدل خطی که یکی از قدیمی‌ترین و پرکاربردترین روش‌های آماری است را بیان کرده و در ادامه مدل‌های خطی تعمیم‌یافته، مدل‌های جمعی و مدل‌های جمعی تعمیم‌یافته برای مکان، مقیاس و شکل (GAMLSS) را شرح می‌دهیم.

### ۲.۱ مدل‌های خطی

مدل رگرسیون خطی<sup>۱</sup> یک مدل ساده اما کاربردی است که سالیان متمادی مبنای تجزیه و تحلیل داده‌ها را تشکیل می‌دهد. این تکنیک زمانی که رابطه بین متغیر وابسته و متغیرهای مستقل خطی است، ابزاری قدرتمند است و تکیه بر آن در بسیاری از کاربردها، محدود کننده است زیرا در بسیاری از پدیده‌ها روابط خطی قابل تعریف نیستند [۱۳]. یک مدل رگرسیون خطی بر پایه‌ی فرض‌های اساسی زیر ارائه می‌شود:

- استقلال مشاهدات

---

<sup>۱</sup>Linear regression model



• وجود رابطه خطی بین متغیرهای پاسخ و توضیحی

• ثابت بودن واریانس مشاهدات

در این مدل  $X$  به عنوان متغیر مستقل یا پیش گو عمل می نماید که مقادیر آن به وسیله آزمایشگر کنترل می شود این متغیر کنترل شده را متغیر پیش بین می نامند و  $Y$  که به متغیر  $X$  وابسته است را متغیر اثر یا پاسخ می نامند. فرض کنید متغیر پاسخ  $Y_i$  با رابطه خطی خطی زیر به متغیر مستقل  $X_i$ ، در ارتباط باشد:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad i = 1, 2, \dots, n \quad (1.1)$$

که در آن  $\beta_0$  و  $\beta_1$  مجموعه پارامترهای نامعلوم اند که  $\beta_0$  عرض از مبدا و  $\beta_1$  شیب خط رگرسیونی است که ضرایب رگرسیونی نامیده می شوند و  $\epsilon$  خطاست که متغیرهای غیرقابل مشاهده اند که دارای میانگین صفر و واریانس نامعلوم  $\sigma^2$  است. چنانچه بیش از یک متغیر مستقل در مدل داشته باشیم، مدل رگرسیون چندجمله ای حاصل می گردد که به صورت زیر است:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_r X_{ir} + \epsilon_i \quad i = 1, 2, \dots, n$$

در عمل، برآورد  $\beta = (\beta_0, \beta_1, \dots, \beta_r)$  را به گونه ای برمی گزینیم که مجموع توان های دوم انحراف ها از خط حقیقی را کمینه کند. با به کارگیری روش کمترین توان های دوم، برآوردهای زیر حاصل می شود:

$$\hat{\beta} = \operatorname{argmin}_{\beta} (Y - X\beta)^T (Y - X\beta)$$

در نتیجه

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad (2.1)$$

همچنین برآورد  $\beta$  از روش ماکزیمم درستنمایی (MLE) به صورت (۲.۱) است. برآورد پارامتر در روش درستنمایی ماکزیمم به توزیع متغیر پاسخ بستگی دارد، چون توزیع آن نرمال است بنابراین برآورد پارامتر از هر دو روش یکسان است. با برآورد پارامترهای مدل، مقادیر برازشی  $\hat{Y}$  از رابطه ی زیر حاصل می شود:

$$\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = HY$$

$H$  ماتریس هت نامیده می شود.

تفاوت بین مقادیر واقعی و مقادیر برازشی، مانده های رگرسیونی نامیده می شود.

$$\hat{\epsilon} = Y - \hat{Y} = Y(I - H)$$

درجه آزادی یک مدل خطی، تعداد پارامترهای  $\beta$  که مستقل‌اند، به عبارتی تعداد ستون‌های مستقل  $X$  نشان دهنده اثر ماتریس هت است [۱].

$$df = tr(H) = p$$

از اهداف اصلی رگرسیون، ارائه یک خلاصه یا کاهش داده‌های مشاهده شده به منظور کشف و ارائه‌ی ارتباط بین متغیرهای پاسخ و توضیحی است و یک هدف مهم دیگر آن، استفاده از مدل (۱.۱) و یافتن رابطه‌ای است که بتوان اثر تغییرات یک یا چند متغیر را بر روی متغیرهای دیگر پیش‌بینی کرد. پیش‌بینی یک جنبه بسیار مهم رگرسیون است. این روش دارای مزایا و معایبی است. یک مزیت این است که این روش را می‌توان به طیف گسترده‌ای از ساختارهای داده‌ها گسترش داد. برای مثال مدل‌های خطی تعمیم‌یافته توسط مک کولاغ و نلدر در سال ۱۹۸۹ معرفی شدند [۱۱].

## ۳.۱ مدل‌های خطی تعمیم‌یافته

رویکرد مدل‌های خطی تعمیم‌یافته<sup>۲</sup>، یک مدل پارامتری بوده و یک تعمیم انعطاف‌پذیری از مدل‌های خطی می‌باشد. هنگامی که متغیر پاسخ  $Y$  شمارشی باشد یا رابطه بین متغیر پاسخ و توضیحی خطی نباشد دیگر مدل‌های خطی کلاسیک جوابگوی این‌گونه مسائل نیستند و برای این منظور، مدل‌های خطی تعمیم‌یافته را به کار می‌بریم، که توسط جان نلدر و رابرت دربرن (۱۹۷۲) ارائه شدند و توسط مک کولاغ و نلدر در سال ۱۹۸۹ توسعه یافتند و برای انجام آن از بسته‌های محاسباتی آماری استفاده کردند. در این مدل متغیر پاسخ  $Y$  دارای توزیع نرمال است ولی در بسیاری از موقعیت‌های عملی، این فرض حتی به‌طور تقریبی برقرار نیست. مدل خطی تعمیم‌یافته به‌طور گسترده مدل‌های آماری پرکاربرد را در برمی‌گیرد از قبیل رگرسیون خطی برای پاسخ‌هایی که به‌صورت نرمال توزیع شده‌اند، مدل لجستیک برای داده‌های دوتایی (دودویی) و مدل‌های لگاریتم خطی برای داده‌های شمارشی (نوعی از داده آماری که در آن مشاهدات صرفاً می‌توانند مقادیر اعداد صحیح غیرمنفی بگیرند). مدل‌های خطی تعمیم‌یافته این امکان را فراهم می‌سازد که متغیر پاسخ از خانواده توزیع نمایی تبعیت کند. خانواده نمایی شامل توزیع‌های نرمال، دوجمله‌ای، پواسن، هندسی، نمایی، گاما است.

در خانواده مدل‌های خطی تعمیم‌یافته دو مؤلفه مهم وجود دارد [۳، ۲۸، ۳۴]:

• توزیع متغیر پاسخ

• مدلی که میانگین متغیر پاسخ را به متغیرهای توضیحی  $X$  ارتباط دهد.

در مدل‌های خطی تعمیم‌یافته، واریانس مشاهدات ثابت فرض شده و تابعی از میانگین است.

<sup>۲</sup> Generalized linear model

### ۱.۳.۱ ساختار مدل‌های خطی تعمیم‌یافته

مدل‌های خطی تعمیم‌یافته دارای ویژگی‌های کلی زیر است:

- متغیرهای پاسخ  $Y_1, Y_2, \dots, Y_n$  مستقل با میانگین‌های  $\mu_1, \dots, \mu_n$  و واریانس‌های  $\sigma_1^2, \dots, \sigma_n^2$  می‌باشد و دارای توزیعی از خانواده توزیع‌های نمایی است.
- پیش‌گوی خطی مدل به صورت

$$\eta_i = \beta_0 + X_i^T \beta = \beta_0 + \beta_1 X_{1i} + \dots + \beta_p X_{pi} = \beta_0 + \sum_{j=1}^p \beta_j X_{ji}$$

تعریف می‌شود که شامل  $X_i$  یک بردار  $(p \times 1)$  از متغیرهای توضیحی و  $\beta$  بردار  $(p \times 1)$  از پارامترها می‌باشد.

- تابع پیوند  $g$  میانگین  $\mu_i$  را با یک تابع معلوم و مشتق‌پذیر از میانگین مدل‌سازی می‌کند. این تابع، میانگین پاسخ را در پیش‌گوی خطی برای مقادیر  $i = 1, \dots, n$  به صورت  $\eta_i = g(\mu_i)$  بیان می‌کند. هم‌چنین داریم:

$$\mu_i = E(Y_i) = g^{-1}(\eta_i)$$

تابع پیوند اغلب غیرخطی و بدون محدودیت است اما میانگین بعضی از توزیع‌ها مانند توزیع دو جمله‌ای محدودیت دارد. در مدل‌های خطی کلاسیک، تابع پیوند همان تابع همانی است، بنابراین  $\mu = \eta$  و هر مقداری از اعداد حقیقی را می‌تواند اختیار کند.

تابع چگالی متغیر پاسخ  $Y$  با میانگین  $\mu = E(Y)$  و پارامترهای  $\theta$  و  $\phi$  متعلق به خانواده‌ی توزیع‌های نمایی است اگر بتوان آن را به صورت زیر نوشت:

$$f_Y(y|\theta, \phi) = \exp\left\{\frac{y\theta - b(\theta)}{\phi} + c(y, \phi)\right\} \quad (3.1)$$

که در آن  $\theta$  پارامتر طبیعی یا پارامتر کانونی و  $\phi$  پارامتر مقیاس یا پراکندگی نامیده می‌شود که به‌طور معمول، معلوم فرض می‌شود و تابعی از میانگین متغیر پاسخ است. در نتیجه

$$E(Y) = b'(\theta) \quad (4.1)$$

$$Var(Y) = \phi V(\mu)$$

با توجه به ویژگی  $\int f(y|\theta, \phi) dy = 1$  و مشتقات اول و دوم نسبت به  $\theta$  از هر دو طرف معادله، نتایج (۴.۱) حاصل می‌شود.

$$\int \phi^{-1}(y - b'(\theta))f(y|\theta, \phi) dy = 0$$

بنابراین

$$\mu = E(Y) = b'(\theta)$$

و با توجه به مشتق دوم داریم:

$$\int \{\phi^{-1}[y - b'(\theta)]^2 - b''(\theta)\} f(y|\theta, \phi) dy = 0$$

در نتیجه

$$Var(Y) = \phi b''(\theta) = \phi b''(b'^{-1}(\mu)) = \phi V(\mu)$$

$V$  تابع واریانس نامیده می شود و تابعی از  $\mu$  است.

در جدول ۱.۱ برخی از توزیع ها که متعلق به خانواده ی توزیع های نمایی به همراه تابع پیوند مورد بررسی قرار گرفته اند [۳، ۲۸، ۳۴]:

جدول ۱.۱: جدول توزیع های خانواده نمایی

| توزیع                    | پیوند    | $\theta(\mu)$             | $\phi$     | $c(y, \phi)$                                 | $b(\theta)$          | $\mu$                           | $Var(Y)$                            |
|--------------------------|----------|---------------------------|------------|--|----------------------|---------------------------------|-------------------------------------|
| نرمال $N(\mu, \sigma^2)$ | همانی    | $\mu$                     | $\sigma^2$ | $\frac{-y^2}{2\phi} - \log(\sqrt{2\pi\phi})$ | $\frac{\theta^2}{2}$ | $\theta$                        | $\sigma^2$                          |
| برنولی $B(1, \mu)$       | لوجیت    | $\log(\frac{\mu}{1-\mu})$ | ۱          | ۰  | $\log(1 + e^\theta)$ | $\frac{e^\theta}{1 + e^\theta}$ | $\frac{e^\theta}{(1 + e^\theta)^2}$ |
| پواسن $P(\mu)$           | لگاریتمی | $\log(\mu)$               | ۱          | $-\log(y!)$                                  | $\exp(\theta)$       | $\exp(\theta)$                  | $\exp(\theta)$                      |

## ۴.۱ مدل جمعی تعمیم یافته

این مدل، یک مدل ناپارامتری بوده و بسط مدل های خطی تعمیم یافته است که خود بسط مدل های خطی می باشند. مدل جمعی تعمیم یافته<sup>۳</sup> روشی بسیار مناسب برای بررسی رابطه بین متغیر پاسخ و متغیرهای مستقل و همچنین تحلیل داده ها ارائه می دهد و به وسیله ی بسته نرم افزاری mgcvR قابل اجرا می باشد. در مدل های جمعی تعمیم یافته بر خلاف مدل رگرسیون خطی، داده ها شکل منحنی پاسخ را مشخص می کنند. در این مدل فرض بر این است که متغیر پاسخ  $Y$  دارای توزیعی از خانواده ی توزیع های نمایی با میانگین  $\mu = E(Y|x_1, \dots, x_p)$  می باشد که از طریق تابع پیوند  $g$  به متغیرهای پیش گو  $X$  متصل می شود. ساختار مدل های جمعی تعمیم یافته به صورت زیر است:

$$\eta = g(\mu) = \beta_0 + s_1(x_1) + \dots + s_J(x_J)$$

<sup>۳</sup> Generalized additive models

که در آن  $s_j$  ها  $j = 1, \dots, J$  توابع ناپارامتری هموار و نامعلوم می‌باشند. در مدل‌های جمعی تعمیم‌یافته  $\sum_{j=1}^J s_j(x_j)$  جانشین توابع خطی  $\sum_{j=1}^p \beta_j X_j$  در مدل‌های خطی تعمیم‌یافته می‌شود. دارا بودن هموارسازها (اسپلاین‌های مکعبی، اسپلاین‌های جریمه‌ای، لوییس...) یکی از مزایای مهم این مدل است که باعث توانایی این مدل در شناسایی روابط غیرخطی شده است. به عبارتی با اطلاعات بیشتری از روابط بین داده‌ها، کیفیت پیش‌بینی متغیر پاسخ را به حداکثر می‌رساند. همچنین به خاطر انعطاف‌پذیر بودن در تعیین نوع و درجه ارتباط به یک مدل محبوب تبدیل شده است [۲۸].

### ۱.۴.۱ مدل‌های جمعی تعمیم‌یافته برای مکان، مقیاس و شکل

مدل‌های جمعی تعمیم‌یافته برای مکان، مقیاس و شکل<sup>۴</sup> (GAMLSS) دارای یک ساختار کلی برای برازش انواع مدل‌های رگرسیونی است که توسط ریگبی و استامینوپولوس<sup>۵</sup> (۲۰۰۱، ۲۰۰۵) معرفی شدند و بعدها برای غلبه بر برخی از محدودیت‌های مرتبط با مدل‌های خطی تعمیم‌یافته و مدل‌های جمعی تعمیم‌یافته گسترش یافت. این مدل در حقیقت قادر به ارائه یک مدل هموار و انعطاف‌پذیر است، فرض می‌شود که متغیر پاسخ از توزیعی پارامتری هموار است که ممکن دم کلفت یا با کشیدگی مثبت یا منفی باشد پیروی می‌کند، علاوه بر این تمام پارامترهای توزیع، مکان (میانگین)، مقیاس (واریانس) و شکل توزیع (کشیدگی یا چولگی) را می‌توان به عنوان توابع خطی یا غیرخطی و یا هموارساز از متغیرهای توضیحی مدل‌سازی کرد [۲۴، ۱۷، ۲۸].

### ساختار مدل

در مدل جمعی تعمیم‌یافته برای مکان، مقیاس و شکل فرض می‌شود برای هر  $i = 1, 2, \dots, n$  متغیر پاسخ  $Y_i$  دارای تابع احتمال  $f_Y(y_i | \mu_i, \sigma_i, \nu_i, \tau_i)$  مشروط بر  $(\mu_i, \sigma_i, \nu_i, \tau_i)$  بردار چهار پارامتر توزیع که هر کدام می‌تواند تابعی برای متغیرهای توضیحی باشد و آن را با  $Y_i \sim D(\mu, \sigma, \nu, \tau)$  نشان می‌دهند.

پارامترهای  $(\mu_i, \sigma_i, \nu_i, \tau_i)$  پارامترهای توزیع می‌باشد. پارامترهای شرطی  $\mu$  و  $\sigma$  به ترتیب برای مکان و مقیاس و همچنین  $\nu$  و  $\tau$  برای شکل توزیع (کشیدگی و چولگی) به کار می‌روند. این مدل می‌تواند برای پارامترهای هر توزیع به کار رود و همچنین می‌تواند به بیش از چهار پارامتر توزیع تعمیم یابد.

با فرض این که  $Y^T = (Y_1, Y_2, \dots, Y_n)$  برداری با طول  $n$  متغیر پاسخ باشد فرموله کردن GAMLSS به صورت زیر است:

<sup>۴</sup> Generalized additive model for location, scale and shape

<sup>۵</sup> Rigby and stasinopoulos

برای  $k = 1, 2, 3, 4$  و فرض  $g_k(\cdot)$  تابع پیوند یکنوا پارامترهای توزیع

$$\theta_k = (\theta_1, \theta_2, \theta_3, \theta_4) = (\mu_i, \sigma_i, \nu_i, \tau_i)$$

به متغیر پیش بینی  $\eta_k$  داریم:

$$g_k(\theta_k) = \eta_k = X_k \beta_k + \sum_{j=1}^{J'_k} h_{jk}(x_{jk}), \quad k = 1, 2, 3, 4$$

که در آن  $\beta_k^T = (\beta_{1k}, \beta_{2k}, \dots, \beta_{J'_k k})$  بردار پارامترها به طول  $J'_k$ ،  $X_k$  ماتریس طرح با بعد  $n \times J'_k$ ،  $h_{jk}$  ماتریس طرح اثرات تصادفی با بعد  $n \times q_{jk}$  و  $x_{jk}$  (پارامترهای اثرات تصادفی) متغیر تصادفی  $q_{jk}$  بعدی که دارای توزیع

$$x_{jk} \sim N(\circ, G_{jk}^{-1})$$

به طوری که  $G_{jk}^{-1}$  معکوس تعمیم یافته ی ماتریس متقارن  $q_{jk} \times q_{jk}$  به شکل  $G_{jk} = G_{jk}(\lambda_{jk})$  است که می تواند به بردار ابر پارامتر  $\lambda_{jk}$  وابسته باشد. به دلیل دستیابی به معادلات بسیار ساده در روند برازش مدل، متغیرهای تصادفی را به توزیع نرمال محدود می کنیم که در حالت کلی این محدودیت ممکن است حذف شود.

اگر عبارت جمعی در پارامترهای توزیع نباشد مدل GAMLSS پارامتری خطی به صورت زیر است:

$$g_1(\mu) = \eta_1 = X_1 \beta_1$$

$$g_2(\sigma) = \eta_2 = X_2 \beta_2$$

$$g_3(\nu) = \eta_3 = X_3 \beta_3$$

$$g_4(\tau) = \eta_4 = X_4 \beta_4$$

با قرار دادن تابع  $h_k(X_k, \beta_k)$  که  $h_k$  که تابع غیرخطی است به جای  $X_k \beta_k$  مدل GAMLSS پارامتری غیرخطی به صورت زیر به دست می آید:

$$g_1(\mu) = \eta_1 = h_1(X_1, \beta_1)$$

$$g_2(\sigma) = \eta_2 = h_2(X_2, \beta_2)$$

$$g_3(\nu) = \eta_3 = h_3(X_3, \beta_3)$$

$$g_4(\tau) = \eta_4 = h_4(X_4, \beta_4)$$

پس از تعیین مدل، بردار پارامترهای  $\beta_k$  و پارامترهای اثرات تصادفی  $x_{jk}$  را می توان با ماکسیم کردن تابع درست نمایی توانیده  $l_p$  که به صورت زیر است برآورد نمود:

$$\ell_p = \ell - \frac{1}{\varphi} \sum_{k=1}^p \sum_{j=1}^{J_k} x'_{jk} G_{jk}(\lambda_{jk}) x_{jk}$$

که  $\ell = \sum_{i=1}^n \log\{D(y_i|\mu_i, \sigma_i, \nu_i, \tau_i)\}$  لگاریتم تابع درستنمایی است.  
 خواص توابع GAMLSS [۲۴، ۱۷، ۲۸]:

- چارچوبی بسیار انعطاف‌پذیر برای انواع مدل‌های رگرسیونی ایجاد می‌کند.
- اجرا و پیاده‌سازی آن‌ها برای توزیع‌های جدید و همچنین عبارات جمعی جدید آسان است.
- در این تابع، متغیر پاسخ می‌تواند هر توزیعی داشته باشد و تمام پارامترهای توزیع می‌توانند به‌عنوان تابعی از متغیرهای توضیحی مدل شوند.

## فصل ۲

# سیر تحول اسپلاین‌ها و انواع آن‌ها

اسپلاین‌ها توابعی شامل چند جمله‌ای‌هایی با مرتبه  $m$  روی زیربازه‌هایی از یک بازه مشخص هستند و با شرایط پیوستگی تعریف شده‌ای به هم می‌پیوندند و از مرزهای مشترک بین دو زیربازه، موسوم به گره، عبور می‌کنند. اسپلاین‌ها ترکیب خطی از توابع پایه و وزن‌هایی برای هر تابع پایه، ساخته می‌شوند. توابع پایه توابعی بر حسب مشاهدات هستند و وزن‌های توابع پایه به عنوان پارامتر قلمداد می‌شوند که می‌توان با روش‌هایی هم‌چون روش کمترین توان‌های دوم خطا یا روش درست‌نمایی ماکسیمم آن‌ها را برآورد کرد. در طی سال‌هایی که از معرفی اسپلاین‌ها می‌گذرد، مبانی نظری آن‌ها توسعه یافته‌اند و نسخه‌های تعمیم‌یافته متفاوتی از آن‌ها معرفی شده‌اند. حوزه کاربردهای این ابزار قدرتمند ناپارامتری نیز به شدت گسترش یافته است.

اسپلاین‌ها انواع مختلفی هم‌چون اسپلاین‌های هموار، رگرسیون اسپلاینی، و B-اسپلاین‌ها دارند. برای اطلاعات جامع‌تر در این زمینه می‌توانید به هیستی و تیپشیرانی [۱۲]، دی بور [۶]، واهبا [۳۲]، گرین و سیلورمن [۱۱] مراجعه کنید.

## ۱.۲ توابع پایه

در جبر خطی منظور از یک پایه در یک فضای برداری، مجموعه‌ای از بردارهای موجود در آن فضا است به طوری که مستقل خطی باشند و هر بردار دیگر در فضای برداری را بتوان از



ترکیب خطی آن بردارها به‌دست آورد. برای مثال چندجمله‌ای درجه دو با ضرایب حقیقی را می‌توان به‌صورت

$$y = 1a + bt + ct^2$$

نوشت که در واقع از ترکیب خطی توابع پایه  $t^2, t, 1$  شکل گرفته است [۳۰].

## ۲.۲ تابع توانی بریده‌شده

یک منحنی اسپالین از ترکیب خطی توابع پایه و پارامترهای اسپالین تشکیل می‌شود که تابع توانی بریده‌شده عضوی از مجموعه توابع پایه است.

فرض کنید تعداد مشاهدات  $x_1, \dots, x_n$  داشته باشیم و یک خط راست روی بازه  $[t_0, t_{m+1}]$  که  $t_0 = \min(x_i)$  و  $t_{m+1} = \max(x_i)$  رسم شده باشد و خط را در نقطه  $t_1$  که  $t_0 < t_1 < t_{m+1}$  به دو قسمت تقسیم می‌کنیم به‌طوری که پیوستگی آن حفظ شود و می‌توان معادله خط که در  $t_1$  شکسته شده را به‌صورت ترکیب خطی از توابع پایه ۱ و  $x$  و  $(x - t_1)_+$  نوشت که  $(x - t_1)_+$  را تابع توانی بریده‌شده می‌نامند که به‌صورت زیر تعریف می‌شود.

$$(x - t_1)_+ = \max\left\{0, (x - t_1)\right\} = \begin{cases} x - t_1 & x \geq t_1 \\ 0 & o.w. \end{cases}$$

برای به‌دست آوردن معادله خط، تابع  $f$  که در  $t_1$  شکسته شده در بازه  $[t_0, t_1]$  دارای یک معادله خط و در بازه  $[t_1, t_{m+1}]$  دارای معادله خط دیگری است که به‌صورت زیر داریم:

$$f(x) = \begin{cases} \beta_0 + \beta_1 x & x < t_1 \\ \beta'_0 + (\beta_1 + \alpha_1)x & x \geq t_1 \end{cases}$$

با توجه به پیوستگی  $f$  در نقطه  $t_1$  داریم:

$$\beta_0 + \beta_1 t_1 = \beta'_0 + (\beta_1 + \alpha_1)t_1$$

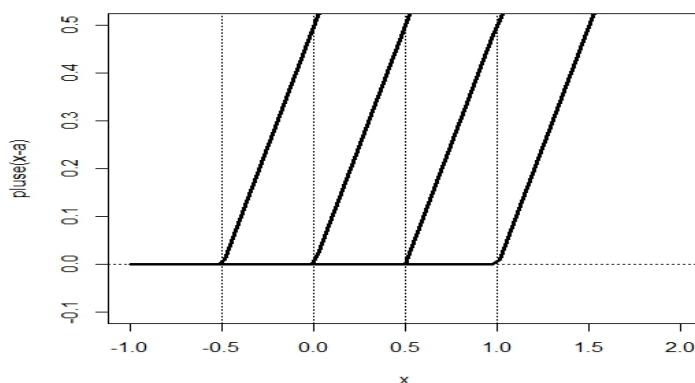
یا به‌طور معادل داریم:

$$\beta'_0 = \beta_0 - \alpha_1 t_1$$

بنابراین معادله کلی  $f$  به‌صورت زیر است.

$$\begin{aligned} f(x) &= (\beta_0 + \beta_1 x)I_{\{x < t_1\}} + (\beta'_0 + (\beta_1 + \alpha_1)x)I_{\{x \geq t_1\}} \\ &= (\beta_0 + \beta_1 x)I_{\{x < t_1\}} + (\beta_0 + \beta_1 x + \alpha_1(x - t_1))I_{\{x \geq t_1\}} \\ &= \beta_0 + \beta_1 x + \alpha_1(x - t_1)I_{\{x \geq t_1\}} \\ &= \beta_0 + \beta_1 x + \alpha_1(x - t_1)_+ \end{aligned}$$

بنابراین تابع  $f$  را می‌توان به صورت ترکیب خطی از توابع پایه ۱ و  $x$  و  $(x - t_1)_+$  نوشت [۱۴]. در شکل ۱.۲ تابع توانی بریده شده از مرتبه یک که در گره‌های  $0, 0.5, 1$  بریده شده است، نشان داده شده‌اند.



شکل ۱.۲: توابع توانی بریده شده مرتبه یک

## ۳.۲ اسپلاین‌ها

اگر برای تحلیل یک مدل رگرسیونی، پذیره خطی بودن رابطه برقرار نباشد و نتوان از یک تابع پارامتری مشخص به سادگی برای مدل استفاده کرد، از روش‌های ناپارامتری رگرسیونی استفاده می‌شود. اسپلاین‌ها به عنوان یکی از ابزارهای درونیایی، از جمله روش‌های ناپارامتری برای مدل بندی رگرسیون ناپارامتری است. مکانیسم عملکرد روش برازش مبتنی بر اسپلاین‌ها به این صورت است که ابتدا کل بازه‌ای که مشاهدات در آن قرار دارند به تعدادی زیربازه تقسیم می‌شود. سپس در هر زیربازه یک منحنی از درجه  $p \geq 1$  به مشاهدات برازش داده می‌شود به طوری که منحنی‌ها در نقاط مرزی بین زیربازه‌ها (یا همان نقاط گره) به هم متصل می‌شوند. فرض کنید تعدادی مشاهده داشته باشیم که در بازه  $[t_0 = \min(x_i), t_{m+1} = \max(x_i)]$  قرار دارند. اگر طول بازه را به فواصل  $[t_0, t_1], [t_1, t_2], \dots, [t_m, t_{m+1}]$  تقسیم کنیم، آن گاه چنانچه منحنی‌های برازش شده  $f$ ، در هر زیربازه به صورت یک چندجمله‌ای از مرتبه  $P$  باشد، یک تابع اسپلاین از مرتبه  $P$  با شرایط زیر داریم:

- تابع  $f$  در گره  $t_j$   $i = 1, 2, \dots, m$  پیوسته است.
- مشتق‌های تابع  $f$  تا مرتبه  $P-1$  وجود دارند و در گره‌ها پیوسته هستند.

ضابطه کلی توابع اسپلاین به صورت زیر است:

$$f(x) = \sum_{k=0}^p \beta_k x^k + \sum_{j=1}^m \beta_{p+j} (x - t_j)_+^p \quad (1.2)$$

که در آن  $\beta = (\beta_0, \dots, \beta_{p+m})$  بردار ضرایب است. می‌توان اسپلاین‌های مرتبه  $p$  با  $t_1, \dots, t_m$  گره، که در هر قطعه یک چندجمله‌ای از مرتبه  $p - 1$  و همچنین دارای مشتقات پیوسته تا مرتبه  $p - 2$  زیر را به صورت

$$f(x) = \sum_{j=1}^m h_j(x) \beta$$

بازنویسی کرد که در آن  $h_j(x) = (1, x, \dots, x^p, (x - t_1)_+^p, \dots, (x - t_m)_+^p)$  توابع پایه هستند. با استفاده از (۱.۲) می‌توان مدل رگرسیون اسپلاینی را به فرم ماتریسی

$$Y = T\beta + \epsilon \quad (2.2)$$

ارائه کرد که در آن  $T = (h_j(t_1), \dots, h_j(t_m))'$  ماتریس طرح  $(n \times (p + m + 1))$  بعدی است. با به کار بردن روش کمترین توان‌های دوم برآوردگر زیر حاصل می‌شود [۳۴، ۸، ۱۴].

$$\hat{\beta} = (T'T)^{-1}T'Y \quad (3.2)$$

یکی از جدی‌ترین معایب اسپلاین‌هایی که از توابع پایه ساخته می‌شوند، ناپایداری عددی است که با افزایش درجه چندجمله‌ای، هم‌خطی بین توابع بریده‌شده ایجاد می‌شود در این حالت معادله‌ی  $T'T\beta = T'Y$  ناجور نامیده می‌شود و  $\hat{\beta}$  ممکن است ناپایدار باشد، بدین مفهوم که تغییرات کوچکی در داده‌ها ممکن است موجب تغییراتی بزرگ در  $\hat{\beta}$  شود. در بخش بعدی نوع دیگری از اسپلاین‌ها را معرفی می‌کنیم که از رابطه‌ی بازگشتی برای ساخت توابع پایه استفاده می‌شود تا ضمن رفع مشکل هم‌خطی، چندجمله‌ای‌های متعامدی را ایجاد می‌کنند [۲]. محققین بسیاری بر روی اسپلاین‌ها، مباحث نظری و کاربردهای آن‌ها پژوهش کرده‌اند. به عنوان چند نمونه می‌توان به هیستی و تیبشیرانی [۱۲]، وند و اورمرد [۳۳] و روپرت و همکارانش [۲۵] اشاره کرد.

## ۴.۲ اسپلاین مکعبی (CS)

اسپلاین مکعبی<sup>۱</sup> یک ابزار قدرتمند برای تجزیه و تحلیل داده‌ها می‌باشد و کاربرد زیادی در هموارسازی دارد. یک اسپلاین درجه سوم با گره‌های  $t_1, t_2, \dots, t_m$  در بازه  $[k_0, k_{m+1}]$  به طوری که  $k_0 < t_1 < \dots < t_m < k_{m+1}$  و مشتقات اول و دوم پیوسته به صورت زیر تعریف می‌گردد:

<sup>۱</sup>Cubic splines

$$f(x) = \sum_{k=0}^3 \beta_k x^k + \sum_{j=0}^m \beta_{3+j} (x - t_j)_+^3$$

قطعه‌های چندجمله‌ای به نحوی برازش می‌یابند که تابع  $f(x)$  و مشتقات اول و دوم آن در هر گره پیوسته باشند در نتیجه روی تمام بازه پیوسته است.

## ۱.۴.۲ درونیابی اسپلاین مکعبی

اگر در هر زیربازه  $[t_i, t_{i+1}]$  برای هر  $i = 0, 1, \dots, m+1$  ،  $s_i(t_i) = f(t_i)$  باشد آن را درونیاب اسپلاین مکعبی می‌نامند و به صورت زیر تعریف می‌شود:

$$S(x) = \begin{cases} s_0(x) & t_0 \leq x \leq t_1 \\ s_1(x) & t_1 \leq x \leq t_2 \\ s_2(x) & t_2 \leq x \leq t_3 \\ \vdots & \vdots \\ s_{m-1}(x) & t_{m-1} \leq x \leq t_m \\ s_m(x) & t_m \leq x \leq t_{m+1} \end{cases}$$

که هر یک از  $s_i$  ها یک چندجمله‌ای درجه سوم است که به صورت

$$s_i(x) = a_i(x - t_i)^3 + b_i(x - t_i)^2 + c_i(x - t_i) + d_i, \quad i = 0, 1, 2, \dots, m+1$$

است. یک روش معمول برای تعیین چندجمله‌ای درجه سوم در هر زیر بازه تعیین ضرایب چندجمله‌ای می‌باشد. ضرایب  $a_i, b_i, c_i, d_i$  مجهولات هر بازه را تشکیل می‌دهند. مشتقات اول و دوم برای  $m+1$  معادله ضروری است و داریم:

$$s'_i(x) = 3a_i(x - t_i)^2 + 2b_i(x - t_i) + c_i \quad (۴.۲)$$

$$s''_i(x) = 6a_i(x - t_i) + 2b_i. \quad (۵.۲)$$

## ۲.۴.۲ خواص اسپلاین‌های مکعبی

اسپلاین‌های مکعبی در شرایط زیر صدق می‌کنند:

- تابع قطعه‌ای  $s(x)$  تمام داده‌ها را در بر می‌گیرد.
- تابع  $s(x)$  در بازه  $[t_0, t_{m+1}]$  پیوسته است، به طوری که

$$s_i(t_i) = s_i(t_i) \quad i = 1, 2, \dots, m$$

• تابع  $s'(x)$  در بازه  $[t_0, t_{m+1}]$  پیوسته است، به طوری که

$$s'_{i-1}(t_i) = s'_i(t_i) \quad i = 1, 2, \dots, m$$

• تابع  $s''(x)$  در بازه  $[t_0, t_{m+1}]$  پیوسته است، به طوری که

$$s''_{i-1}(t_i) = s''_i(t_i) \quad i = 1, 2, \dots, m$$

• شرط‌های زیر به صورت اختیاری هستند:

$$\begin{aligned} s''(t_0) &= s''(t_{m+1}) = 0 \\ s'(t_0) &= f'(t_0) \quad s'(t_{m+1}) = f'(t_{m+1}) \end{aligned}$$

ضرایب چندجمله‌ای در هر زیربازه و چندجمله‌ای اسپلاین مکعبی با حل این دستگاه به دست می‌آیند:

چون تابع  $S(x)$  تمام مشاهدات را درونیابی می‌کند، در نتیجه

$$s_i(t_i) = f(t_i)$$

$$s_i(t_i) = a_i(t_i - t_i)^3 + b_i(t_i - t_i)^2 + c_i(t_i - t_i) + d_i = d_i$$

$$f(t_i) = d_i$$

و با پیوستگی تابع  $S(x)$  در تمام فواصل داریم

$$s_i(t_i) = s_{i-1}(t_i)$$

$$s_{i-1}(t_i) = a_{i-1}(t_i - t_{i-1})^3 + b_{i-1}(t_i - t_{i-1})^2 + c_{i-1}(t_i - t_{i-1}) + d_{i-1}.$$

بنابراین

$$d_i = a_{i-1}(t_i - t_{i-1})^3 + b_{i-1}(t_i - t_{i-1})^2 + c_{i-1}(t_i - t_{i-1}) + d_{i-1}$$

با فرض  $h = (t_i - t_{i-1})$  رابطه‌ی زیر حاصل می‌گردد:

$$d_i = a_{i-1}h^3 + b_{i-1}h^2 + c_{i-1}h + d_{i-1} \quad (۶.۲)$$

با تساوی مشتق اول در نقاط میانی با توجه به رابطه (۴.۲) روابط زیر به دست می‌آید:

$$s'_i(t_i) = c_i$$

$$s'_i(t_i) = s'_{i-1}(t_i)$$

$$s'_i(t_i) = 3a_{i-1}(t_i - t_{i-1})^2 + 2b_{i-1}(t_i - t_{i-1}) + c_{i-1}.$$

در نتیجه

$$c_i = 3a_{i-1}h^2 + 2b_{i-1}h + c_{i-1} \quad (7.2)$$

و با تساوی مشتق دوم در نقاط میانی و با توجه به رابطه‌ی (۵.۲) روابط زیر به دست می‌آید:

$$s''_{i+1}(t_i) = 2b_{i+1},$$

$$s''_{i+1}(t_i) = s''(t_{i+1})$$

$$s''_i(t_{i+1}) = 6a_i(t_{i+1} - t_i) + 2b_i$$

$$2b_{i+1} = 6a_ih + 2b_i.$$

پس با جایگزین کردن  $s''_i(t_i) = M_i$  رابطه‌ی زیر حاصل می‌شود:

$$s''_i(t_i) = M_i$$

$$2b_i = M_i \implies b_i = \frac{M_i}{2}$$

اکنون  $a_i$  را می‌توان به دست آورد:

$$a_i = \frac{2b_{i+1} - 2b_i}{6h} = \frac{2(\frac{M_{i+1}}{2}) - 2(\frac{M_i}{2})}{6h} \quad (8.2)$$

با استفاده از رابطه‌ی (۶.۲) و با جایگزین کردن  $a_i, b_i, d_i$  مقدار  $c_i$  حاصل می‌شود:

$$c_i = \frac{f_{i+1} - f_i}{h} - (\frac{M_{i+1} + 2M_i}{6})h \quad (9.2)$$

اکنون که تمامی ضرایب اسپلاین درجه سه محاسبه شدند به راحتی می‌توان توابع اسپلاین را به دست آورد. ضرایب چندجمله‌ای به دست آمده به صورت زیر است:

$$\begin{cases} a_i = \frac{M_{i+1} - M_i}{6h} \\ b_i = \frac{M_i}{2} \\ c_i = \frac{f_{i+1} - f_i}{h} - (\frac{M_{i+1} + 2M_i}{6})h \\ d_i = f_i \end{cases}$$

با قرار دادن ضرایب به دست آمده در رابطه‌ی  $c_{i+1} = 3a_ih^2 + 2b_ih + c_i$  می‌توان فرم ماتریسی معادلات را به دست آورد:

$$\begin{aligned} 3(\frac{M_{i+1} - M_i}{6h})h^2 + 2(\frac{M_i}{2})h + \frac{f_{i+1} - f_i}{h} - (\frac{M_{i+1} + 2M_i}{6})h &= \frac{f_{i+2} - f_{i+1}}{h} - (\frac{M_{i+2} + 2M_{i+1}}{6})h \\ \frac{h}{6}(M_i + 4M_{i+1} + M_{i+2}) &= \frac{f_i - 2f_{i+1} + f_{i+2}}{h} \\ M_i + 4M_{i+1} + M_{i+2} &= 6(\frac{f_i - 2f_{i+1} + f_{i+2}}{h^2}) \end{aligned}$$

برای  $i = 1, 2, \dots, m-2$  برقرار است.

$$\begin{bmatrix} 1 & 4 & 1 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & 1 & 4 & 1 & \dots & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 4 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 4 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & 1 & 4 & 1 & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 & 4 & 1 \end{bmatrix} \begin{bmatrix} M_1 \\ M_2 \\ M_3 \\ M_4 \\ \vdots \\ M_{m-3} \\ M_{m-2} \\ M_{m-1} \\ M_m \end{bmatrix} = \frac{6}{h^2} \begin{bmatrix} f_1 - 2f_2 + f_3 \\ f_2 - 2f_3 + f_4 \\ f_3 - 2f_4 + f_5 \\ \vdots \\ f_{m-4} - 2f_{m-3} + f_{m-2} \\ f_{m-3} - 2f_{m-2} + f_{m-1} \\ f_{m-2} - 2f_{m-1} + f_m \end{bmatrix}$$

ذکر این نکته لازم است که برای هر بازه‌ی بین دو نقطه‌ی متوالی یک تابع اسپلاین به‌دست می‌آید که با رسم آن‌ها به‌صورت متوالی کل تابع تقریب در کل بازه به‌دست می‌آید [۱۹].

### ۳.۴.۲ اسپلاین مکعبی طبیعی (NCS)

حالت خاصی از اسپلاین‌های مکعبی با اضافه کردن محدودیتی همراه است. اسپلاین‌های مکعبی طبیعی<sup>۲</sup> شامل قاعده‌ای است که باید مشتق دوم در بازه‌های انتهایی برابر صفر باشد.

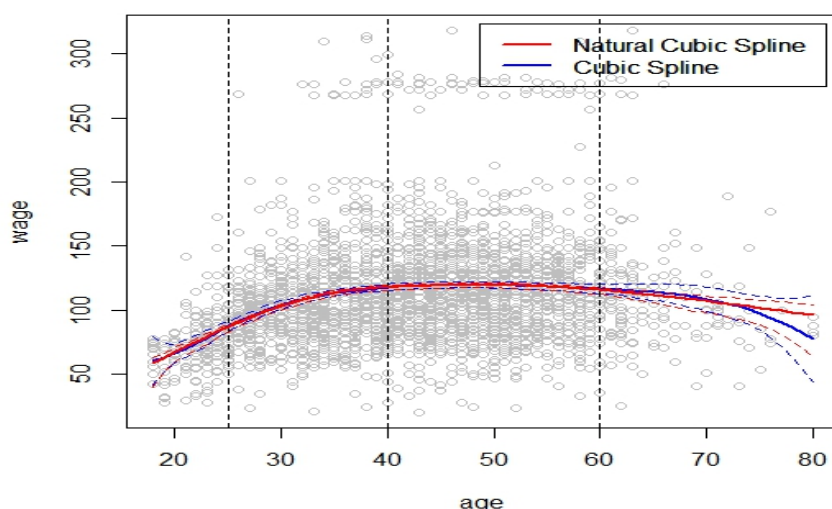
$$M_1 = M_m = 0$$

در نتیجه  $S$  یک اسپلاین مکعبی منحصربه‌فرد است اگر در شرایط مرزی آزاد  $S''(a) = S''(b)$  صدق کند در نتیجه  $a_0 = b_0 = c_0 = d_0 = 0$  و منحنی برازشی در بازه‌های انتهایی به‌صورت خطی است. همچنین فرم ماتریسی آن به‌صورت زیر است:

$$\begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & 1 & 4 & 1 & \dots & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 4 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 4 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & \dots & 1 & 4 & 1 & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ M_2 \\ M_3 \\ M_4 \\ \vdots \\ M_{m-3} \\ M_{m-2} \\ M_{m-1} \\ 0 \end{bmatrix} = \frac{6}{h^2} \begin{bmatrix} f_1 - 2f_2 + f_3 \\ f_2 - 2f_3 + f_4 \\ f_3 - 2f_4 + f_5 \\ \vdots \\ f_{m-4} - 2f_{m-3} + f_{m-2} \\ f_{m-3} - 2f_{m-2} + f_{m-1} \\ f_{m-2} - 2f_{m-1} + f_m \end{bmatrix}$$

<sup>۲</sup>Natural cubic spline

برای سادگی محاسبات می توان ستون اول و آخر را که صفر می باشند حذف نمود [۱۹]. در شکل ۲.۲ مثالی از توابع اسپلاین مکعبی و مکعبی طبیعی رسم شده است. طبق شکل زیر تفاوت زیادی بین اسپلاین مکعبی و مکعبی طبیعی دیده نمی شود، اسپلاین های مکعبی در کران ها انعطاف پذیرتر هستند و اسپلاین های مکعبی طبیعی دارای محدودیت هایی در کران ها می باشند که باعث ایجاد برآوردهای پایدارتری در کران ها می شود و فواصل اطمینان متناظر آن هم محدودتر هستند.



شکل ۲.۲: نمایش اسپلاین های مکعبی و مکعبی طبیعی

## ۵.۲ B-اسپلاین ها

توابع پایه B-اسپلاین، چندجمله ای هایی از مرتبه  $k = p - 1$  هستند که در گره ها پیوستگی دارند و هر کدام از این توابع پایه، نواحی کوچکی از کل مشاهدات را پوشش می دهند. مکانیزم کار در این نوع اسپلاین، مشابه اسپلاین های معمولی است، تنها تفاوت در نحوه ساختن توابع پایه است. اگر  $t = (t_{(k-1)}, \dots, t_{(m+k)})$  یک دنباله غیر نزولی از گره ها باشد، توابع پایه B-اسپلاین به صورت زیر

$$B_j^{k-1}, \quad j = -(k-1), \dots, m$$

تعریف می شود. هم چنین تعداد توابع پایه B-اسپلاینی  $p + m + 1$  است. اولین تابع پایه B-اسپلاینی از مرتبه صفر ( $k = 1$ ) به صورت زیر تعریف می شود:



$$B_j^\circ(x) = \begin{cases} 1 & t_j \leq x \leq t_{j+1} \\ 0 & o.w. \end{cases}$$

دی بور (۱۹۷۸) برای ساخت توابع پایه مراتب بالاتر از توابع پایه مرتبه‌های پایین‌تر استفاده کرد. بر اساس یک رابطه بازگشتی، می‌توان B-اسپلاین‌هایی از هر درجه دلخواه را به صورت زیر محاسبه کرد:

(۱۰.۲)

$$B_j^p(x) = B_j^{k-1}(x) = \frac{x - t_j}{t_{j+k-1} - t_j} B_j^{k-2}(x) + \frac{t_{j+k} - x}{t_{j+k} - t_{j+1}} B_{j+1}^{k-2}(x), j = -(k-1), \dots, m$$

و با توجه به معادله (۱۰.۲)، B-اسپلاین‌های مرتبه یک ( $k=2$ ) را به دست می‌آوریم که بر اساس توابع پایه درجات پایین‌تر تعریف شده است.

$$B_j^1(x) = \frac{x - t_j}{t_{j+1} - t_j} B_j^\circ(x) + \frac{t_{j+2} - x}{t_{j+2} - t_{j+1}} B_{j+1}^\circ(x) \quad (11.2)$$

$$= \begin{cases} \frac{x - t_j}{t_{j+1} - t_j} & t_j \leq x < t_{j+1} \\ \frac{t_{j+2} - x}{t_{j+2} - t_{j+1}} & t_{j+1} \leq x < t_{j+2} \end{cases}$$

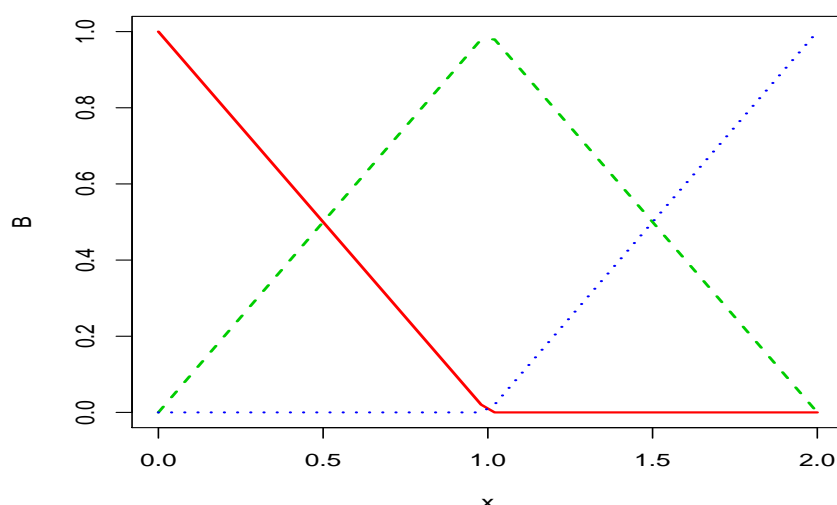
با جایگذاری گره‌ها در رابطه (۱۱.۲) تعدادی معادلات درجه یک به دست می‌آید. به عنوان مثال در بازه‌ی  $[0, 2]$  با در نظر گرفتن یک گره ۱ تعداد توابع پایه B-اسپلاینی  $p+m+1=3$  برای  $j=-1, 0, 1$  به صورت زیر به دست می‌آید:

$$B_{-1}^1(x) = \begin{cases} 1 - x & 0 \leq x < 1 \\ 0 & o.w. \end{cases}$$

$$B_0^1(x) = \begin{cases} x & 0 \leq x < 1 \\ 2 - x & 1 \leq x < 2 \\ 0 & o.w. \end{cases}$$

$$B_1^1(x) = \begin{cases} x - 1 & 1 \leq x < 2 \\ 0 & o.w. \end{cases}$$

نمودار این توابع در شکل ۳.۲ رسم شده‌اند:



شکل ۳.۲: نمایش توابع پایه B-اسپلاینی مرتبه یک

توابع پایه مراتب دو و سه بین محققان محبوب‌تر واقع شدند، دلیل آن هم ویژگی‌ها و انعطاف‌پذیری خوبی است که دارا هستند. برای اطلاعات بیشتر در زمینه B-اسپلاین‌ها به دایرکس [۷] و دی بور [۶] مراجعه کنید.

برخی از ویژگی‌های مهم B-اسپلاین مرتبه  $p = k - 1$  به شرح زیر هستند [۸، ۱۴]:

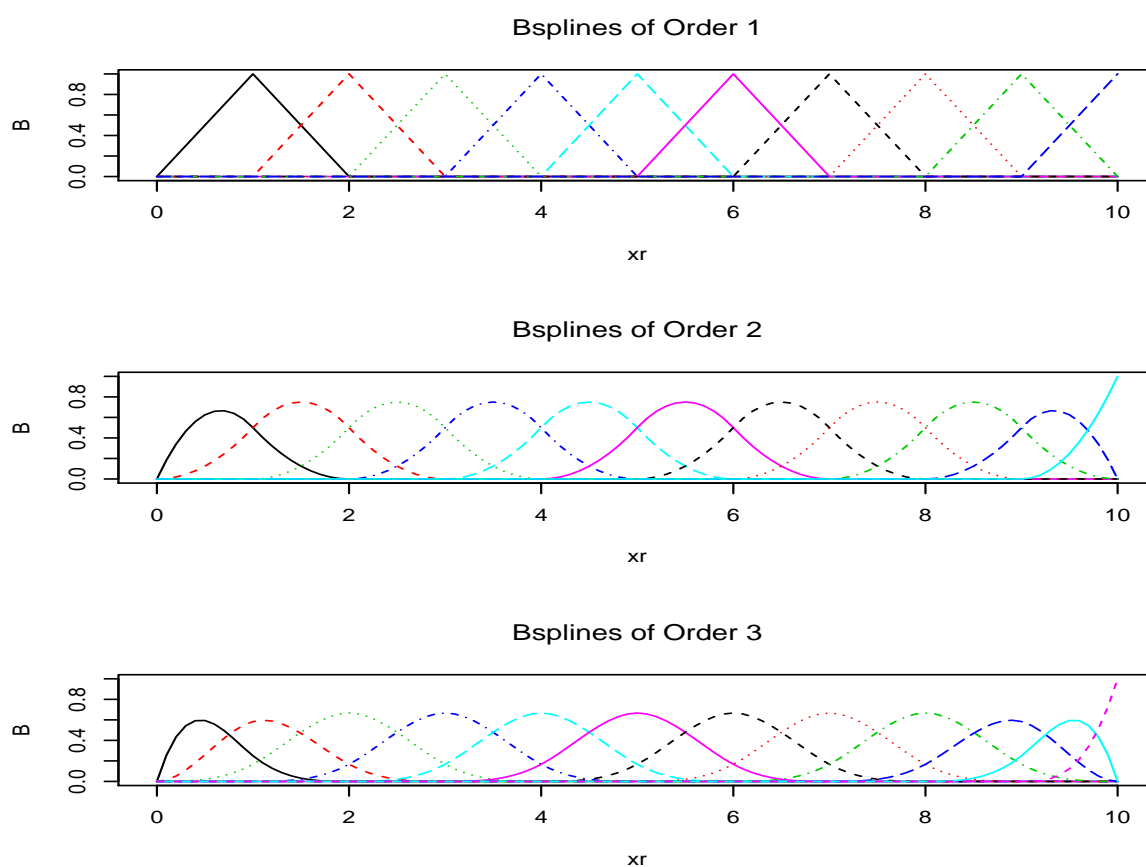
- هر تابع پایه با درجه  $p$  شامل  $p+1$  تکه از درجه  $p$  است.
- تکه‌های چندجمله‌ای مجاور هم در گره‌ها پیوسته هستند.
- در گره‌ها از مرتبه ۱ تا مرتبه  $p-1$  پیوسته و مشتق‌پذیر است.
- توابع پایه B-اسپلاین‌ها هم‌پوشانی دارند، به‌طور مثال اسپلاین‌ها با مرتبه یک با دو همسایگی هم‌چنین B-اسپلاین‌ها با مرتبه دو با چهار همسایگی هم‌پوشانی دارند و در قسمت انتهایی سمت راست و چپ اسپلاین‌ها دارای کمترین هم‌پوشانی است.
- توابع پایه روی تکیه‌گاه‌شان مثبت هستند، یعنی

$$B_j^{k-1}(x) > 0 \quad x \in [t_j, t_{j+m}]$$

- جمع مقادیر این توابع برابر یک است.

$$\sum_{j=-(k-1)}^m B_j^{k-1}(x) = 1$$

در شکل ۴.۲ اسپلاین‌های مرتبه یک، مرتبه دو و مرتبه سه نشان داده شده‌اند.



شکل ۴.۲: نمایش توابع پایه B-اسپلاین

## ۶.۲ هموارسازی اسپلاینی

یکی از کاربردهای اسپلاینها برآورد روند موجود در داده‌ها است. برای این کار می‌توان از هموارسازیهای اسپلاینی استفاده کرد. این نوع چندجمله‌ای‌ها در متون آماری به‌عنوان درونیاب یاد می‌شوند. اگر  $x_0, x_1, \dots, x_n$ ،  $(n+1)$  نقطه دو به دو متمایز و  $f$  تابعی نامعلوم در این نقاط باشد چندجمله‌ای درونیاب  $f$  یک چندجمله‌ای حداکثر از درجه  $n$  است که از نقاط یا گره‌های فوق بگذرد به‌طوری که

$$s(x_k) = f(x_k) \quad k = 0, 1, \dots, n$$

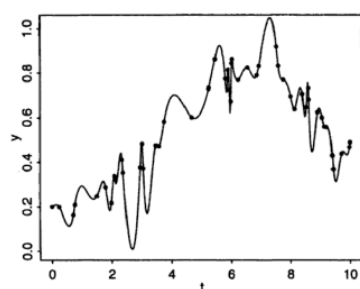
$$e(x_k) = f(x_k) - s(x_k) \quad \text{تابع خطا}$$

ساده‌ترین درونیابی، درونیابی خطی است، نمودار این توابع خط شکسته است که تمام نقاط را به هم وصل می‌کند و در عمل مورد استفاده قرار نمی‌گیرد زیرا خطای آن بسیار زیاد است. اگر تابع درونیاب به‌صورت چندجمله‌ای باشد دارای مزیت‌های بسیاری است زیرا به راحتی می‌توان مشتق و انتگرال این توابع را محاسبه نمود. راه‌های زیادی برای به دست آوردن چندجمله‌ای‌های درونیاب وجود دارد مثل روش لاگرانژ، روش نیوتون،... یا درونیابی به‌صورت گویا که تابع در آن به‌صورت کسری است که در صورت و مخرجش چندجمله‌ای‌ها قرار دارند. درونیابی به روش‌های اخیر دارای دو مشکل بزرگ است:

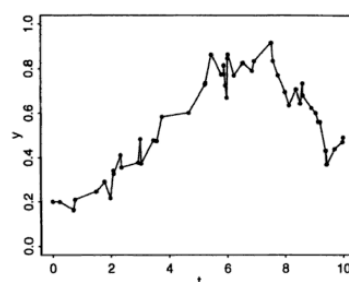
- ممکن است نوسانات تابع درونیاب (به ویژه برای درونیابی چندجمله‌ای‌ها با درجه بالا) خیلی زیاد شود و خطای درونیابی به شدت افزایش یابد.
- اگر یک نقطه به نقاط اولیه اضافه کنیم فرم تابع کلاً تغییر پیدا می‌کند.

استفاده از اسپلاینها بهترین گزینه است زیرا این توابع اطلاعات را به‌طور موضعی درونیابی می‌کنند، به عبارتی بازه را به زیربازه‌هایی تقسیم کرده و تا جای ممکن درجه چندجمله‌ای درونیاب را کاهش می‌دهند، که این رهیافت تقریب قطعه به قطعه با چندجمله‌ای‌ها نامیده می‌شود. در هموارسازی اسپلاینی نوعی از اسپلاین ایجاد می‌شود که علاوه بر این که مجموع مربعات خطا را کمینه کند، باید تابعی هموار باشد یعنی دارای نوسانات سریع نباشد و بتواند از نزدیک داده‌ها عبور کند، نه فقط مشروط بر این که آن‌ها را درونیابی کند بلکه باید شرط مشتق‌پذیری تا مرحله معینی را نیز داشته باشد. بنابراین به برآوردی از یک منحنی دست خواهیم یافت که به‌طور تقریباً همواری از بین توده داده‌ها می‌گذرد که به این عمل هموارسازی اسپلاینی می‌گویند [۱۱، ۱۳].

در شکل ۵.۲ قاب سمت راست، داده‌ها با خط مستقیم به هم متصل شده‌اند که کاملاً شکسته و بریده است و منحنی کاملاً ناهمواری را ایجاد نموده است و در قاب سمت چپ، منحنی با مشتق دوم پیوسته از تمام نقاط عبور کرده و منحنی هموارتری را ایجاد نموده است. شکل



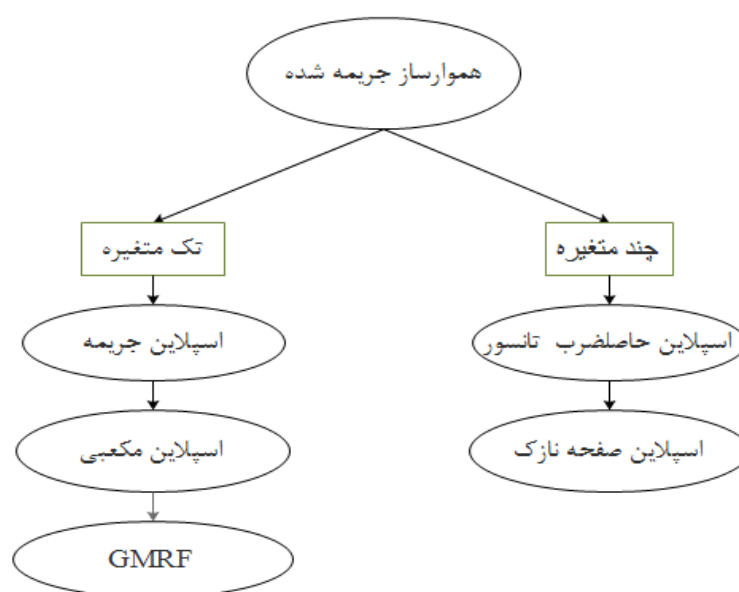
(ب) داده‌ها با منحنی با مشتق دوم پیوسته به هم متصل شده‌اند



(آ) داده‌ها با خط مستقیم به هم متصل شده‌اند

شکل ۵.۲: هموارسازی

۶.۲ روش‌های هموارسازی تک‌متغیره و چندمتغیره را نشان می‌دهد که در بخش‌های بعد این روش‌ها معرفی می‌شوند.



شکل ۶.۲: فلوچارت روش‌های هموارساز

## ۷.۲ اسپلاین‌های جریمه‌ای (p-اسپلاین)

فرض کنید تعداد  $n$  زوج داده داشته باشیم و بخواهیم روند برازشی این مشاهدات را با استفاده از B-اسپلاین‌ها برآورد کنیم. تابع B-اسپلاینی بهتر است که مقدار کمیت زیر را کمینه کند:

$$S = \sum_{i=1}^n \left\{ y_i - \sum_{j=-(k-1)}^m \alpha_j B_j^k(x) \right\}^2.$$

اگر تعداد گره‌ها را نسبتاً زیاد در نظر بگیریم، منحنی برازشی بیشترین تغییرات تعدیل شده توسط داده‌ها را نشان خواهد داد، به عبارتی از تعداد مشاهدات بیشتری عبور خواهد کرد و در نتیجه مقدار مجموع توان‌های دوم خطا ( $S$ ) کمتر خواهد شد. بنابراین تابعی که کمینه می‌شود با تعداد گره‌ها رابطه معکوس دارد. اما اگر تعداد گره‌ها افزایش یابد، منجر به زیربازه‌های بیشتری می‌شود و در نتیجه توابع پایه بیشتری خواهیم داشت و پیرو آن باید پارامترهای بیشتری را برآورد کنیم. از آنجایی که تعداد گره‌ها نامعلوم است، معمولاً تعداد گره‌ها را برابر با تعداد مشاهدات در نظر می‌گیرند که در این صورت ممکن است با بیش‌برازشی مواجه شویم. به عبارتی نوسانات منحنی زیاد می‌شود و منحنی حالت هموار بودن خود را از دست می‌دهد. بنابراین باید بین برازش مناسب و کاهش پیچیدگی مدل (بر حسب تعداد گره‌های کافی) یک تعادل ایجاد شود. روبرت (۲۰۰۲) برای انتخاب تعداد و موقعیت گره‌ها روشی معرفی کرد که در آن تعداد گره‌ها برابر است با  $\min(\frac{N}{4}, 35)$ . برای کنترل و جلوگیری از بیش‌برازشی اسپلوان (۱۹۸۶ و ۱۹۸۸) پیشنهاد داد که جریمه‌ای برای منحنی برازشی اضافه شود. او انتگرال توان دوم مشتق دوم منحنی برازش‌شده را به‌عنوان جریمه هموارسازی به تابع توان دوم خطای برازش به‌صورت زیر افزود:

$$S = \sum_{i=1}^n \left\{ y_i - \sum_{j=-(k-1)}^m \alpha_j B_j^k(x_i) \right\}^2 + \lambda \int_{x_{\min}}^{x_{\max}} \left\{ \sum_{j=-(k-1)}^m \alpha_j B_j''(x) \right\}^2 dx$$

انتگرال مشتق دوم به‌عنوان جریمه هموارساز رایج شد. هیچ تأکید خاصی برای مشتق دوم وجود ندارد در حقیقت می‌توان از مشتقات مراتب بالاتر یا پایین‌تر استفاده نمود در حالی که مشتق اول منجر به معادلات ساده و برازش خطی می‌شود و مشتق مراتب بالاتر منجر به محاسبات طولانی و پیچیده خواهد شد. مشتق تابع B-اسپلاین در پیوست ۲.۴ شرح داده شده است. اسپلاین‌هایی که در برازش آن‌ها از جریمه هموارساز استفاده می‌شود به‌عنوان اسپلاین‌های جریمه‌ای معرفی شدند که اجزای اصلی آن B-اسپلاین‌ها و جریمه ناهماری است. در واقع اسپلاین‌های جریمه‌ای ترکیبی از یک پایه‌ی B-اسپلاینی و یک جمله جریمه روی اختلاف ضرایب جملات مجموعی از توابع اسپلاین است. جمله جریمه بخش اساسی و نقطه قوت در اسپلاین‌های جریمه‌ای است.

ایلرز و مارکس (۱۹۹۶) پیشنهاد دادند که جریمه ناهمواری به صورت تفاضلات متناهی ضرایب مجاور هم قرار داده شوند. یعنی

$$S = \sum_{i=1}^n \left\{ y_i - \sum_{j=-(k-1)}^m \alpha_j B_j^k(x) \right\}^2 + \lambda \sum_{j=-(k-1)}^m (\Delta^k \alpha_j)^2 \quad (12.2)$$

این روش پیشنهادی موجب کاهش حجم محاسبات و همچنین نوشتن معادلات به شکل ماتریسی می‌شود.

## ۱.۷.۲ صورت ماتریسی و برآورد ضرایب اسپلاین جریمه‌ای با استفاده از کمترین توان‌های دوم

فرم ماتریسی جمله جریمه در معادله (۱۲.۲) به صورت زیر است:

$$\sum_{j=-(k-1)}^m (\Delta^k \alpha_j)^2 = \|(\Delta^k \alpha_{-(k-1)} \dots \Delta^k \alpha_m)^T\|^2 = \alpha^T D^T D \alpha$$

$D$  ماتریس جریمه  $k$ -امین مرتبه تفاضلات ضرایب مجاور است. تفاضل مرتبه اول، دوم و به همین ترتیب  $k$ ام به صورت زیر است:

$$\Delta \alpha_j = \alpha_j - \alpha_{j-1}$$

$$\Delta^2 \alpha_j = \Delta(\Delta \alpha_j) = \alpha_j - 2\alpha_{j-1} + \alpha_{j-2}$$

$\vdots$

$$\Delta^k \alpha_j = \Delta(\Delta^{k-1} \alpha_j)$$

تفاضل مرتبه دوم به ازای هر  $j$  را می‌توان به صورت زیر نوشت:

$$j = 2 \Rightarrow \alpha_2 - 2\alpha_1 + \alpha_0 = \Delta^2 \alpha_2$$

$$j = 3 \Rightarrow \alpha_3 - 2\alpha_2 + \alpha_1 = \Delta^2 \alpha_3$$

$\vdots$

در نتیجه ماتریس  $D_{(m+k-2)(m+k)}$  حاصل می‌شود:

$$D = \begin{pmatrix} 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \end{pmatrix}$$

$D^T D$  ماتریسی با بعد  $(m+k) \times (m+k)$  است را محاسبه می‌کنیم:

$$D^T D = \begin{pmatrix} 1 & -2 & 1 & & & & & \\ & 5 & -4 & 1 & & & & \\ 1 & -4 & 6 & -4 & & & & \\ & 1 & -4 & 6 & & & & \\ & & 1 & -4 & & & & \\ & & & 1 & & & & \\ & & & & 1 & & & \\ & & & & & 6 & -4 & 1 \\ & & & & & -4 & 5 & -2 \\ & & & & & 1 & -2 & 1 \end{pmatrix}$$

در نتیجه فرم ماتریسی رابطه‌ی (۱۲.۲) را می‌توان به صورت زیر نوشت:

$$\begin{aligned} S &= (y - B\alpha)^T (y - B\alpha) + \lambda(\alpha^T D^T D \alpha) \\ &= y^T y - 2y^T B\alpha + \alpha^T B^T B\alpha + \lambda(\alpha^T D^T D \alpha) \end{aligned} \quad (13.2)$$

برای مینیم کردن عبارت بالا نسبت به  $\alpha$  مشتق می‌گیریم در نتیجه داریم:

$$\begin{aligned} \frac{\partial S}{\partial \alpha} &= -2\beta^T y + 2\beta^T \beta \hat{\alpha} + 2\lambda D^T D \hat{\alpha} = 0 \\ \beta^T y &= \hat{\alpha}(\beta^T \beta + \lambda D^T D) \end{aligned}$$

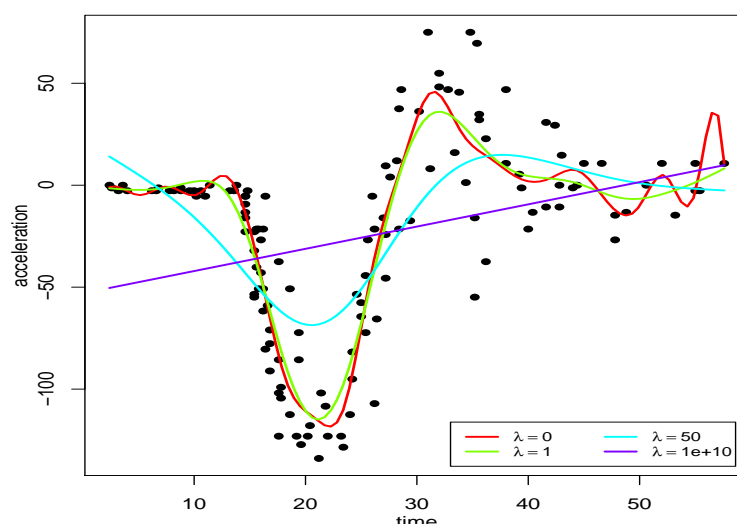
در نتیجه ضرایب برآورد شده به صورت

$$\hat{\alpha} = (\beta^T y)(\beta^T \beta + \lambda D^T D)^{-1} \quad (14.2)$$

است.

پارامتر  $\lambda$  نیز پارامتر هموارساز نامیده می‌شود که مقداری مثبت است و برای کنترل پیوستگی منحنی مورد استفاده قرار می‌گیرد. اگر مقدار  $\lambda = 0$  باشد، منحنی برازش شده می‌تواند تابعی باشد که داده‌ها را به طور دقیق درونیابی می‌کند به عبارتی از روی تمام نقاط عبور می‌کند و مجموع مربعات خطا سهم اصلی را در  $S$  دارد، و اگر مقدار  $\lambda = \infty$  در نظر گرفته شود، جمله جریمه به صفر میل می‌کند و منحنی برازشی یک خط رگرسیون خواهد بود. در شکل ۷.۲ کاملاً مشهود می‌باشد.





شکل ۷.۲: هموارسازی تابع رگرسیونی برای مقادیر مختلف  $\lambda$

## ۲.۷.۲ خواص اسپلاین جریمه‌ای

اسپلاین‌های جریمه‌ای دارای خواص خوبی هستند که تا حدودی از B-اسپلاین‌ها به ارث برده‌اند، برخی از خواص آن‌ها عبارتند از [۸]:

- اسپلاین‌های جریمه‌ای چندجمله‌ای‌هایی را به‌طور دقیق به داده‌ها برازش می‌دهند. فرض کنید  $(x_i, y_i)$  یک مجموعه داده باشد اگر  $y_i$  تابع چندجمله‌ای از  $x$  از مرتبه  $k$  باشد آنگاه B-اسپلاین‌ها از درجه  $k$  یا بالاتر به‌طور دقیق به داده‌ها برازش داده می‌شوند [۵]. اگر مرتبه جریمه  $(k+1)$  یا دارای مرتبه‌های بالاتری باشد، صرف‌نظر از مقادیر مختلف  $\lambda$  (پارامتر هموارسازی)، برای اسپلاین‌های جریمه‌ای برقرار می‌باشد. برای اثبات این مطلب به ایلرز و مارکس ۱۹۹۶ مراجعه شود.

- اسپلاین‌های جریمه‌ای گشتاور داده‌ها را حفظ می‌کنند. برای یک مدل خطی اسپلاین‌های توانیده از مرتبه  $k+1$  و مرتبه جریمه  $k+1$  یا بالاتر رابطه زیر برقرار است:

$$\sum_{i=1}^m x^k y_i = \sum_{i=1}^m x^k \hat{y}_i$$

برای همه مقادیر  $\lambda$  (پارامتر هموارسازی)،  $\hat{y}_i = \sum_{j=1}^n b_{ij} \hat{\alpha}_j$  مقادیر برازش داده‌شده هستند و برای مدل‌های خطی تعمیم‌یافته به‌صورت زیر به‌دست می‌آیند:

$$\sum_{i=1}^m x^k y_i = \sum_{i=1}^m x^k \hat{\mu}_i.$$

این ویژگی خصوصا در هموارسازی تابع چگالی، مفید است زیرا میانگین و واریانس تابع چگالی برآورد شده، با میانگین و واریانس داده‌ها برابر است.

## ۸.۲ اسپلاین حاصل ضرب تانسور

اسپلاین حاصل ضرب تانسوری<sup>۳</sup> سالیان متمادی برای تقریب تک‌متغیره و روش‌های درونیایی دومتغیره مورد استفاده قرار گرفته‌اند و تعمیمی ساده و مستقیم از اسپلاین‌های  $n$  بعدی را فراهم می‌کنند. تقریب چندمتغیری از توابع توسط حاصل ضرب تانسوری دارای جهت‌گیری قوی در امتداد خطوط موازی با جهت محور می‌باشد. رویه‌های حاصل ضرب تانسوری به‌شدت به توابع اسپلاینی وابسته می‌باشند. در این روش می‌توان برای ساختار توابع هموارگر از هر تعداد متغیر استفاده نمود. ساده‌ترین روش برای ساخت یک تابع هموارگر استفاده از سه متغیر  $x$ ،  $z$  و  $\nu$  است که هر یک از متغیرها دارای تابع هموارگر  $f_x$ ،  $f_z$  و  $f_\nu$  با توابع پایه می‌باشند. که به‌صورت زیر قابل نمایش است:

$$f_x(X) = \sum_{i=1}^I \alpha_i a_i(x) \quad f_z(Z) = \sum_{l=1}^L \delta_l d_l(z) \quad f_\nu(\nu) = \sum_{k=1}^K \beta_k b_k(\nu)$$

که در آن  $\alpha_i$ ،  $\delta_l$  و  $\beta_k$  پارامترها و  $a_i(x)$ ،  $d_l(z)$  و  $b_k(\nu)$  توابع پایه معلوم می‌باشند و  $L$ ،  $I$  و  $K$  می‌توانند بازه‌ای نسبتاً بزرگ اختیار نمایند. تابع هموارگر  $f_x(X)$  می‌تواند به یک تابع هموارگر  $f(x, z)$  تبدیل شود، برای این کار لازم است که  $f(x, z)$  به‌طور مساوی با  $z$  تغییر کند در نتیجه باید پارامتر  $\alpha_i$  به‌طور مساوی با  $z$  تغییر کند و با به‌کار بردن توابع پایه تابع هموارگر  $z$  داریم:

$$\alpha_i(z) = \sum_{l=1}^L \delta_{il} d_l(z)$$

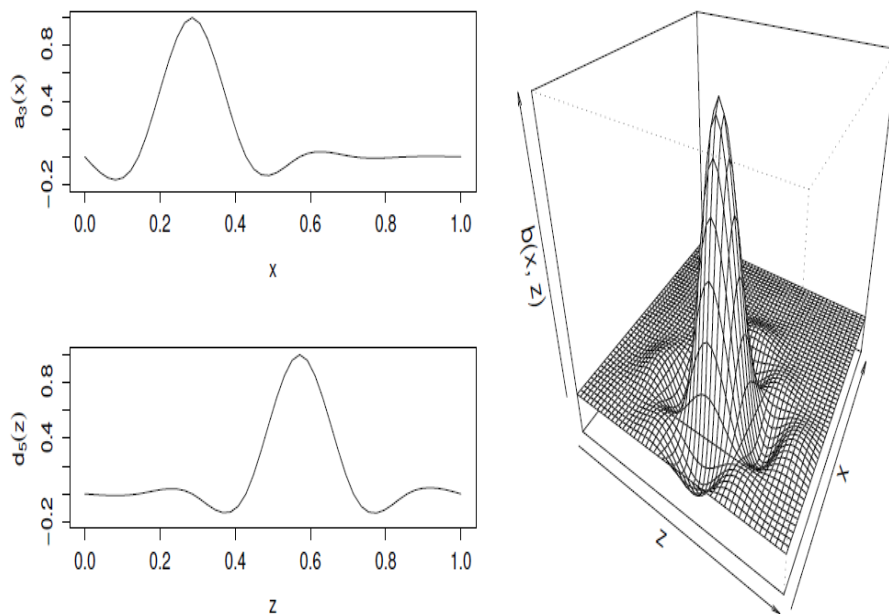
بنابراین ساختار پایه‌های حاصل ضرب تانسوری برای یک سطح به‌صورت زیر حاصل می‌شود:

$$f_{xz}(x, z) = \sum_{i=1}^I \alpha_i(z) a_i(x) = \sum_{i=1}^I \sum_{l=1}^L \delta_{il} d_l(z) a_i(x) \quad (۱۵.۲)$$

شکل ۸.۲ ساختار معادله (۱۵.۲) را نشان می‌دهد. با ادامه این روند می‌توان تابع هموارگر  $x$ ،  $z$  و  $\nu$  را با تغییر دادن  $f_{xz}$  به‌طور یکسان با  $\nu$  به‌دست آورد. پس با همان استدلال قبل داریم [۳۶، ۳۴]:

$$f_{x\nu}(x, z, \nu) = \sum_{i=1}^I \sum_{l=1}^L \sum_{k=1}^K \beta_{ilk} b_k(\nu) d_l(z) a_i(x)$$

<sup>۳</sup> Tensor product splines



شکل ۸.۲: حاصل ضرب دو تابع پایه‌ی حاشیه‌ای برای توابع هموار  $f_x$  و  $f_z$

این ساختار می‌تواند برای تعداد متغیرهای بیشتری تعمیم داده شود. می‌توان میزان انعطاف‌پذیری این توابع را اندازه‌گیری نمود. فرض کنید که هر هموارگر حاشیه‌ای دارای تابعی مرتبط است که میزان انعطاف‌پذیری را اندازه‌گیری می‌کند و می‌تواند به فرم درجه دوم از پارامترهای حاشیه‌ای بیان شود که به صورت زیر می‌باشد:

$$J_x(f_x) = \int f''(x)^2 dx = \alpha^T S_x \alpha$$

$$J_z(f_z) = \int f''(z)^2 dz = \delta^T S_z \delta$$

$$J_\nu(f_\nu) = \int f''(\nu)^2 d\nu = \beta^T S_\nu \beta$$

$S$  ماتریسی از ضرایب معلوم و  $\alpha, \delta, \beta$  بردار ضرایب از هموارگر حاشیه‌ای است. اکنون در نظر بگیرید  $f_{x|z}(x)$  که به  $f_{xz}(x, z)$  به عنوان تابعی از  $x$  با  $z$  ثابت تبدیل شده است. در نتیجه  $J_x(f_{x|z})$  میزان انعطاف‌پذیری  $f_{x|z}$  را اندازه‌گیری می‌کند و  $\int J_x(f_{x|z}) dz$  با میانگین انعطاف‌پذیری در جهت  $x$  متناسب است. بنابراین یک جریمه مناسب به صورت معادله زیر می‌باشد:

$$J_x(f_{xz}) = \lambda_x \int_z J_x(f_{x|z}) dz + \lambda_z \int_x J_z(f_{z|x}) dx$$

به طور مشابه می‌توان برای توابع  $f_{z|x\nu}(z)$  و  $f_{\nu|xz}(\nu)$  تعریف نمود. بنابراین یک روش برای اندازه‌گیری انعطاف‌پذیری تابع  $f_{xz\nu}$  به صورت معادله زیر می‌باشد.

$$J(f_{xz\nu}) = \lambda_x \int_{z\nu} J_x(f_{x|z\nu}) dz d\nu + \lambda_z \int_{x\nu} J_z(f_{z|x\nu}) dx d\nu + \lambda_\nu \int_{xz} J_\nu(f_{\nu|xz}) dx dz$$

که در آن  $\lambda$  پارامتر هموارسازی است که انعطاف‌پذیری تابع را در جهت‌های مختلف کنترل می‌کند. به‌عنوان مثال با در نظر گرفتن تابع جریمه در جهت  $x$  تابع  $f_{x|z\nu}(x)$  به‌صورت زیر بیان می‌شود:

$$f_{x|z\nu}(x) = \sum_{i=1}^I \alpha_i(z, \nu) a_i(x).$$

می‌توان ماتریس ضرایب  $M_{z,\nu}$  را تعریف نمود به‌طوری که  $\alpha(z, \nu) = M_{z\nu}\beta$  که  $\beta$  بردار پارامتر برای  $f(x, z, \nu)$  است. از این رو

$$J_x(f_{x|z,\nu}) = \alpha(z, \nu)^T S_x \alpha(z, \nu) = \beta^T M_{z\nu}^T S_x M_{z\nu} \beta.$$

بنابراین

$$\int_{z,\nu} J_x(f_{x|z,\nu}) dz d\nu = \beta^T \int_{z,\nu} M_{z\nu}^T S_x M_{z\nu} dz d\nu \beta.$$

انتگرال فوق را می‌توان محاسبه نمود اما بسیار زمانبر است. به خصوص زمانی که تعداد متغیرها افزایش یابد. و این رویکرد برای تمام اجزای جریمه اعمال می‌شود. با این حال، با به‌کارگیری بازپارامتریدن ساده می‌توان تقریبی برای شرایط جریمه ارائه داد که به‌خوبی عمل می‌کند و از محاسبه انتگرال عددی جلوگیری می‌کند. برای نشان دادن عملکرد این رویکرد تابع حاشیه‌ای  $f_x$  را در نظر بگیرید و قرار دهید  $\{x_i^* : i = 1, \dots, I\}$ ، مجموعه‌ای از مقادیر  $x$  که به‌طور یکنواخت در دامنه مقادیر مشاهده شده  $x$  پخش شده‌اند. با بازنویسی پارامترهای تابع حاشیه‌ای  $f_x$  داریم:

$$\alpha_i'' = f_x(x_i^*).$$

تحت شرایط بازپارامتری کردن  $\alpha' = \Gamma \alpha$  و  $\Gamma_{ij} = a_i(x_j^*)$ . بنابراین ماتریس مدل حاشیه‌ای به‌صورت  $X'_x = X_x L$  و ماتریس ضرایب جریمه به  $S'_x = l^{-T} S_x l^{-1}$  تبدیل می‌شود و داریم:

$$\int_{z\nu} J_x(f_{x|z\nu}) dz d\nu \approx h \sum_{lk} J_x(f_{x|z_l^* \nu_k^*}),$$

که  $h$  مقداری ثابت از نسبت فاصله‌ی بین  $z_l^*$  و  $\nu_k^*$  است. نشان دادن این تقریب در مجموع بالا سرراست است و داریم [۳۶، ۳۴]:

$$J_x^*(f_{xz\nu}) = \beta^T \tilde{S}_x \beta \quad \tilde{S}_x = S'_x \otimes I_L \otimes I_K.$$

با تعاریف مشابه برای دیگر اجزای جریمه معادلات زیر به‌دست آمده است:

$$J_z^*(f_{xz\nu}) = \beta^T \tilde{S}_z \beta \quad \tilde{S}_z = I_I \otimes S'_z \otimes I_K$$

9

$$J_\nu^*(f_{xz\nu}) = \beta^T \tilde{S}_\nu \beta \quad \tilde{S}_\nu = I_I \otimes I_L \otimes S'_\nu.$$

نماد  $\otimes$  نمایانگر ضرب کرونکر است که در پیوست ۱.۴ شرح داده شده است. بنابراین

$$J(f_{xz\nu}) \approx J^*(f_{xz\nu}) = \lambda_x J^*(f_{xz\nu}) + \lambda_z J^*(f_{xz\nu}) + \lambda_\nu J^*(f_{xz\nu})$$

## ۹.۲ اسپلاین صفحه نازک

اسپلاین صفحه نازک<sup>۴</sup> یک نوع از هموارسازی اسپلاینی است که برای تجسم روابط پیچیده بین متغیرهای پاسخ و پیش‌گوی پیوسته به کار می‌رود و به دلیل ظاهر چند بعدی برای بررسی اثر ترکیبی دو متغیر پیش‌گوی پیوسته در یک نتیجه واحد قابل استفاده است. اسپلاین صفحه نازک تکنیکی مبتنی بر اسپلاین برای درونیابی و هموارسازی مجموعه‌ای از نقاط کنترل است. در این روش یک سطح از میان نقاط کنترل عبور می‌کند. نام اسپلاین صفحه نازک اشاره‌ای به حالت فیزیکی صفحه نازک فلزی که قابل خم شدن است را دارد. داچون و ماینگت پایه و اساس اسپلاین صفحه نازک را طراحی کردند که سطحی قابل انعطاف است [۱۰، ۳۵].

اسپلاین صفحه نازک برای اولین بار در سال ۱۹۷۰ در علوم کامپیوتر در زمینه محاسبات هندسی مورد استفاده قرار گرفتند. این رویکرد مدل‌سازی امروزه در بسیاری از زمینه‌ها از جمله مهندسی (طراحی ساختاری)، اکولوژی (رشد جمعیت) و شناخت الگو (شناسایی اثر انگشت) کاربرد دارد. استفاده رایج آن در تحقیقات علمی در زمینه سلامت، شامل توسعه و تجزیه تحلیل تکنیک‌های تصویربرداری پزشکی است [۲۰].

با در نظر گرفتن  $n$  نقطه متمایز و متغیرهای ورودی  $\{(x_i, y_i), i = 1, \dots, n\}$  و متغیر پاسخ  $z_i$  مطلوب‌ترین برازش، تعمیمی از اسپلاین‌هایی است که انحنای سطح برازشی را کمینه کند. در این حالت سطح  $f(x, y)$  طوری تعیین شود که فاصله آن با نقاط کنترل کمینه شود، به عبارتی رابطه‌ی زیر مینماید.

$$E_{Tps} = \sum_{i=1}^n (z_i - f(x_i, y_i))^2 + \lambda j[f(x, y)] \quad (۱۶.۲)$$

که  $\lambda$  پارامتر هموارساز است هنگامی که  $\lambda = 0$  باشد فقط داده‌ها بدون هموارسازی درونیابی می‌شوند و وقتی  $\lambda \rightarrow \infty$  مساله یافتن سطحی است که فاصله آن از نقاط کنترل کمینه گردد. تابع جریمه را می‌توان با معیارهای مختلفی از جمله معیار اعتبارسنجی متقابل انتخاب نمود که فرم کلی آن به صورت زیر است:

$$j[f(x, y)] = \iint_{\mathbb{R}^2} \left[ \left( \frac{\partial^2 f}{\partial x^2} \right)^2 + 2 \left( \frac{\partial^2 f}{\partial x \partial y} \right)^2 + \left( \frac{\partial^2 f}{\partial y^2} \right)^2 \right] dx dy \quad (۱۷.۲)$$

$f(x, y)$  یک نوع از توابع پایه شعاعی<sup>۵</sup> است. توابع پایه شعاعی برای درونیابی استفاده می‌شوند،

<sup>۴</sup>Thin plate spline

به صورت معادله‌ی زیر می‌باشد.

$$f(x, y) = a_0 + a_1 x + a_2 y + \sum_{i=1}^n W_i \varphi(\|x - x_i\|). \quad (18.2)$$

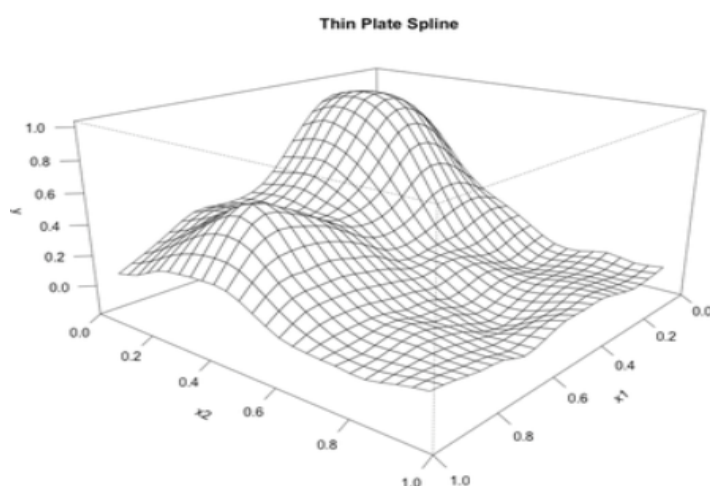
سه جمله اول مربوط به قسمت خطی است. صفحه‌ای تخت که به بهترین نحو بر نقاط کنترل منطبق باشد، تعریف می‌شود و عبارت آخر هم مربوط به نیروی خمشی ناشی از نقاط کنترل است.  $\|\cdot\|$  بیانگر فاصله اقلیدسی و  $w_i$  مجموعه‌ای از ضرایب نگاشت که هر کدام متعلق به یکی از مراکز داده‌ها است،  $x_i, i = 1, \dots, n$  مجموعه نقاط کنترل و  $\varphi$  تابع شعاعی کرنل که به صورت زیر تعریف می‌شود:

$$\varphi(r) = r^2 \log(r) \quad (19.2)$$

$r$  فاصله اقلیدسی در فضای دو بعدی است و  $\varphi(r)$  وابسته به فاصله اقلیدسی از مبدا  $x_i$  است. برازش اسپلاین صفحه نازک با مدل جمعی تعمیم‌یافته (GAM) به صورت زیر است:

$$g(E(Y)) = \beta_0 + f(x) + \epsilon.$$

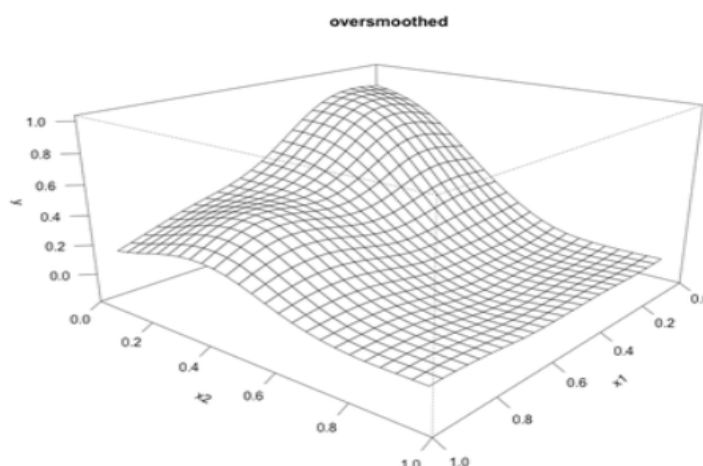
$\beta_0$  مقداری ثابت است،  $f(x)$  تابعی انعطاف پذیر از  $x$  و  $\epsilon$  عبارت خطا است و تابعی هموار بر اساس روش حداقل مربعات جریمه‌ای را فراهم می‌کند که با افزایش  $\epsilon$  تابع هموارتری نسبت به توابع اسپلاین به دست می‌آید [۲۰، ۳۵]. شکل ۹.۲ نمایش اسپلاین صفحه نازک دو بعدی است که هر متغیر بر روی یک محور ترسیم شده است و یک سطح دو بعدی در سه بعد قابل نمایش است.



شکل ۹.۲: نمایش اسپلاین صفحه نازک

عبارت خطا را می‌توان به عنوان کشش معرفی کرد یا میزان فشاری که برای خمیدگی یک ورق نازک فلزی لازم است. با کشش بالاتر مقاومت صفحات نازک برای خمیدگی بیشتر می‌شود و در نتیجه تاثیر متغیرهای پیش‌گو بر متغیر پاسخ افزایش یافته و در پی آن تابع هموارتری ظاهر خواهد شد.

در شکل ۱۰.۲ با افزایش  $\epsilon$  تابع هموارتری نسبت به شکل ۹.۲ حاصل شده است.

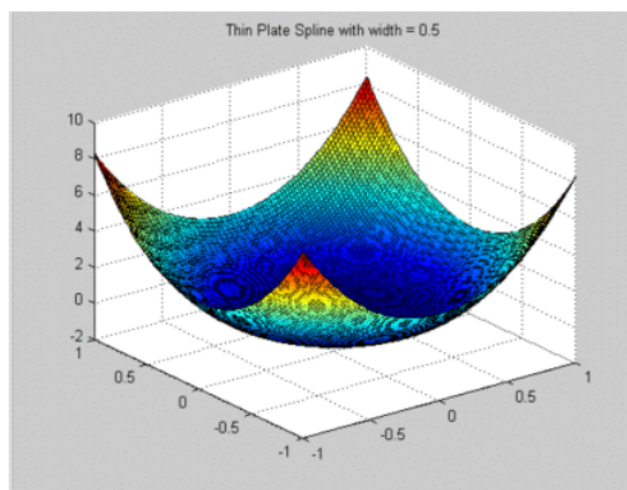


شکل ۱۰.۲: نمایش اسپلاین صفحه نازک هموارتر

## ۱.۹.۲ مزایای اسپلاین صفحه نازک

مزایای اسپلاین صفحه نازک مانند سایر اسپلاین‌های هموارساز است که در آن مدل‌های جمعی تعمیم‌یافته (GAM) نیازی به هیچ گونه اطلاعات پیشینی از فرم عملکرد داده‌ها یا روابط بین آن‌ها ندارند. تعیین تعداد گره و موقعیت آن‌ها چالشی در اسپلاین‌ها است که به‌طور موثر به‌عنوان بخشی از توابع هموارساز اسپلاین صفحه نازک است. ماهیت اسپلاین صفحه نازک آن‌ها را به یک ابزار قدرتمند و جذاب برای تجسم روابط پیچیده بین متغیرهای پیش‌گو و پاسخ تبدیل کرده است و با افزودن نقشه‌های گرمایی<sup>۶</sup> به اسپلاین صفحه نازک که بیشتر بر انحنای رویه تاکید دارند آن را به یک گرافیک بصری جذاب تبدیل می‌کنند، مانند شکل ۱۱.۲ [۱۱، ۱۶].

<sup>۶</sup>Heat maps



شکل ۱۱.۲: نمایش اسپلاین صفحه نازک





## فصل ۳

# برخی از مباحث نظری اسپلاین‌ها

در این فصل نحوه‌ی انتخاب پارامتر هموارسازی، روش‌های انتخاب تعداد و موقعیت گره‌ها و انتخاب بهترین مدل بر اساس معیارهای  $GCV, CV, AIC$ ، ... شرح داده شده است.

### ۱.۳ نحوه‌ی انتخاب پارامتر هموارسازی

پارامتر هموارسازی برای کنترل همواری منحنی استفاده می‌شود. هیستی و تیبشیرانی (۱۹۹۰) نماد  $\lambda$  را برای نشان دادن پارامتر همواری به کار بردند. در هر حال  $\lambda$  باید به‌طور محافظه‌کارانه انتخاب شود تا مسئله بیش‌همواری یا کم‌همواری پیش نیاید و یک تابع هموار مناسب را نتیجه دهد. برای انتخاب بهینه‌ی پارامتر همواری روش‌های مختلفی وجود دارد. یک روش، انتخاب یک مقدار دلخواه برای پارامتر همواری است که با تغییر این مقدار و با توجه به نوع موضوع مورد بررسی یک مقدار برای آن انتخاب می‌گردد به‌طوری که برآورد حاصل، بیشترین تطبیق را با داده‌ها داشته باشد.

از دیگر روش‌ها برای انتخاب بهینه پارامتر هموارسازی استفاده از معیار اطلاع آکائیک<sup>۱</sup>، روش اعتبارسنجی متقابل<sup>۲</sup> و اعتبارسنجی متقابل تعمیم‌یافته<sup>۳</sup> که در ذیل به تفصیل آن‌ها می‌پردازیم.

---

<sup>۱</sup> Akaike information criterion

<sup>۲</sup> Cross validation

<sup>۳</sup> Generalized cross-validation

## ۱.۱.۳ معیار اطلاع آکائیک

ایده‌ی اصلی معیار اطلاع آکائیک به این صورت است که لگاریتم تابع درست‌نمایی مدل برازشی را برای تعداد پارامترهای مؤثر بدست می‌آورد. در واقع این معیار تعادلی بین پیچیدگی مدل و دقت آن را برقرار می‌کند. زمانی که برای مقایسه مدل‌های برازشی مختلف از این معیار استفاده می‌شود نشان می‌دهد که استفاده از یک مدل آماری به چه میزان باعث از دست رفتن اطلاعات می‌شود. بنابراین مدلی که کمترین مقدار آکائیک را داشته باشد به‌عنوان بهترین مدل انتخاب می‌شود. فرم کلی معیار آکائیک به‌صورت زیر است:

$$AIC = \log \left( \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right) + \frac{2 \times \text{tr}(H)}{n}$$

$H$  ماتریس هت نامیده می‌شود.

ایلرز و مارکس (۱۹۹۶) برای به‌دست آوردن معیار آکائیک از انحراف به‌جای لگاریتم درست‌نمایی استفاده نمودند و معیار اطلاع آکائیک را به‌صورت زیر تعریف کردند:

$$AIC(\lambda) = dev(y, \alpha, \lambda) + 2 \times \dim(\alpha, \lambda)$$

در آن  $\dim(\alpha, \lambda)$  بعد (مؤثر) بردار پارامترها است و  $dev(y, \alpha, \lambda)$  میزان انحراف مدل برازشی نسبت به داده‌ها است. محاسبه انحراف سراسر است اما تعیین بعد مؤثر پارامترها در برازش مدل‌های اسپلاین جریمه‌ای پیچیده و مشکل است. هستای و تیشیرانی (۱۹۹۰) برای این مساله روشی را پیشنهاد کردند به این صورت که اثر ماتریس هموارگر خطی را به‌عنوان یک تقریب برای ابعاد مؤثر مدل به‌کار بردند.

$$\dim(\alpha) = \text{tr}(H)$$

مدل اسپلاین جریمه‌ای را در نظر بگیرید:

$$Y = \beta\alpha + \epsilon$$

با استفاده از

$$\hat{\alpha} = (\beta'Y)(\beta'\beta + \lambda D'D)^{-1}$$

و مقادیر برازش داده‌شده به‌صورت زیر محاسبه می‌شوند:

$$\hat{Y} = \beta\hat{\alpha} = \underbrace{\beta(\beta^T W \beta + \lambda D^T D)^{-1} \beta^T W}_{H} y = Hy$$

که ماتریس  $H$  را ماتریس  $\hat{H}$  می‌نامند و  $W$  ماتریس قطری از وزن‌ها است که به‌صورت زیر به‌دست می‌آید:

$$W_{ii} = \frac{1}{\nu_i} \left( \frac{\partial y_i}{\partial \mu} \right)$$

$\nu_i$  واریانس  $y_i$  است. همچنین  $tr(H) = \sum_{i=1}^n h_{ii}$  می‌توان اثر ماتریس را بدون در نظر گرفتن قطر اصلی محاسبه نمود که در زیر شرح داده شده است. فرض کنید  $n$  تا زوج داده به صورت  $(x_i, y_i)$  داریم اگر  $y_i$  یک چندجمله‌ای از درجه  $k$  در  $x$  باشد، آن‌گاه با افزایش  $\lambda$  اثر ماتریس  $H$  به  $k$  نزدیک خواهد شد اثبات به صورت زیر برقرار است. ابتدا یادآوری می‌شود که برای ماتریس‌های ضرب پذیر (سازگار) داریم:

$$tr(AB) = tr(BA).$$

فرض کنید

$$Q_\lambda = \lambda D^T D \quad Q_\beta = \beta^T W \beta$$

$$\begin{aligned} \Rightarrow tr(H) &= tr\{(Q_\beta + Q_\lambda)^{-1} Q_\beta\} \\ &= tr\{Q_\beta^{-1/2} (Q_\beta + Q_\lambda)^{-1} Q_\beta^{1/2}\} \\ &= tr\{(I + Q_\beta^{-1/2} Q_\lambda Q_\beta^{-1/2})^{-1}\} \end{aligned}$$

اگر

$$L = Q_\beta^{-1/2} Q_\lambda Q_\beta^{-1/2}$$

آن‌گاه

$$tr(H) = tr\{(I + \lambda L)^{-1}\} = \sum_{j=1}^n \frac{1}{1 + \lambda \gamma_j} \quad \gamma_j \quad (1.3)$$

$\gamma_j$  مقادیر ویژه  $L$  هستند. در رابطه (۱.۳) اگر مقدار  $\lambda$  بزرگ باشد کسر مورد نظر صفر می‌شود و فقط جملاتی صفر نمی‌شوند که  $\gamma_j$  مربوط به آن‌ها برابر با صفر باشد در نتیجه کسر مربوط به آن‌ها یک می‌شود و در نتیجه مجموع برابر با  $k$  می‌شود.

## ۲.۱.۳ اعتبار سنجی متقابل

یکی دیگر از روش‌های انتخاب پارامتر همواری استفاده از معیار اعتبارسنجی متقابل است که به اختصار با  $CV$  نشان داده می‌شود. روشی ساده و بصری برای برآورد پیش‌گویی یک مدل است. در این روش مبنای انتخاب  $\lambda$  براساس مینیمم کردن  $CV(\lambda)$  است. در این روش داده‌ها را به دو قسمت نمونه آزمون و نمونه آموزشی (مجموعه داده‌هایی که برای برازش مدل استفاده می‌شود) تقسیم می‌شوند. فرض کنید  $n$  تا مشاهده  $y_1, \dots, y_n$

داریم ابتدا مشاهده  $y_i$  را حذف کرده و مدل را با  $n - 1$  مشاهده‌ی باقیمانده برازش داده و سپس خطای پیش‌گویی  $MSE_i = E(y_i - \hat{y}_i)^2$  محاسبه می‌کنیم. این فرایند را تا زمانی که هر مشاهده یک بار خارج شود انجام می‌دهیم و با ادامه این روند برای سایر مشاهدات و میانگین گرفتن از آن‌ها برآورد خطای پیش‌گویی به صورت معادله‌ی زیر حاصل می‌شود.

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i. \quad (2.3)$$

به‌طور کلی  $CV$  دارای پتانسیل بالایی برای پیاده‌سازی می‌باشد زیرا مدل  $n$  بار برازش داده می‌شود اما اگر  $n$  (تعداد مشاهدات) زیاد باشد محاسبه  $CV$  بسیار زمان‌بر و پرهزینه است. اما برای مدل‌های خطی یا رگرسیون چندجمله‌ای با استفاده از ماتریس  $hat$  محاسبه آن بسیار ساده و آسان است. اگر مقادیر روی قطر اصلی ماتریس  $H$  را با  $h_{ii}$  نشان دهیم، آماره اعتبارسنجی متقابل به صورت معادله‌ی زیر قابل محاسبه است:

$$CV(\lambda) = \sum_{i=1}^n \left( \frac{Y_i - \hat{Y}_i}{1 - h_{ii}} \right)^2 \quad (3.3)$$

و یک نتیجه بسیار شگفت‌آور این که برای محاسبه  $CV$  تنها یک‌بار مدل به کل داده‌ها برازش داده می‌شود [۴، ۱۳].

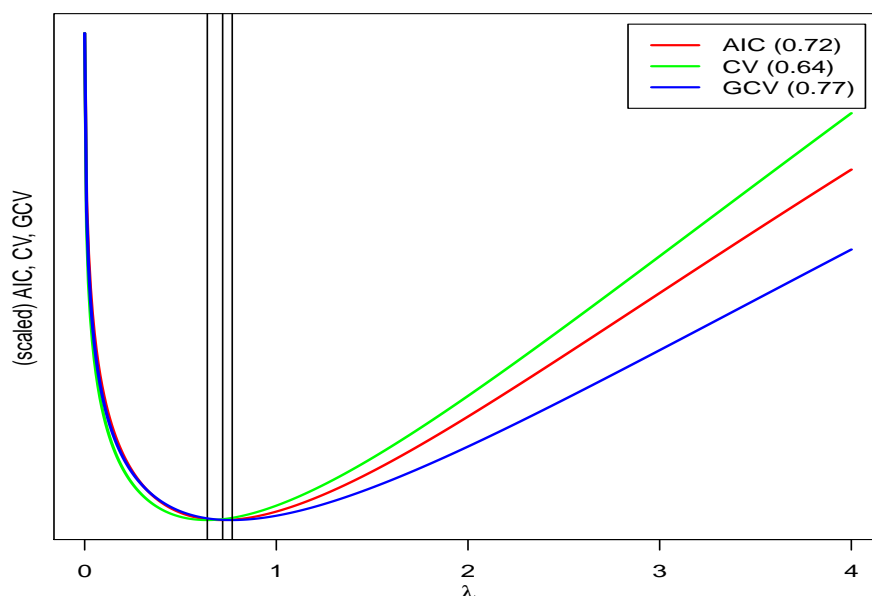
### ۳.۱.۳ اعتبارسنجی متقابل تعمیم‌یافته

روش دیگر برای انتخاب بهینه  $\lambda$  استفاده از معیار اعتبارسنجی متقابل تعمیم‌یافته که به صورت زیر است:

$$GCV(\lambda) = \sum_{i=1}^n \frac{(Y_i - \hat{Y}_i)^2}{(n - \sum_{i=1}^n h_{ii})^2} \quad (4.3)$$

این معیار به گونه‌ای است که کمترین مقدار اعتبارسنجی متقابل تعمیم‌یافته، بهترین پارامتر هموارسازی را ارائه می‌دهد [۸، ۳۲].

در شکل ۱.۳ برای داده‌های موتور سیکلت شامل  $n = 133$  مشاهده برای آزمودن کلاه ایمنی که در آن متغیر وابسته، زمان پس از تاثیر (بر حسب میلی ثانیه) و متغیر مستقل شتاب سر (بر حسب گرم) است سیلورمن (۱۹۸۵)، نتایج مربوط به انتخاب پارامتر هموارسازی بر حسب معیارهای  $CV$ ،  $AIC$  و  $GCV$  نشان داده شده است. مقادیر به دست آمده هر سه معیار، دقیقاً یک مقدار را نشان نمی‌دهند اما نزدیک  $\lambda = 0.75$  مقدار مینیمم را اختیار نموده‌اند.



شکل ۱.۳: پارامتر هموارسازی در مقابل معیار انتخاب برای تابع رگرسیونی P-اسپلاینی از مرتبه  $m = 4$

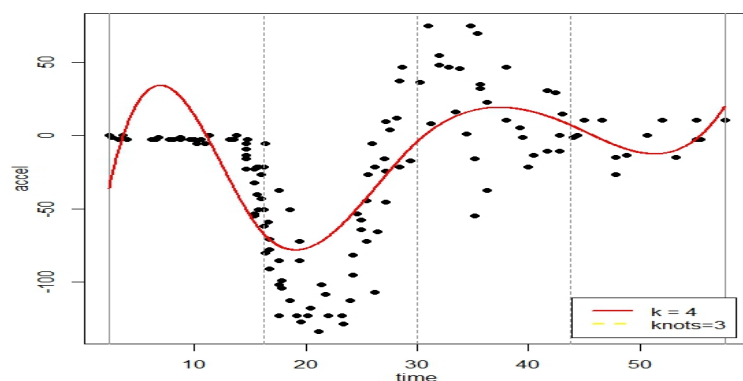
## ۲.۳ روش‌های انتخاب تعداد و موقعیت گره

یکی دیگر از چالش‌ها در رگرسیون اسپلاینی انتخاب تعداد و موقعیت گره‌ها می‌باشد. در رگرسیون اسپلاینی محدوده‌هایی انعطاف‌پذیری بیشتری دارند که تعداد گره‌ها در آن بیشتر باشد، زیرا در این محدوده‌ها تعداد ضرایب چندجمله‌ای به سرعت تغییر می‌کند و اثر زیادی روی برازش اسپلاینی دارد در نتیجه منحنی برازشی ناهموار خواهد بود. در مکان‌هایی که تعداد گره‌ها کمتر است به نظر می‌رسد منحنی پایدارتر است [۱۳].

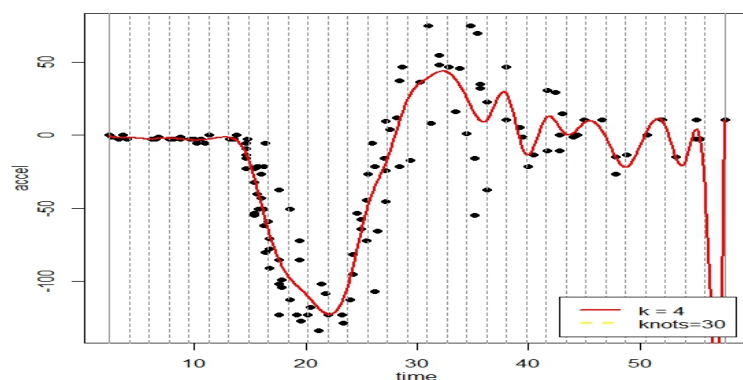
در شکل ۲.۳ و ۲.۴ داده‌های موتور سیکلت با استفاده از P-اسپلاین مرتبه  $k = 4$  با تعداد گره‌های مختلف برازش داده شده است که در شکل ۲.۳ با در نظر گرفتن تعداد کم گره، داده‌ها کم برازش شده‌اند به عبارتی اگر براساس تعداد پارامترهای کم، عمل برازش را انجام دهیم، مدل دچار کم‌برازشی خواهد شد و قدرت پیش‌بینی را از دست خواهد داد و به علت کمبود پارامترهای به کار رفته، خطای زیادی در مدل به وجود می‌آید. البته رفع مشکل کم‌برازشی با افزایش تعداد پارامترهای مدل قابل رفع است، ولی نکته‌ای که باید رعایت کرد، انتخاب تعداد پارامترهای مناسب در مدل است که آن را دچار بیش‌برازشی نکند.

در شکل ۲.۳ با در نظر گرفتن تعداد زیاد گره‌ها داده‌ها بیش برازش شده‌اند. مدل بیش‌برازش، مدلی بسیار پیچیده‌ای است به این معنی که در تحلیل رگرسیونی، مدلی با بیشترین پارامترها ایجاد می‌شود. در این گونه موارد، اگر مدل رگرسیون به دست آمده، برای پیش‌بینی نمونه دیگری به کار رود، مقدارهای پیش‌بینی شده اصلاً مناسب به نظر نخواهند

رسید. انتظار ما از یک تحلیل رگرسیون مناسب، ایجاد مدلی است که نه تنها بتواند برای داده‌های مربوط به نمونه برازش مناسب را انجام دهد، بلکه برای داده‌هایی جدید نیز امکان برآورد مناسب وجود داشته باشد.

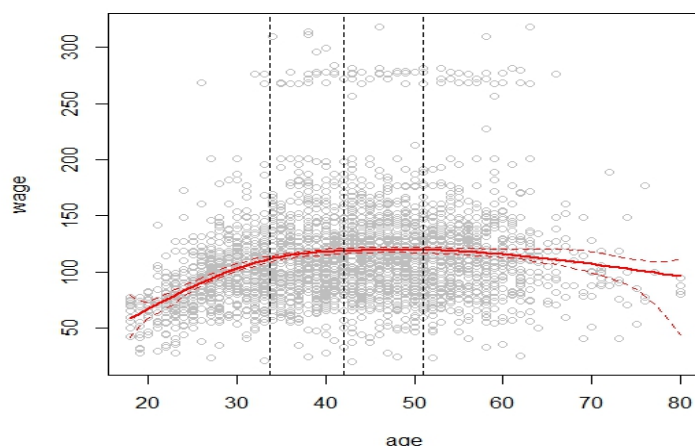


شکل ۲.۳: برازش P-اسپلاین مرتبه  $k = 4$  با تعداد گره کم



شکل ۳.۳: برازش P-اسپلاین مرتبه  $k = 4$  با تعداد گره زیاد

یک روش برای انتخاب تعداد گره‌های مناسب این است که تعداد گره‌های مختلف را به کار ببریم و بهترین منحنی تولید شده را انتخاب کنیم. در عمل بیشتر رایج است که چندک‌ها را به عنوان گره در نظر می‌گیرند و با مشخص کردن درجه آزادی، نرم افزار محل مربوط به تعداد گره‌ها را در چندک داده‌ها قرار می‌دهد. در شکل ۴.۳ یک اسپلاین مکعبی طبیعی با چهار درجه آزادی به داده‌ها برازش داده شده است و موقعیت گره‌ها به طور خودکار در ۲۵، ۵۰ و ۷۵ درصد از داده‌ها قرار گرفته است.



شکل ۴.۳: برازش اسپلاین مکعبی طبیعی با چهار درجه آزادی و قرار گرفتن موقعیت گره‌ها در چندک داده‌ها

اکنون یکی از چالش‌ها، انتخاب تعداد مطلوب گره و انتخاب درجه آزادی مناسب برای تابع اسپلاینی است. یک روش آن است که تعداد گره‌های مختلف را به کار ببریم و بهترین منحنی تولید شده متناظر با آن را مشاهده کنیم. روش دیگر استفاده از اعتبار سنجی متقابل است که معرفی شد، به این صورت که بخش‌هایی از داده‌ها را حذف کنیم (به‌طور مثال ۱۰٪) و یک تابع اسپلاینی با تعداد گره مشخص به داده‌های باقیمانده برازش داده و سپس مقدار خطای پیش‌گویی را به دست آوریم. این فرایند برای تعداد مختلف  $k$  گره می‌تواند تکرار شود تا زمانی که هر مشاهده یک‌بار خارج شود، سپس تعداد گره‌ها برابر کوچکترین مقدار خطای پیش‌گویی خواهد بود [۱۳].

روش دیگر استفاده از الگوریتمی که توسط روپرت و کارول در سال ۲۰۰۰ ارائه شد که به صورت زیر است:

- ۱- تعداد گره  $k = 5$  قرار بده و مقدار  $\hat{\alpha}_{GCV}$  را محاسبه کن.
- ۲- تعداد گره  $k = 10$  قرار بده و مقدار  $\hat{\alpha}_{GCV}$  متناظر آن را محاسبه کن.
- ۳- اگر  $\hat{\alpha}_{GCV} > 0.98 \hat{\alpha}_{GCV}$  آن‌گاه  $k = 10$  نهایی شود. در غیر این صورت مقدار  $k$  افزایش می‌یابد.

## ۳.۳ انتخاب مدل

انتخاب مدل بخش مهمی از هر تحلیل آماری است. فرآیندی است که بهترین مدل را از مجموعه مدل‌های کاندید انتخاب می‌کند. هدف از انتخاب مدل، انتخاب یک مدل مناسب



است که به طور کامل داده‌ها را توضیح دهد و مقدار خطا را به حداقل برساند [۴]. برای گزینش بهترین مدل از معیارهای انتخاب مدل شامل  $AIC$ ،  $BIC$ ،  $CV$  و  $C_p$  ... استفاده می‌شود. مدلی که پایین‌ترین مقادیر معیارهای ذکر شده را به خود اختصاص دهد به عنوان مدل با دقت برآزش بالا انتخاب می‌شود و مدلی که از نظر معیار مذکور بالاترین مقادیر را داشته باشد به عنوان مدل با عملکرد پایین معرفی می‌شود. به دو معیار اطلاع اکائیک و اعتبارسنجی اشاره شد اکنون در مورد معیار اطلاع بیزی و آماره‌ی  $C_p$  بحث می‌کنیم.

### ۱.۳.۳ معیار اطلاع بیزی (شوارتز)

یکی دیگر از معیارهایی که می‌تواند برای بررسی و انتخاب بهترین مدل به کار رود، معیار اطلاع بیزی<sup>۴</sup> ( $BIC$ ) است که توسط شوارتز در ۱۹۷۸ ارائه شد و براساس رابطه‌ی زیر قابل محاسبه می‌باشد:

$$BIC_p = n \ln\left(\frac{SSE_p}{n}\right) + p(\ln n)$$

که  $p$  تعداد پارامترهای قابل برآورد در مدل و  $n$  تعداد مشاهدات در نمونه تصادفی است.

### ۲.۳.۳ آماره $C_p$ مالو

آماره‌ی  $C_p$  برای قضاوت در مورد یک مدل، باید به جای میانگین مربعات انحراف از مدل، میانگین مربعات خطای مقدار پیش‌بینی شده در نظر گرفته شود. میانگین مربعات خطای پیش‌بینی استاندارد شده  $j_p$  برای داده‌های مشاهده شده به صورت زیر محاسبه می‌شود.

$$j_p = \frac{1}{\sigma^2} \sum_{i=1}^n MSE(\hat{y}_i)$$

که در آن  $MSE(\hat{y}_i)$  میانگین مربعات خطای  $i$ امین مقدار پیش‌بینی شده از یک معادله‌ی  $p$  جمله‌ای و  $\sigma$  واریانس مانده‌ها است. برای برآورد  $j_p$  مالو (۱۹۷۳) آماره‌ی زیر را به کار برد:

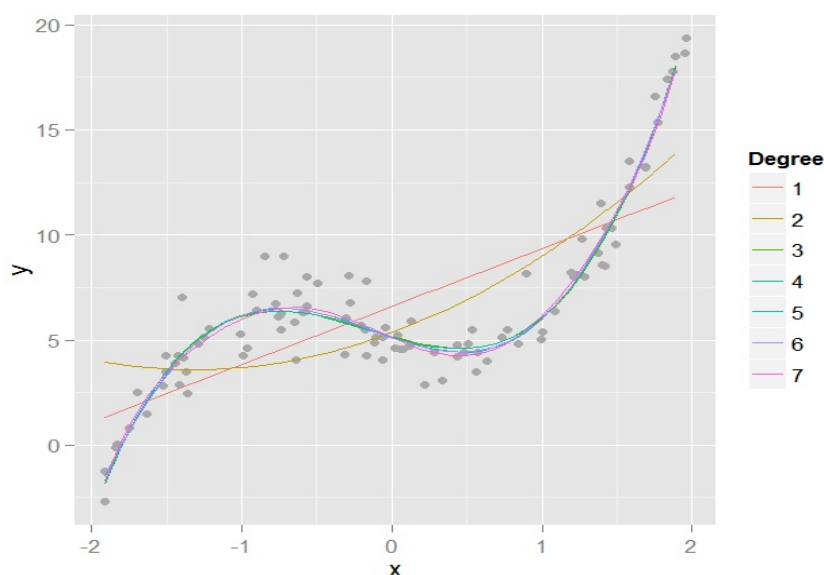
$$C_p = \frac{SSE_p}{\hat{\sigma}^2} + (2p - n)$$

که در آن  $\hat{\sigma}^2$  برآورد  $\sigma^2$  است. بر اساس این معیارها، مدلی که دارای کمترین مقدار معیار باشد به عنوان بهترین مدل انتخاب می‌شود.

<sup>۴</sup> Bayesian information criterion

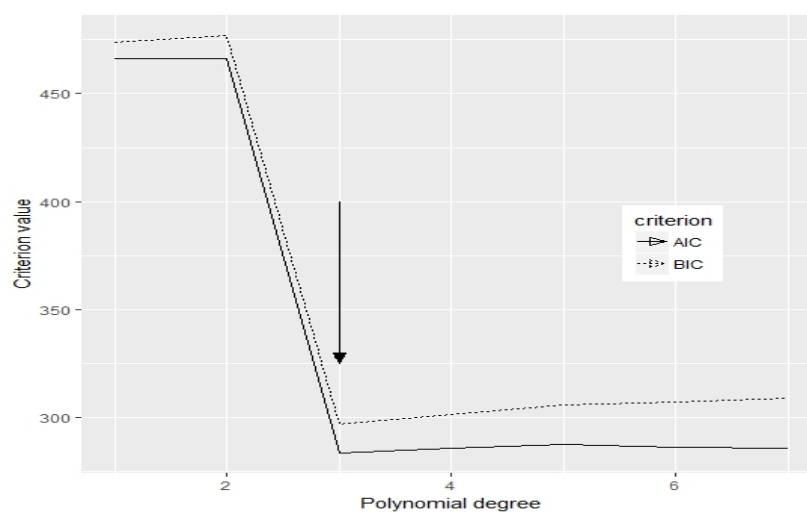
استون (۱۹۹۷) نشان داد که در مدل‌های خطی، معیار اطلاع آکائیک و اعتبارسنجی متقابل با هدف پیش‌گویی معادل هستند. برای بررسی این ایده، مجموعه داده‌ای از رگرسیون چندجمله‌ای شبیه‌سازی شده است و چندجمله‌ای از مرتبه یک تا هفت را به داده‌ها برازش داده و با استفاده از معیارهای فوق صحت آن را می‌آزماییم. همچنین بر حسب علاقه مقدار آماره  $C_p$  و معیار اطلاع بیزی  $BIC$  را محاسبه نموده‌ایم.

شکل ۵.۳ مربوط به برازش مدل چندجمله‌ای از مرتبه یک تا هفت برای داده‌های شبیه‌سازی شده می‌باشد که در آن چندجمله‌ای‌های مرتبه‌ی سه تا هفت برازش بهتری نسبت به چندجمله‌ای با درجات یک و دو دارند و روند داده‌ها را بهتر توصیف می‌کنند. اکنون با استفاده از معیار اطلاع آکائیک، اعتبارسنجی و معیار اطلاع بیزی بهترین مدلی که می‌تواند به داده‌ها برازش یابد را مشخص می‌کنیم.



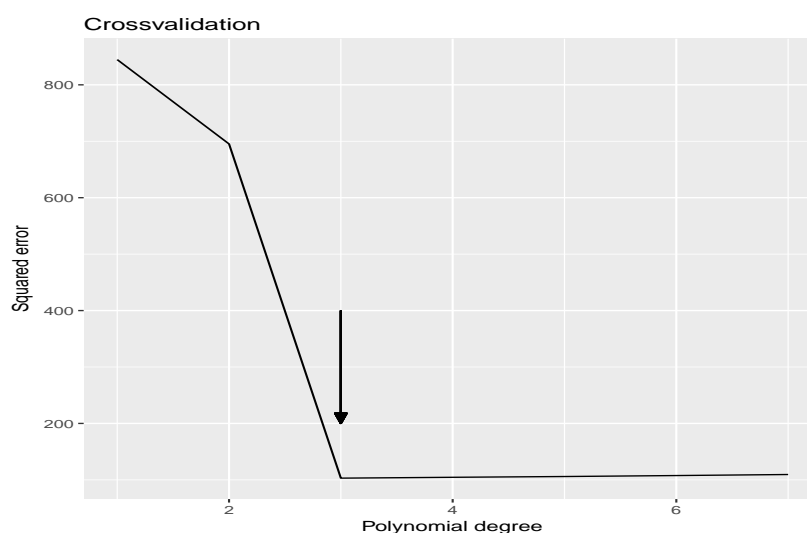
شکل ۵.۳: برازش چندجمله‌ای از مرتبه یک تا هفت برای مجموعه داده‌های شبیه‌سازی شده

طبق شکل ۶.۳ مقدار معیار اطلاع آکائیک و بیزی برای چندجمله‌ای‌ها تا درجه هفت به دست آمده را رسم نموده و در نمودار رسم شده مشهود می‌باشد که چندجمله‌ای درجه سه دارای کمترین مقدار معیارهای فوق می‌باشد در نتیجه بهترین مدلی است که می‌تواند داده‌ها را توصیف کند.



شکل ۶.۳: مقایسه اثربخشی معیار اطلاع اکائیک و بیزی در انتخاب بهترین مدل

همچنین با توجه به نمودار ۷.۳ مقدار  $CV$  برای چندجمله‌ای تا درجه هفت رسم شده گویای این است که مدل چندجمله‌ای با درجه سه به‌عنوان بهترین مدل برازشی است.



شکل ۷.۳: انتخاب بهترین مدل با استفاده از معیار اعتبارسنجی متقابل

در جدول ۱.۳ مقدار معیارهای معیار اطلاع آکائیک، اعتبارسنجی، معیار اطلاع بیزی و آماره‌ی  $C_p$  برای مدل چندجمله‌ای تا درجه هفت به دست آمده و نتایج حاصل حاکی از آن می‌باشد که مدل چندجمله‌ای درجه سه کمترین مقادیر معیارهای فوق را دارا می‌باشد، در نتیجه به‌عنوان مطلوب‌ترین مدل انتخاب می‌گردد.

جدول ۱.۳: جدول مقادیر معیارهای  $AIC$ ,  $CV$ ,  $BIC$ ,  $C_p$  برای انتخاب بهترین مدل

| درجه آزادی | مجموع توان دوم $CV$ | $BIC$    | $C_p$    | $AIC$    |
|------------|---------------------|----------|----------|----------|
| ۱          | ۶۱۸/۰۰۴۸            | ۴۷۳/۸۹۱۷ | ۵۲۵/۵۹۴۸ | ۴۶۶/۰۷۶۲ |
| ۲          | ۶۴۰/۲۸۹۵            | ۴۷۶/۸۰۳۲ | ۵۱۷/۱۵۵۵ | ۴۶۶/۳۸۲۵ |
| ۳          | ۱۰۰/۱۵۳۳            | ۲۹۶/۸۰۶۶ | ۴/۴۷۹۸   | ۲۸۳/۷۸۰۷ |
| ۴          | ۱۰۴/۴۸۸۹            | ۳۰۱/۴۱۰۷ | ۶/۴۷۸۸   | ۲۸۵/۷۷۹۷ |
| ۵          | ۱۰۸/۳۷۶۵            | ۳۰۵/۸۹۴۲ | ۸/۳۶۱۴   | ۲۸۷/۶۵۸۰ |
| ۶          | ۱۰۷/۱۴۳۴            | ۳۰۷/۰۸۶۶ | ۷/۱۲۸۳   | ۲۸۶/۲۴۵۳ |
| ۷          | ۱۱۰/۳۴۲۸            | ۳۰۹/۱۰۵۲ | ۶/۷۵۰۴   | ۲۸۵/۶۵۸۷ |



## فصل ۴

### برخی از کاربردهای اسپلاین‌ها

در این فصل کاربرد مدل‌های GAMLSS که چارچوبی بسیار عمومی برای تحلیل داده‌ها به منظور مدل‌سازی و یادگیری در زمینه‌های مختلف را فراهم می‌سازد ارائه شده است. این توابع در حال حاضر به‌طور گسترده‌ای در مدل‌سازی آماری مورد استفاده قرار می‌گیرند. از مزایای بسیار عالی این توابع فراهم کردن انتخاب گسترده‌ای از مدل‌های رگرسیونی و برازش مدلی انعطاف‌پذیر به داده‌ها می‌باشد. در این مدل متغیر پاسخ می‌تواند هر توزیعی داشته باشد و تمام پارامترهای توزیع می‌توانند مدل شوند. توابع GAMLSS در بسیاری از زمینه‌های کاربردی مورد استفاده قرار گرفته است از جمله علوم اجتماعی، زیست‌شناسی، انرژی، هواشناسی، بارش باران و .... مثال‌های استفاده شده زیر نمونه کوچکی برای نشان دادن قابلیت توابع GAMLSS با استفاده از نرم افزار R می‌باشد. این فصل برگرفته از کتاب استامینوپولوس و همکارانش [۲۸] می‌باشد.

مثال اول مربوط به مجموعه داده‌های اجاره که طریقه‌ی استفاده از توابع GAMLSS را نشان می‌دهد و با استفاده از این توابع یک مدل پیچیده سه پارامتری به داده‌ها برازش می‌یابد. مثال دوم مجموعه داده‌های گونه‌های مختلف ماهی که مثالی ساده از برازش توزیع‌های شمارشی مختلف با استفاده از توابع GAMLSS به داده‌ها است. مثال سوم مجموعه داده‌های مدت بستری که متغیر پاسخ دارای توزیع دوجمله‌ای است. این داده‌ها در توسعه توابع GAMLSS نقش زیادی داشته‌اند، چون اولین مجموعه داده بودند که محققان پی به وجود توابع GAMLSS بردند. توزیع‌های دو و سه پارامتری به کار رفته در مثال‌ها در پیوست در جدول ۱.آ و ۲.آ با

جزئیات بیشتری شرح داده شده است.

## ۱.۴ مدل‌های جمعی

• مدل جمعی تعمیم‌یافته (GAM)

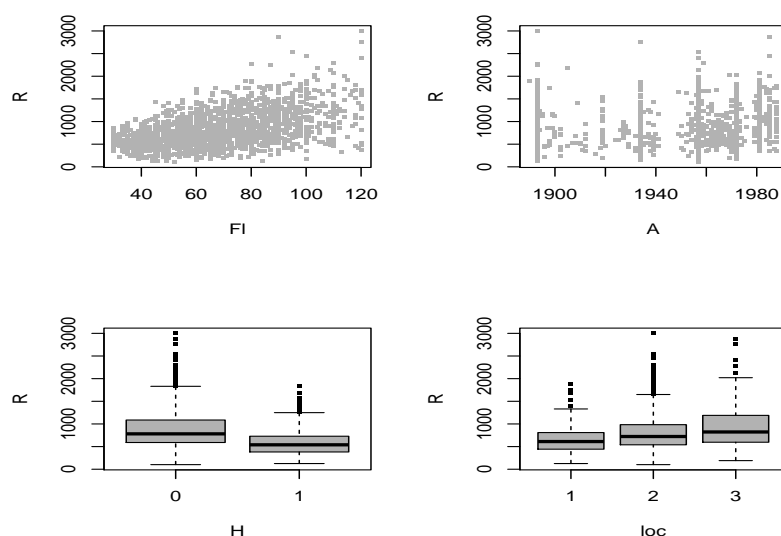
• مدل جمعی تعمیم‌یافته برای مکان، مقیاس و شکل (GAMLSS)

مدل جمعی جایگزینی برای مدل‌های خطی می‌باشد و می‌تواند برای تعیین روابط بین متغیرها مورد استفاده قرار گیرد. این مدل می‌تواند کیفیت پیش‌بینی متغیر پاسخ را به حداکثر برساند و روابط غیرخطی بین متغیر پاسخ و مجموعه متغیرهای پیش‌گو را کشف کند. در شرایطی که پراکندگی داده‌ها زیاد باشد، برای ساخت مدلی دقیق‌تر بهتر است از مدل GAMLSS استفاده شود تا از خانواده‌ی توزیع‌های عمومی به‌جای خانواده‌ی توزیع‌های نمایی استفاده شود. مدل GAMLSS این قابلیت را دارد که می‌تواند داده‌های پراکنده را با خطای کمتری برازش نماید. از مزایای این مدل نسبت به مدل جمعی تعمیم‌یافته که فقط به مدل‌سازی میانگین شرطی متغیرهای پاسخ محدود هستند این است که برای هر پارامتر توزیع رابطه‌ی رگرسیونی ایجاد می‌کند. در فصل اول مدل‌های GAM و GAMLSS به‌طور کامل شرح داده شده است.

### ۱.۱.۴ کاربرد اول: داده‌های اجاره مونیخ

داده‌های مربوط به اجاره بر اساس یک نظرسنجی در آوریل ۱۹۹۳ توسط اینفراتست سوزیال جمع‌آوری شد که در آن یک نمونه تصادفی از محل اقامت با موافقت‌نامه اجاره جدید یا افزایش اجاره طی چهار سال گذشته در مونیخ انتخاب شد. این مجموعه داده‌ها شامل،  $n = 1969$ ، مشاهده است که اجاره خالص ماهیانه نقش متغیر پاسخ و متغیرهای توضیحی متراژ (F1)، سال ساخت (A)، گرمایش مرکزی (H) شامل دو سطح، وجود دارد (۰) و عدم وجود (۱) و عاملی که نشان دهنده‌ی موقعیت ساختمان (Loc) شامل سه سطح متوسط رو به پایین (۱)، متوسط (۲)، متوسط رو به بالا (۳).

شکل ۱.۴ نمودار متغیر پاسخ، میزان اجاره (R)، در مقابل هر یک از متغیرهای توضیحی را نشان می‌دهد.

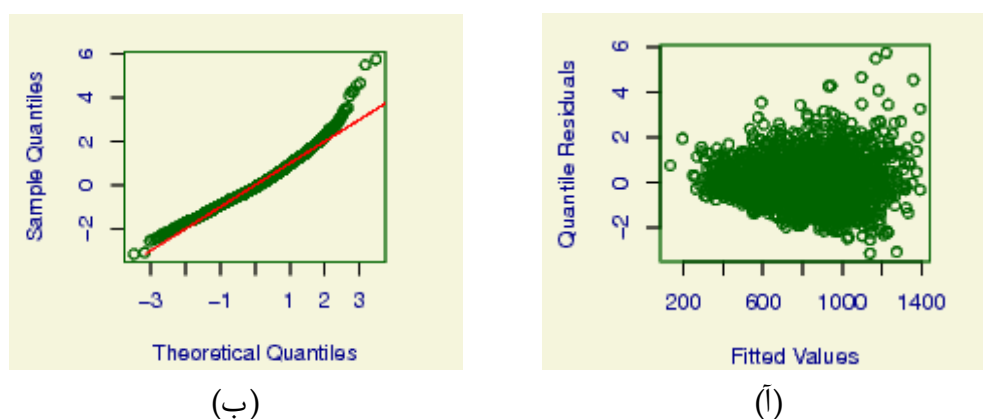


شکل ۱.۴: نمودار متغیر پاسخ (میزان اجاره R) در مقابل هریک از متغیرهای توضیحی

نمودار ترسیم شده‌ی میزان اجاره (R) در مقابل متراژ (FI) نشان می‌دهد که با افزایش متراژ میزان اجاره نیز افزایش می‌یابد و رابطه‌ای مثبت دارند در نتیجه فرض همگنی واریانس نقض خواهد شد و ممکن است واریانس داده‌های اجاره به میانگین یا به متغیرهای توضیحی بستگی داشته باشد. همچنین چولگی مثبتی هم در توزیع داده‌های اجاره نیز وجود دارد. نمودار ترسیم شده میزان اجاره در مقابل سال ساخت و ساز نشان می‌دهد که تا سال ۱۹۶۰ میزان اجاره تقریباً ثابت است اما با ساخت و سازهای بعد آن سال، روندی افزایشی در میزان اجاره وجود دارد. دو نمودار جعبه‌ای نشان می‌دهد که میزان اجاره طبق عامل‌های توضیحی متفاوت است یعنی میزان متوسط اجاره اگر گرمایش مرکزی و موقعیت تغییر کند از متوسط رو به کم به متوسط و به متوسط رو به بالا افزایش می‌یابد.

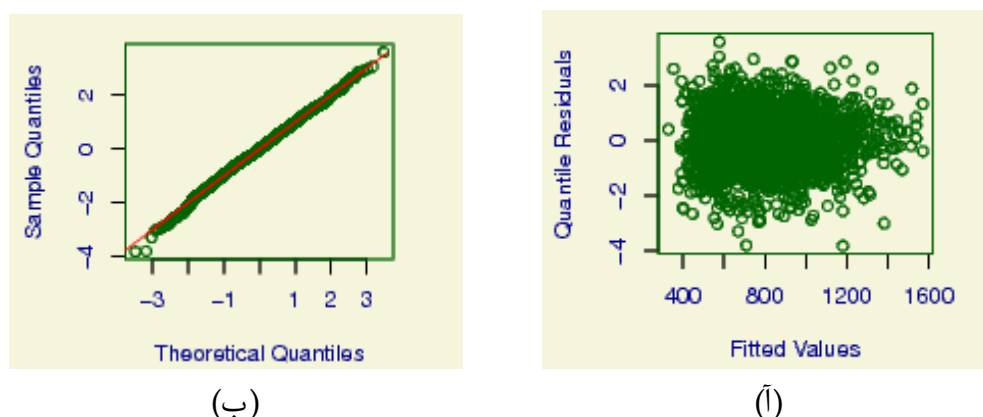
اکنون باید مدل برازشی به داده‌های اجاره قادر به ارائه مدلی باشد که رابطه‌ی غیرخطی بین متغیر پاسخ و متغیرهای توضیحی، عدم همگنی واریانس داده‌ها و همچنین چولگی مثبت توزیع داده‌های اجاره را بیان کند. مدلهایی را به داده‌های اجاره برازش می‌دهیم و با استفاده از تابع GAMLSS در نرم افزار R به تحلیل داده‌های اجاره می‌پردازیم. برای ارزیابی کفایت یک مدل به یک مجموعه داده، می‌توان به بررسی مانده‌ها پرداخت. تابع GAMLSS از مانده‌های استاندارد شده استفاده می‌کند. مدل خطی  $R \sim FI + A + H + loc$  را به داده‌ها برازش می‌دهیم. شکل ۲.۴ نمودار مانده‌های حاصل از برازش مدل خطی است.





شکل ۲.۴: نمودار مانده‌های برازش مدل خطی

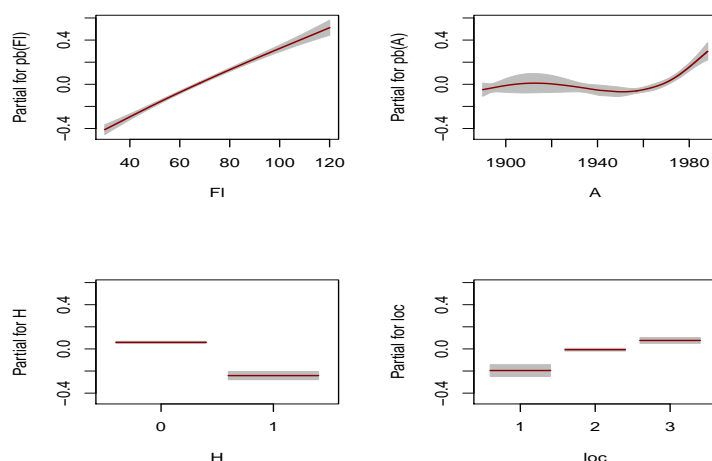
باتوجه به نمودار چندک-چندک شکل (ب)، هرچه نقاط نمودار به خط نزدیک‌تر باشند، توزیع داده‌ها به نرمال نزدیک‌تر است در نتیجه فرض نرمال بودن توزیع داده‌ها پذیرفته نمی‌شود. نمودار شکل (ا)، مانده‌ها در مقابل مقادیر برازشی بر روی خطی افقی در صفر به‌طور تصادفی پراکنده نشده‌اند و این نشان از ناهمگونی واریانس دارد یعنی واریانس با افزایش میانگین تغییر می‌کند. با توجه به بررسی‌های مدل، فرض نرمال بودن توزیع داده‌ها پذیرفته نمی‌شود. فرض می‌کنیم مجموعه داده‌های اجاره دارای توزیع گاما که مدل خطی تعمیم‌یافته است و برای بررسی کفایت این مدل به مانده‌ها می‌پردازیم. در شکل ۳.۴ نمودارهای ترسیمی نمایش داده شده است.



شکل ۳.۴: نمودار مانده‌های برازش مدل گاما

مانده‌ها در شکل ۳.۴ وضعیت بسیار بهتری نسبت به شکل ۲.۴ دارند در آن ناهمگنی در مانده‌ها در مقابل مقادیر برازشی از بین رفته‌اند و همچنین خمیدگی نمودار چندک-چندک به شدت کاهش یافته است، اگرچه ناهمگونی کمی در مانده‌ها وجود دارد.

در ادامه مدل‌های جمعی تعمیم‌یافته را برازش می‌دهیم که برای بررسی رابطه‌ی بین متغیرهای پیش‌گو و پاسخ توابع هموارساز (pb) را به کار می‌برد که مدلی انعطاف‌پذیر از هموارسازهای اسپلاین جریمه‌ای ارائه می‌دهد.



شکل ۴.۴: نمودار مقادیر برازشی برای مدل (GAM)

در نمودارهای ۴.۴ مشهود است که میزان اجاره تقریباً با افزایش مترآژ (F1) به طور خطی افزایش می‌یابد و با سال ساخت (A) رابطه‌ای غیرخطی دارد و تا سال ۱۹۶۰ تقریباً ثابت است اما بعد آن رو به افزایش است و میزان تاثیر دو عامل H و loc همان مقدار مورد انتظار است، یعنی میزان اجاره کم است اگر عامل  $H=1$  باشد و همچنین نسبت به سطوح موقعیت، میزان اجاره از کم به متوسط و به زیاد افزایش پیدا می‌کند.

روش دیگر برای ارزیابی کفایت مدل، استفاده از نمودار مارپیچ<sup>۱</sup> برای تحلیل باقیمانده‌ها است، می‌تواند به عنوان نمودار چندک-چندک در نظر گرفته شود که توسط برن و فریدریک در سال ۲۰۰۱ معرفی شد. ابزاری تشخیصی است برای مقایسه مدل‌ها و مشخص نمودن مدل آماری متناسب با داده‌ها و پیدا کردن مکان‌هایی که مدل می‌تواند بهبود پیدا کند. دارای ویژگی‌هایی است:

- نقاط نمودار مارپیچ هر چه به خط افقی نزدیک‌تر باشد توزیع باقیمانده‌ها به توزیع نرمال استاندارد نزدیک‌تر است.
- در صورتی که مدل کفایت داشته باشد انتظار می‌رود که به طور تقریبی ۹۵٪ نقاط بین دو نمودار بیضوی و ۵٪ خارج از آن باشد.
- شکل (الگو) منحنی برازشی به نقاط نمودار مارپیچ، بازتابنده‌ی عدم کفایت مدل برازشی است که در جدول ۱.۴ شرح داده شده است. به عنوان مثال اگر سطح نقاط نمودار مارپیچ

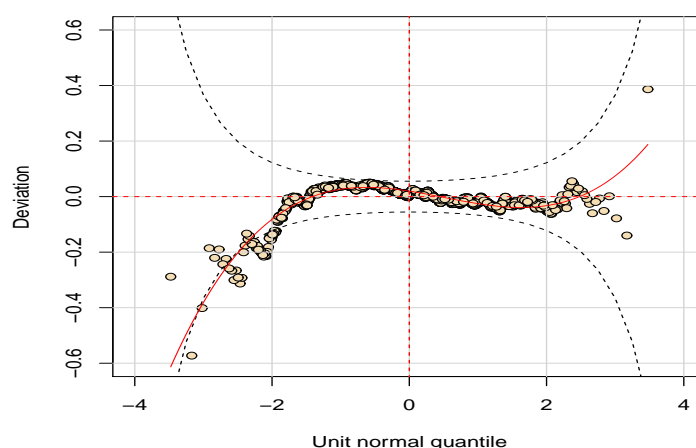
<sup>۱</sup>Worm plot

بالتر از خط افقی در مبدا باشد این نشان دهنده میانگین بسیار بالای مانده‌ها است، بدان معنی است که موقعیت توزیع برازشی بسیار پایین است که می‌توان با افزایش پارامتر  $\mu$  اگر پارامتری مکانی است مدل را بهبود بخشید یا توزیع مدل را تغییر داد.

جدول ۱.۴: تفسیر الگوهای مختلف در نمودار مارپیچ

| شکل منحنی برازشی                  | گشتاور                 | توزیع برازش داده‌شده                |
|-----------------------------------|------------------------|-------------------------------------|
| سطح نقاط اگر بالای مبدا باشد      | میانگین بسیار بالا است | موقعیت برازش بسیار پایین است        |
| سطح نقاط اگر زیر مبدا باشد        | میانگین بسیار کم است   | موقعیت برازش بسیار بالا است         |
| خط دارای شیب مثبت باشد            | واریانس بسیار بالا است | مقیاس برازش داده‌شده بسیار کم است   |
| خط دارای شیب منفی باشد            | واریانس بسیار کم است   | مقیاس برازش داده‌شده بسیار بالا است |
| U-شکل باشد                        | چولگی مثبت دارد        | توزیع به سمت چپ کشیده شده است       |
| ∩-شکل باشد                        | چولگی منفی دارد        | توزیع به سمت راست کشیده شده است     |
| S-شکل باشد که به پایین خم شده است | کشیده                  | توزیع برازشی بسیار دم‌باریک است     |
| S-شکل باشد که به بالا خم شده است  | خم پخ                  | توزیع برازشی بسیار دم‌کلفت است      |

با توجه به نمودار ۵.۴ برای داشتن مدل برازشی مطلوب انتظار داریم که نقاط نزدیک به خط افقی میانه باشد و به‌طور تقریبی ۹۵٪ نقاط بین نمودار بیضوی بالا و پایین باشد و ۵٪ خارج از آن باشد.



شکل ۵.۴: نمودار مارپیچ برای مدل (GAM)

هم‌چنین برای ارزیابی عملکرد پیش‌بینی مدل برازشی می‌توان از معیار AIC استفاده نمود. دو مدل خطی تعمیم‌یافته و جمعی تعمیم‌یافته‌ی برازشی را با استفاده از معیار AIC مقایسه

می‌کنیم طبق جدول ۱۰.۴ کمترین AIC مربوط به مدل جمعی تعمیم‌یافته به میزان ۲۷۷۰۵/۶۵ است. در نتیجه در مقایسه با دو مدل دیگر مدل بهتری است.

جدول ۲.۴: گزینش بهترین مدل با استفاده معیارهای AIC

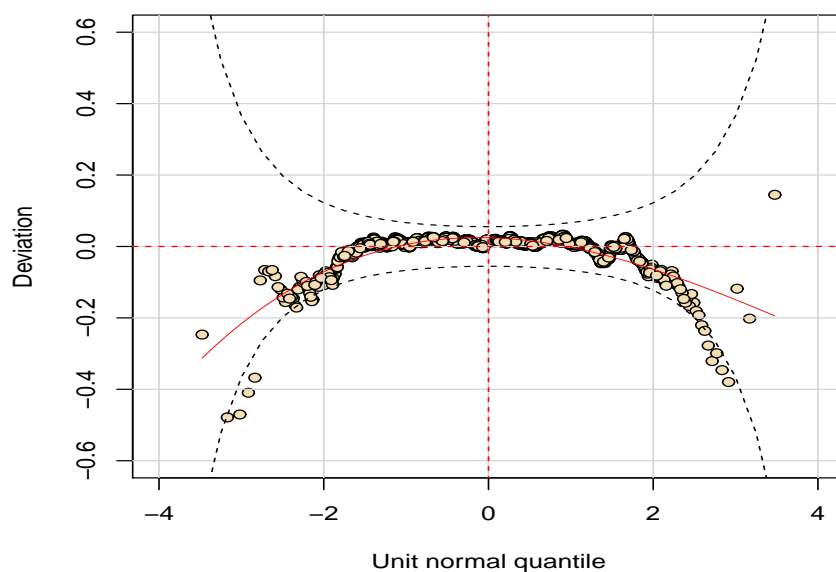
| AIC      | پیش‌گوی جمعی       | تابع پیوند |       | مدل                         |
|----------|--------------------|------------|-------|-----------------------------|
|          |                    | $\sigma$   | $\mu$ |                             |
| ۲۸۱۷۳    | Fl+A+H+loc         | log        | ident | نرمال                       |
| ۲۷۷۷۸/۵۹ | Fl+A+H+loc         | log        | log   | گاما (مدل خطی تعمیم‌یافته)  |
| ۲۷۷۰۵/۶۵ | pb(Fl)+pb(A)+H+loc | log        | log   | گاما (مدل جمعی تعمیم‌یافته) |

برای بهبود بخشیدن به مانده‌ها می‌توان پارامتر  $\sigma$  را مدل‌سازی کرد. داده‌های اجاره مونیخ را با استفاده از توزیع‌های گامای دو پارامتری و گاوسی وارون مدل می‌کنیم، مقدار AIC برای این دو مدل فوق و گاما که پارامتر مقیاس آن ثابت فرض شده است محاسبه شده است. با توجه به جدول ۳.۴ طبق معیار AIC مدلی که با استفاده از توزیع گاما دو پارامتری که پارامتر مقیاس آن مدل شده است، مدل مطلوب‌تری است.

جدول ۳.۴: گزینش بهترین مدل با استفاده معیارهای AIC

| AIC      | پیش‌گوی جمعی                                   | تابع پیوند |       | مدل         |
|----------|--|------------|-------|-------------|
|          |  | $\sigma$   | $\mu$ |             |
| ۲۷۷۰۵/۶۵ | sigma.fo=1 pb(Fl)+pb(A)+H+loc                  | log        | log   | گاما        |
| ۲۷۶۱۴/۷۸ | sigma.fo=pb(Fl)+pb(A)+H+loc pb(Fl)+pb(A)+H+loc | log        | log   | گاما        |
| ۲۷۷۱۶/۶۶ | sigma.fo=pb(Fl)+pb(A)+H+loc pb(Fl)+pb(A)+H+loc | log        | log   | گاوسی وارون |

برای بررسی کفایت مدل گامای دو پارامتری، نمودار مارپیچ در زیر رسم شده است. با توجه به نمودار مارپیچ شکل ۶.۴ تعداد کمی از نقاط وجود دارد که خارج از فاصله اطمینان ۹۵٪ نمودار نقطه‌نقطه می‌باشد، که این نشان دهنده‌ی عدم کفایت مدل می‌باشد، علاوه بر این شکل U وارونه‌ی مانده‌ها اشاره‌ای به چولگی منفی مانده‌ها دارد که حاکی از این می‌باشد که توزیع گاما انعطاف‌پذیری کافی برای از بین بردن چولگی در داده‌ها را ندارد.



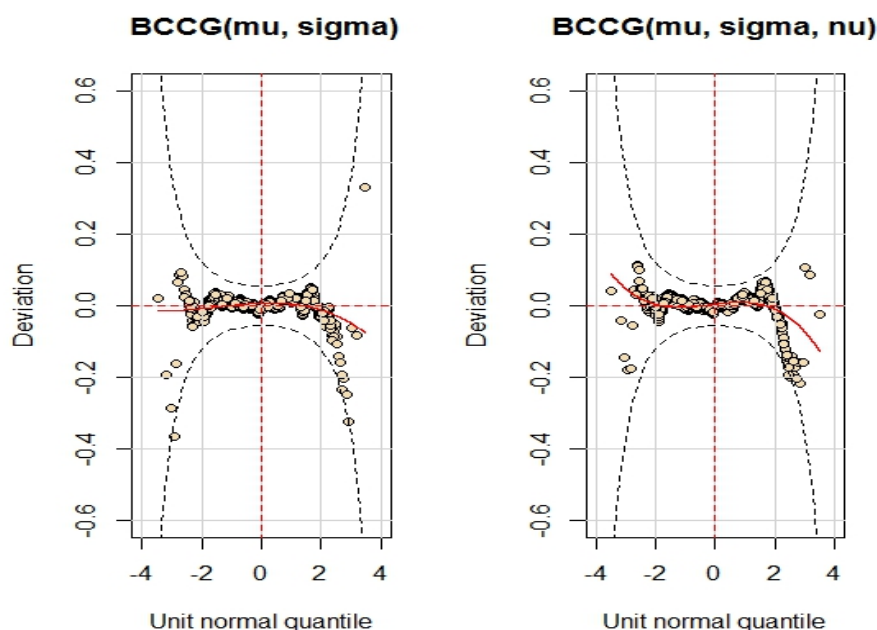
شکل ۶.۴: نمودار مارپیچ برای مدل گامای دو پارامتری

مدل باکس کاکس کول و گرین (BCCGO) که مدلی جمعی تعمیم یافته برای مکان، مقیاس و شکل (GAMLSS) را در نظر می گیریم که انعطاف پذیری بیشتری نسبت به مدل های پیشین دارد، توزیعی سه پارامتری است که هر پارامتر با استفاده از تابع هموار ساز ناپارامتری مدل سازی شده است. دو مدل به داده های اجاره برازش داده، در یک مدل  $\nu$  ثابت در نظر گرفته شده و در مدل دیگر  $\nu$  مدل سازی شده و با گامای دو پارامتری مقایسه شده است. با توجه به جدول ۴.۴ طبق معیار AIC مدلی که پارامتر  $\nu$  مدل سازی شده است بهبودی اندکی در برازش مدل ایجاد نموده است.

جدول ۴.۴: گزینش بهترین مدل با استفاده معیارهای AIC

| AIC      | پیش گوی جمعی  | تابع پیوند |          |       | مدل                  |
|----------|---|------------|----------|-------|----------------------|
|          |   | $\nu$      | $\sigma$ | $\mu$ |                      |
| ۲۷۶۱۱/۰۲ | pb(Fl)+pb(A)+H+loc<br>sigma.fo=pb(Fl)+pb(A)+H+loc<br>nu.fo=1                  | ident      | log      | log   | باکس کاکس<br>(BCCGO) |
| ۲۷۶۰۸/۱۵ | pb(Fl)+pb(A)+H+loc<br>sigma.fo=pb(Fl)+pb(A)+H+loc<br>nu.fo=pb(Fl)+pb(A)+H+loc | ident      | log      | log   | باکس کاکس<br>(BCCGO) |

با استفاده از نمودار مارپیچ شکل ۷.۴ کفایت دو مدل برازشی را بررسی می‌کنیم.



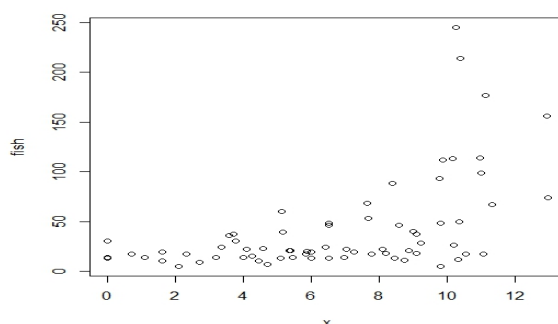
شکل ۷.۴: نمودار مارپیچ برای مدل سه پارامتری

دو نمودار ترسیمی حاکی از کفایت مدل توزیع سه پارامتری می‌باشند.

## ۲.۱.۴ کاربرد دوم: داده‌های گونه‌های ماهی

داده‌های گونه‌های ماهی توسط استین و جیرتیز [۲۹] با استفاده از توزیع معکوس پواسن گاوسی  $PIG(\mu, \sigma)$  با در نظر گرفتن یک مدل خطی برای  $\log(\mu)$  و یک مقدار ثابت برای  $\sigma$  تولید شده‌اند، ارائه شد. این داده‌ها در بسته نرم افزاری gamlss.data موجود می‌باشد که تعداد گونه‌های مختلف ماهی در ۷۰ دریاچه جهان (fish) به‌عنوان متغیر پاسخ و منطقه دریاچه (lake) متغیر توضیحی می‌باشد. در شکل ۸.۴ نمودار مربوط به تعداد گونه‌های ماهی در مقابل متغیر توضیحی  $x=lake$  نشان داده شده است.

داده‌ها را با استفاده از توزیع‌های گسسته مختلف مانند پواسن (PO)، پواسن دوگانه (DPO)، دوجمله‌ای منفی نوع یک و دو (NBI, NBII)، معکوس پواسن گاوسی (PIG)، دلاپورت (DEL)، سیچل (SICHEL) که در بسته‌ی gamlss.dist با تابع پیوندی پیش فرض موجود می‌باشد مدل‌سازی می‌کنیم. متغیر توضیحی  $x=lake$  را یک‌بار با مدل خطی و سپس مدل چندجمله‌ای درجه دوم در مقابل متغیر پاسخ برازش می‌دهیم و با استفاده از معیار AIC مدل‌های برازشی را مقایسه می‌کنیم.



شکل ۸.۴: نمودار گونه‌های ماهی

مقادیر به‌دست آمده، با استفاده از معیار AIC در جدول ۵.۴ آورده شده است.

جدول ۵.۴: مقادیر به‌دست آمده با استفاده از AIC

| مقدار AIC برای برازش متغیر پاسخ در مقابل متغیر توضیحی با فرض خطی بودن            |          |          |          |          |           |          |
|--|----------|----------|----------|----------|-----------|----------|
| PO   | DPO      | NBI      | NBII     | PIG      | DEL       | SICHEL   |
| ۱۹۰۰/۱۵۶۲  | ۶۵۴/۱۶۱۶ | ۶۲۵/۸۴۴۳ | ۶۴۷/۵۳۵۹ | ۶۲۳/۴۶۳۲ | ۴۶۲۶/۲۳۳۰ | ۶۲۵/۳۹۲۳ |
| مقدار AIC برای برازش متغیر پاسخ در مقابل متغیر توضیحی با فرض چندجمله‌ای درجه دوم |          |          |          |          |           |          |
| PO   | DPO      | NBI      | NBII     | PIG      | DEL       | SICHEL   |
| ۱۸۵۵/۲۹۶۵  | ۶۵۵/۲۵۲۰ | ۶۲۲/۳۱۷۳ | ۶۴۵/۰۱۲۹ | ۶۲۱/۳۴۵۹ | ۶۲۳/۵۸۱۶  | ۶۲۳/۰۹۹۵ |

با توجه به مقادیر به‌دست آمده در جدول فوق مقدار AIC برای مدل پواسن نسبت به سایر توزیع‌ها خیلی بزرگ است، بنابراین می‌توان نتیجه گرفت که داده‌ها بیش از حد پراکنده شده‌اند. نتایج نشان می‌دهد که متغیر توضیحی با مدل چندجمله‌ای درجه دوم نسبت به مدل خطی عملکرد بهتری دارد البته بجز توزیع پواسن دوگانه. تا این مرحله مدل معکوس پواسن گاوسی (PIG) با کمترین مقدار AIC به عنوان مدل مطلوب انتخاب می‌گردد.

اکنون در مدل معکوس پواسن گاوسی به‌جای برازش متغیر توضیحی با چندجمله‌ای درجه دوم، هموارگر اسپلینی استفاده می‌کنیم و عملکرد این مدل را مورد ارزیابی قرار می‌دهیم. مقدار AIC برای این مدل به میزان ۶۲۳/۴۶۳۲ است، که مقدار آن افزایش یافت در نتیجه به‌کار بردن هموارگر اسپلینی موثر واقع نشده است. برای بهبود عملکرد مدل می‌توان  $\log(\sigma)$  را مدل‌سازی نمود، اکنون  $\log(\sigma)$  را به‌عنوان یک تابع خطی از  $x$  در توزیع‌های فوق بجز توزیع پواسن که پارامتر  $\sigma$  ندارد مدل کرد.

در جدول ۶.۴ مقادیر AIC برای توزیع‌های فوق به‌دست آمده نشان داده شده است.

جدول ۶.۴: مقادیر به‌دست آمده با استفاده از معیار AIC

| مقدار AIC با مدل کردن پارامتر $\sigma$ |          |          |          |          |          |
|--|----------|----------|----------|----------|----------|
| DPO                                    | NBI      | NBH      | PIG      | DEL      | SICHEL   |
| ۶۲۶/۴۰۵۶                               | ۶۱۴/۹۵۶۵ | ۶۱۵/۱۲۵۰ | ۶۱۲/۳۶۶۷ | ۶۱۴/۶۰۵۹ | ۶۱۳/۷۳۲۷ |

نتایج بیانگر این است که مدل کردن  $\log(\sigma)$  به‌عنوان تابعی از  $x$  عملکرد تمام توزیع‌ها را بهبود بخشیده است و مدل معکوس پواسن گاوسی (PIG) با کمترین مقدار AIC به‌عنوان بهترین مدلی است که می‌تواند داده‌ها را توصیف کند. توزیع‌های سیچل و دلاپورت سه پارامتری می‌باشند که می‌توان برای بهبود بخشیدن به مدل، پارامتر سوم ( $\nu$ ) را به‌عنوان تابعی خطی از  $x$  مدل کرد. توزیع سیچل به‌عنوان پیش‌فرض تابع پیوندی از تابع همانی و توزیع دلاپورت از تابع لوجیت استفاده می‌کند.

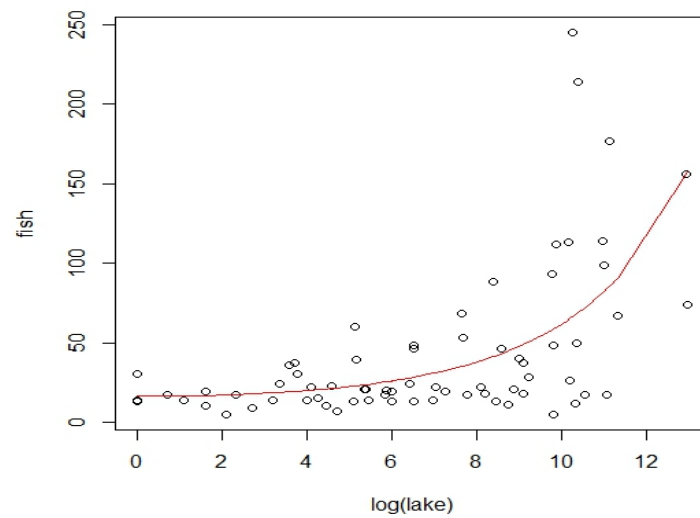
در جدول ۷.۴ مقادیر AIC برای توزیع‌های سیچل و دلاپورت به‌دست آمده نشان داده‌شده است.

جدول ۷.۴: مقادیر به دست آمده با استفاده از معیار AIC

| مقدار AIC برای توزیع‌های سیچل و دلاپورت با مدل کردن پارامتر ( $\nu$ ) |          |
|---|----------|
| DEL   | ۶۱۴/۷۳۷۶ |
| SICHEL  | ۶۱۱/۶۳۴۶ |

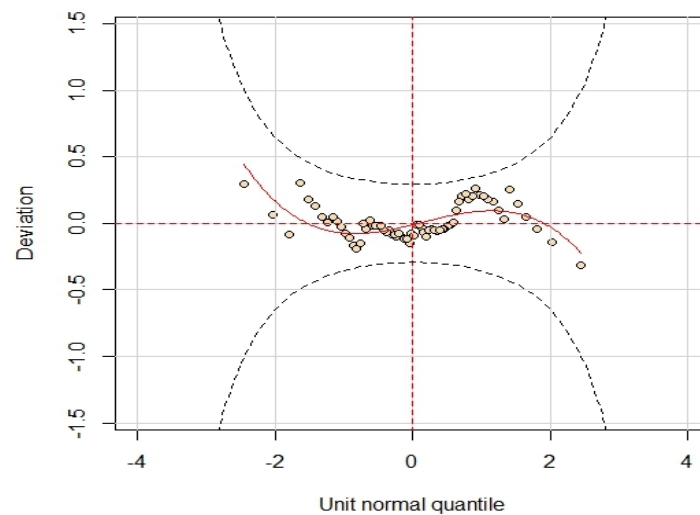
با توجه به نتایج، مدل کردن پارامتر ( $\nu$ ) به‌عنوان تابعی از  $x$ ، مدل سیچل (SIC) را نسبت به مدل معکوس پواسن گاوسی (PIG) بهبود بخشیده است و دارای مقدار AIC کمتری نسبت به آن شده است اما مدل دلاپورت (DEL) بهبودی نیافته است. در شکل ۹.۴ برازش مدل سیچل به داده‌های گونه‌های ماهی نشان داده‌شده است که به‌عنوان بهترین مدلی است که داده‌ها را توصیف می‌کند.





شکل ۹.۴: نمودار برازش مدل سیچل به داده‌های گونه‌های ماهی

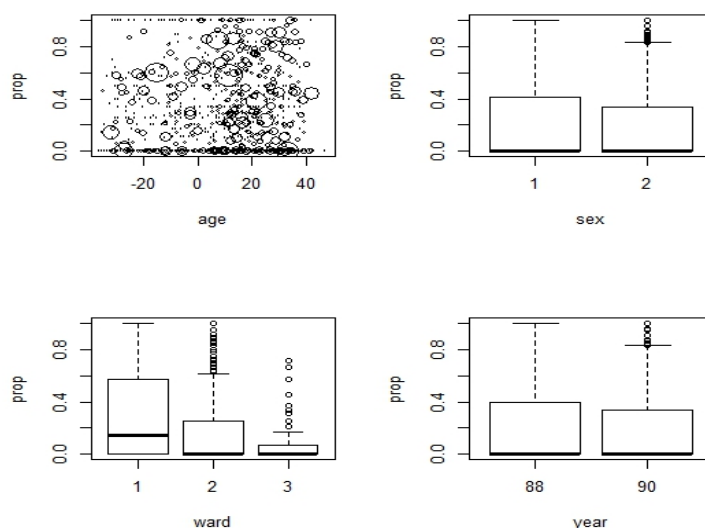
شکل ۱۰.۴ نمودار مارپیچ مدل برازشی سیچل را نشان می‌دهد که نمودار حاکی از کفایت مدل می‌باشد.



شکل ۱۰.۴: نمودار مارپیچ مدل برازشی سیچل به داده‌های گونه‌های ماهی

### ۳.۱.۴ کاربرد سوم: داده‌های مدت بستری در بیمارستان

داده‌های مدت بستری در بیمارستان aep در بسته gamlss.data موجود می‌باشد که دارای ۱۳۸۳ مشاهده، توسط گنگ و همکارانش در طول سال‌های ۱۹۸۸ و ۱۹۹۰ در بیمارستان دلمر بارسلونا جمع‌آوری شده است. متغیر پاسخ، تعداد روزهای غیرمفید از مجموع روزهایی که بیماران در بیمارستان بستری‌اند (noinap) و جنسیت بیمار (sex)، نوع بخش بیمارستان (ward)، سن بیمار (age)، سال ۱۹۸۸ یا ۱۹۹۰ (year)، لگاریتم (مجموع روزهای مدت بستری  $\div 10$ ) (loglos) به‌عنوان متغیرهای توضیحی می‌باشد. هر بیمار برای ماندن غیر مفید در بیمارستان توسط دو پزشک با استفاده از پروتکل ارزیابی مناسب<sup>۲</sup> (aep) مورد بررسی قرار گرفته‌اند. در شکل ۱۱.۴ نمودار نرخ نامناسب متغیر پاسخ در مقابل age, sex, ward, yaer نشان داده‌شده است.



شکل ۱۱.۴: نمودار نرخ نامناسب متغیر پاسخ در مقابل متغیرهای توضیحی در داده‌های مدت بستری

گنگ و همکارانش مدل رگرسیون دوجمله‌ای را برای تعداد روزهای غیرمفید و توزیع‌های دوجمله‌ای و بتا دوجمله‌ای را برای متغیر پاسخ به‌کار بردند و دریافتند که مدل بتا دوجمله‌ای با پارامترهای بازپارامتری شده برازش بهتری را به‌دست می‌آورد. آن‌ها هر دو پارامتر میانگین ( $\mu$ ) و مقیاس ( $\sigma$ ) توزیع بتا دوجمله‌ای را با تابع پیوند لجوجیت و پارامتر مقیاس را با تابع پیوند همانی به‌عنوان تابعی از متغیرهای توضیحی مدل کردند.

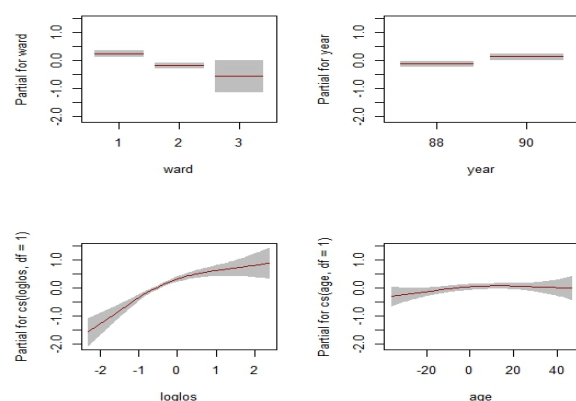
در جدول ۸.۴ چهار مدل به داده‌ها aep برازش داده شده که انعطاف‌پذیری مدل‌های gamlss را نشان می‌دهد و براساس معیارهای AIC و BIC بهترین مدل برازشی انتخاب می‌شود.

<sup>۲</sup> Appropriateness evaluation protocol

جدول ۸.۴: گزینش بهترین مدل با استفاده معیارهای AIC و BIC

| مدل | تابع پیوند           | پیش‌گوی جمعی                           | AIC    | BIC    |
|-----|----------------------|--|--------|--------|
| M1  | $\text{logit}(\mu)$  | ward+loglos+year                       | ۴۵۳۳/۴ | ۴۵۷۰/۱ |
|     | $\text{log}(\sigma)$ | year                                   |        |        |
| M2  | $\text{logit}(\mu)$  | ward+loglos+year                       | ۴۵۰۱/۰ | ۴۵۴۸/۱ |
|     | $\text{log}(\sigma)$ | year+ward                              |        |        |
| M3  | $\text{logit}(\mu)$  | ward+cs(loglos,df=1)+year              | ۴۴۷۹/۴ | ۴۵۳۱/۷ |
|     | $\text{log}(\sigma)$ | year+ward                              |        |        |
| M4  | $\text{logit}(\mu)$  | ward+cs(loglos,df=1)+year+cs(age,df=1) | ۴۴۷۸/۴ | ۴۵۴۱/۱ |
|     | $\text{log}(\sigma)$ | year+ward                              |        |        |

نتایج چهار مدل برازشی به داده‌ها در جدول ۸.۴ نمایش داده شده است. مدل M1 مدل منتخب گنگ و همکارانش است. در مدل M2 متغیر ward برای  $\sigma$  وارد مدل شده است و در مدل M3 و M4 برای مدل کردن  $\mu$  تابع اسپلاین مکعبی هموارساز با یک درجه آزادی استفاده شده است. با توجه به مقدار AIC مدل M4 بهترین مدل است و با اضافه کردن تابع هموارساز برای  $\mu$  در مدل M4 مدل بهبود یافته و مقدار AIC کاهش یافته است. طبق معیار BIC بهترین مدل، مدل M3 است که با اضافه کردن تابع هموارساز برای age در مدل M4 مقدار BIC افزایش یافته در نتیجه مدل مطلوب نمی‌باشد. در زیر نمودار تمام توابع برازشی برای  $\mu$  در مدل M4 نشان داده شده است.

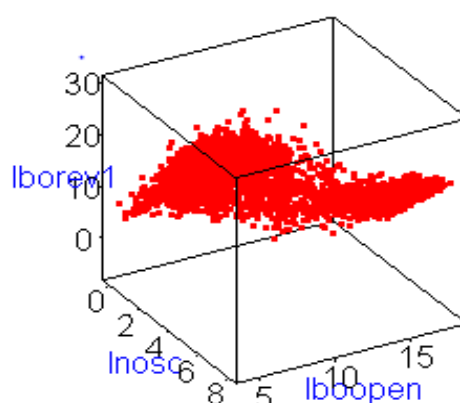

 شکل ۱۲.۴: توابع برازشی برای  $\mu$  در مدل M4

## ۲.۴ اسپلین دو بعدی

روش‌های زیادی برای درونیایی و مدل‌سازی سطوح به کار می‌رود. زمانی که بیش از یک متغیر مستقل موجود باشد، مساله‌ی برازش یک رویه‌ی چندبعدی پیش می‌آید و برای این که روند داده‌ها بهتر نشان داده شود، می‌توان روش‌های هموارسازی اسپلین را تعمیم داد. اسپلین حاصل ضرب تانسور و اسپلین صفحات نازک از جمله روش‌هایی برای مدل‌سازی سطوح است. داده‌های فیلم را برای نشان دادن ویژگی‌های توابع gamlss و به‌عنوان مثالی برای برازش سطوح هموارساز دوبعدی برای یک متغیر پاسخ پیوسته به کار می‌بریم.

### ۱.۲.۴ داده‌های فیلم

داده‌های فیلم توسط ویدوریس و همکارانش [۳۷] برای نشان دادن بسیاری از ویژگی‌های تابع gamlss و بررسی اثرات متقابل هموارساز متغیرهای توضیحی به کار برده شدند. داده‌های فیلم شامل ۴۰۱۵ مشاهده و چهار متغیر می‌باشد. لگاریتم تعداد پرده‌های سینما که فیلم نمایش داده شده است (lnosc)، لگاریتم سود هفتگی صندوق (lboopen) و عاملی است که نشان می‌دهد آیا توزیع کننده‌ی فیلم مستقل یا عمده می‌باشد (dist)، به‌عنوان متغیرهای توضیحی و لگاریتم سود هفتگی بعد از اولین هفته (lborev1)، متغیر پاسخ می‌باشد. داده‌ها در دو یا سه بعد قابل نمایش می‌باشند. با استفاده از بسته rgl می‌توان داده‌ها را در فضای سه بعدی ترسیم نمود.



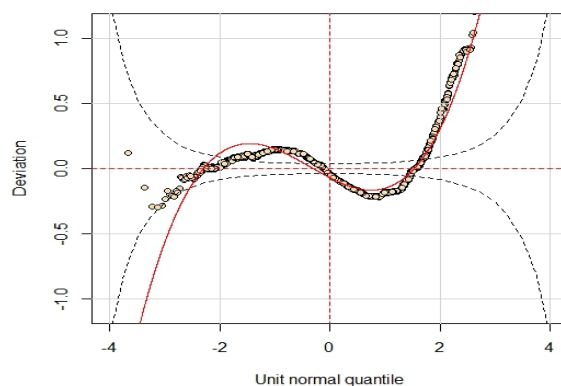
شکل ۱۳.۴: ترسیم نمودار سه بعدی داده‌های فیلم

بعد از تجزیه و تحلیل مقدماتی داده‌ها، با در نظر گرفتن توزیع نرمال برای متغیر پاسخ، مدل‌های برازشی به داده‌ها در جدول ۹.۴ نمایش داده شده است. در دو مدل اول برای هریک از متغیرهای توضیحی مدلی با اثر متقابل خطی، در دو مدل دوم هموارساز جمعی و در دو مدل سوم رویه‌ای هموار در نظر گرفته شده است. برای گزینش بهترین مدل از معیارهای AIC و BIC استفاده شده است.

جدول ۹.۴: گزینش بهترین مدل با استفاده معیارهای AIC و BIC

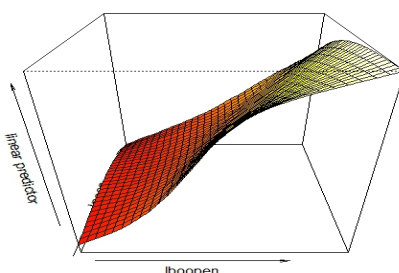
| مدل                     | توزیع متغیر پاسخ | پیش‌گوی جمعی                | AIC      | BIC      |
|-------------------------|------------------|-----------------------------|----------|----------|
| M1                      | normal           | boopen*lnosc                | ۱۲۲۲۶/۸۴ | ۱۲۲۵۸/۳۵ |
| مدل با اثرات متقابل خطی |                  |                             |          |          |
| M2                      | normal           | lboopen*lnosc+dist          | ۱۲۰۸۰/۹۹ | ۱۲۱۱۸/۸۰ |
| M3                      | normal           | pb(lnosc)+pb(lboopen)       | ۱۱۹۰۸/۷۳ | ۱۲۰۲۲/۹۶ |
| مدل جمعی                |                  |                             |          |          |
| M4                      | normal           | pb(lboopen)+ pb(lnosc)+dist | ۱۱۸۴۳/۷۸ | ۱۱۹۴۵/۳۹ |
| M5                      | normal           | ga(te(lboopen,lnosc))       | ۱۱۸۲۸/۵۹ | ۱۱۹۴۴/۰۶ |
| مدلی با برازش رویه      |                  |                             |          |          |
| M6                      | normal           | ga(te(lboopen,lnosc))+dist  | ۱۱۷۷۹/۷۶ | ۱۱۸۸۰/۶۹ |

با توجه به مقادیر به دست آمده در جدول ۹.۴ مدل M6 که رویه‌ای را برای متغیرهای توضیحی برازش می‌دهد به عنوان بهترین مدل برازشی می‌باشد. نتایج حاصله حاکی از این می‌باشد که در نظر گرفتن توزیع نرمال مناسب نمی‌باشد که در نمودار مارگون رسم شده‌ی مدل M6 کاملاً مشخص است. زیرا بیشتر نقاط خارج از بین دو نمودار بیضوی بالا و پایین می‌باشد، که این نشان دهنده‌ی عدم کفایت مدل می‌باشد.



شکل ۱۴.۴: نمودار مارگون برای مدل M6 با برازش رویه برای پارامتر  $\mu$

با استفاده از بسته mgcv می‌توان نمودار رویه‌ی برازشی برای مدل M6 را رسم نمود.



شکل ۱۵.۴: نمودار رویه‌ی برازشی برای مدل M6

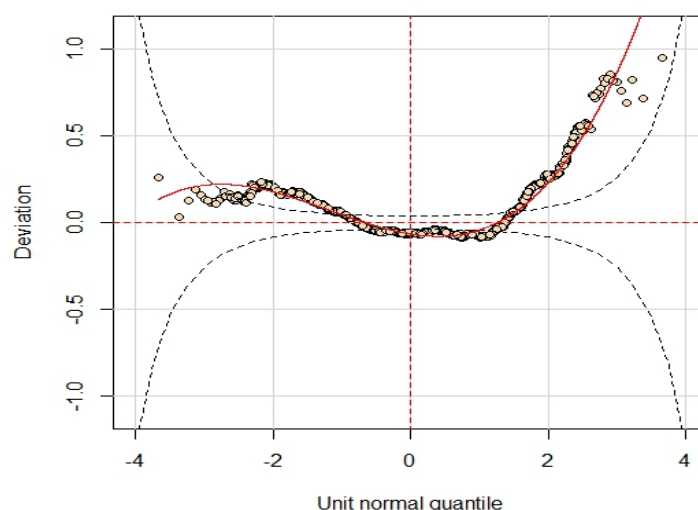
اکنون با مدل کردن پارامتر مقیاس به عنوان تابعی از متغیرهای توضیحی می‌توان بهبود یافتن مدل را بررسی نمود.

جدول ۱۰.۴: گزینش بهترین مدل با استفاده معیارهای AIC

| مدل | توزیع متغیر پاسخ | پیش‌گوی جمعی  | AIC      | BIC      |
|-----|------------------|---|----------|----------|
| M7  | normal           | $ga(te(lboopen, lnosc)) + dist$<br>$sigma.fo = ga(te(lboopen, lnosc)) + dist$ | ۱۰۲۱۸/۱۲ | ۱۰۰۴۳/۸۹ |

میزان AIC مدل M7 نسبت به مدل M6 مقدار کمتری است، در نتیجه مدل بهتری نسبت به

مدل M6 می‌باشد. برای بررسی کفایت مدل، نمودار مارگون مدل M7 رسم شده است و مشهود می‌باشد که این مدل کفایت لازم را ندارد و نمی‌تواند متغیر پاسخ را به‌خوبی توضیح دهد.



شکل ۱۶.۴: نمودار مارگون برای مدل M7 با برازش رویه برای پارامترهای  $\mu$  و  $\sigma$

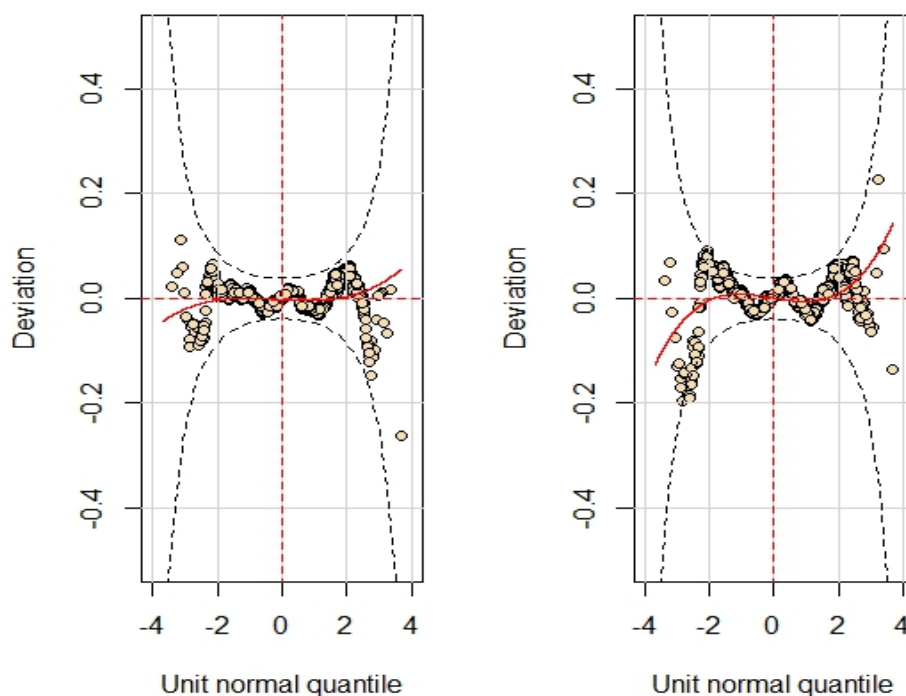
اکنون پارامتر پاسخ را با استفاده از توزیع BCPE<sup>۳</sup> که توزیعی از خانواده‌ی توزیع‌های موجود در بسته gamlss.dist می‌باشد برازش می‌دهیم. توزیعی چهار پارامتری است که پارامترهای  $\mu$  و  $\nu$  به‌طور پیش فرض دارای تابع پیوندی همانی و پارامترهای  $\sigma$  و  $\tau$  دارای تابع پیوندی لگاریتمی می‌باشد. در مدل M8 به هر چهار پارامتر مدل جمعی و در مدل M9 رویه‌ی هموارساز برازش یافته و در جدول ۱۱.۴ نمایش داده‌شده است.

<sup>۳</sup>BoxCox power exponential

جدول ۱۱.۴: گزینش بهترین مدل با استفاده معیارهای AIC

| AIC      | پیش‌گوی جمعی                            | تابع پیوند |       |          |       | مدل |
|----------|---|------------|-------|----------|-------|-----|
|          |   | $\tau$     | $\nu$ | $\sigma$ | $\mu$ |     |
| ۹۹۸۰/۹۴۸ | pb(lboopen)+pb(lnosc) + dist            |            |       |          |       | M8  |
|          | sigma.fo = pb(lboopen)+pb(lnosc) + dist | log        | ident | log      | ident |     |
|          | nu.fo = pb(lboopen)+pb(lnosc) + dist    |            |       |          |       |     |
|          | tau.fo = pb(lboopen)+pb(lnosc) + dist   |            |       |          |       |     |
| ۹۹۸۰/۹۴۸ | ga(te(lboopen,lnosc)) + dist            |            |       |          |       | M9  |
|          | sigma.fo =ga(te(lboopen,lnosc)) + dist  | log        | ident | log      | ident |     |
|          | nu.fo =ga(te(lboopen,lnosc)) + dist     |            |       |          |       |     |
|          | tau.fo =ga(te(lboopen,lnosc)) + dist    |            |       |          |       |     |

با توجه به مقادیر معیارها، عملکرد مدل M9 کمی بهتر است، با توجه به شکل ۱۷.۴ نمودار مارگون ترسیمی دو مدل فوق به نظر می‌رسد که مدل M8 بهتر از مدل M9 است.

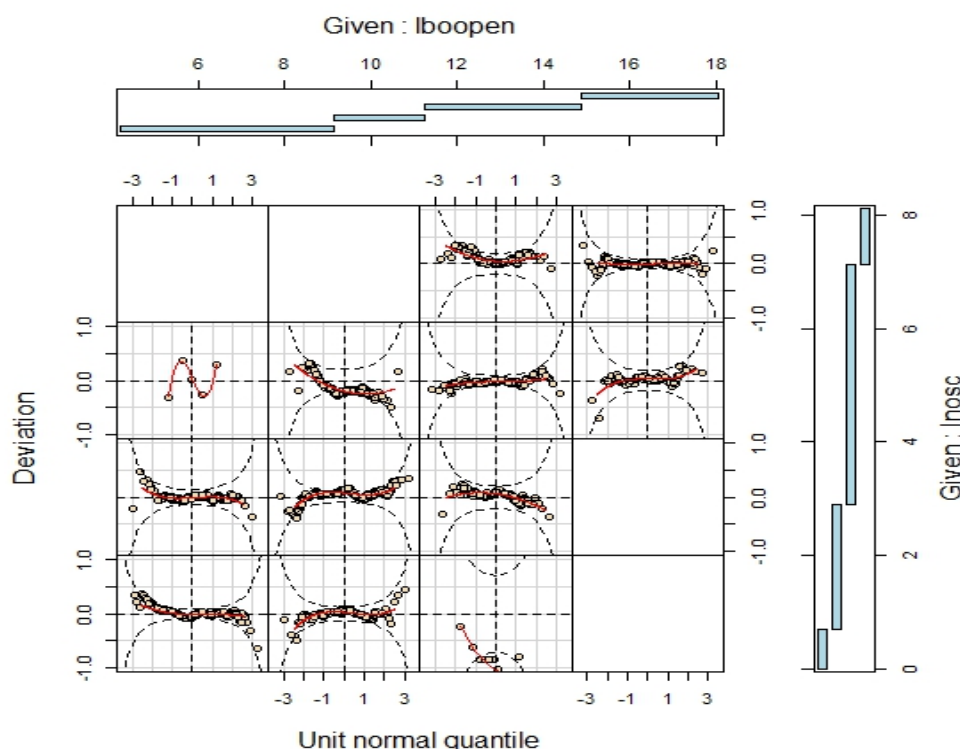


شکل ۱۷.۴: نمودار مارگون قاب سمت چپ مدل M8 و قاب سمت راست مدل M9

تحت این شرایط انتخاب مدل بهتر، دشوار می‌باشد، لذا می‌توان در نواحی مشترک دو متغیر



Iboopen و Inosc مدلی را برازش داد و نمودار مارگون برای هر دو متغیر توضیحی را به کار برد که در زیر رسم شده است.



شکل ۱۸.۴: نمودار مارگون برای مدل M9 با دو متغیر توضیحی

چهار ستون مربوط به چهار ناحیه از متغیر Iboopen در بالای نمودار و چهار ردیف مربوط به چهار ناحیه از متغیر Inosc در سمت راست نمودار رسم شده است. در داخل نمودار، ۱۶ نمودار مارگون در ناحیه‌های پیوستگی دو متغیر وجود دارد. در بعضی از ناحیه‌ها مشاهداتی موجود نمی‌باشد. به طور کلی نمودارهای مارگون در ناحیه‌های پیوستگی نشان دهنده‌ی کفایت مدل برازشی است.

### ۳.۴ مدل میدان‌های تصادفی گاوسی

مدل میدان‌های تصادفی گاوسی (GMRF) به طور گسترده‌ای در آمار فضایی استفاده می‌شود. این مدل زمانی که یک عامل اطلاعات جغرافیایی را دربر دارد بسیار مناسب می‌باشند. به عنوان مثال اطلاعات مکان یا ناحیه‌ای مشخص، مناطق، .... این داده‌ها علاوه بر این که معرف اطلاعات ناحیه‌ای مشخص هستند، بین آن‌ها وابستگی وجود دارد. مشاهداتی که در یک همسایگی قرار دارند، مشابه‌تر هستند. لذا به کارگیری روش‌های معمول تحلیل که در آن‌ها فرض بر استقلال مشاهدات است، برای تحلیل این گونه داده‌ها که داده‌های فضایی شبکه‌ای

نامیده می‌شوند، مناسب نمی‌باشد. عموماً از یک مدل میدان تصادفی مارکوفی گاوسی که یک خانواده مهم از توزیع‌ها است برای مدل‌بندی چنین داده‌هایی استفاده می‌شود. یک مدل (GMRF) را می‌توان به عنوان یک مدل اسپلاین جریمه‌ای با توجه به معادلات (۱۳.۲) (۱۴.۲) در نظر گرفت که در آن ماتریس پایه  $B$  ماتریسی شاخص است که مشاهدات آن مربوط به منطقه فضایی می‌باشد و ماتریس  $D^T D$  حاوی اطلاعات همسایگی فضایی می‌باشد. در نرم افزار R با استفاده از بسته `gamlss.spatia` دستور (GMRF) برای برازش مدل‌های فضایی به کار گرفته می‌شود. برای نشان دادن کاربرد مدل‌های (GMRF) داده‌های جرم و جنایت (columb) در شهر کلمبوس ایالت اوهایو را به کار می‌بریم.

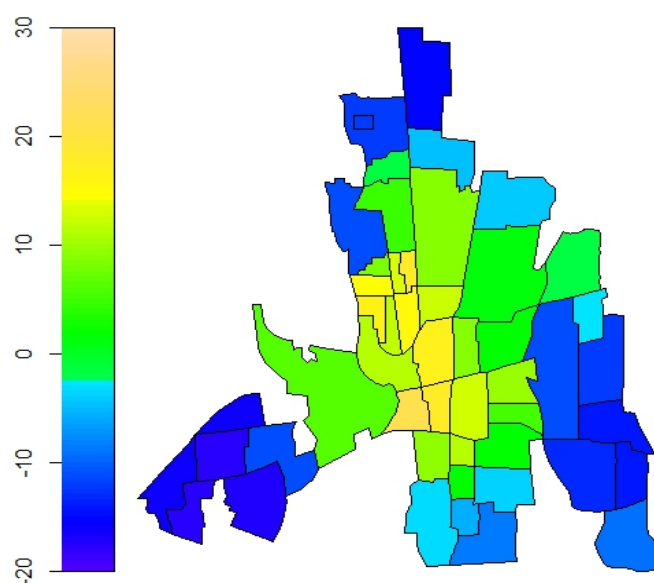
### ۱.۳.۴ داده‌های جرم و جنایت

داده‌های جرم تعداد ۴۹ مشاهده و دارای متغیرهای توضیحی، منطقه مسکونی (area)، قیمت مسکن (home value)، درآمد خانوار (income)، میزان فضای باز ناحیه (open.space)، کد شناسایی منطقه و مطابقت نام برای مجموعه داده‌ی `columb.polys` که حاوی چند ضلعی‌هایی برای مناطق می‌باشد (district) و سرقت‌های مسکونی و خودکار در هر ۱۰۰۰ خانوار (crime) به عنوان متغیر پاسخ می‌باشد. در جدول ۱۲.۴ مدل برازشی به داده‌های نمایش داده شده است.

جدول ۱۲.۴: مدل برازشی به داده‌های `columb` با استفاده از تابع `gmrf`

| مدل   | پیش‌گوی جمعی |
|---|--------------|
| مدل جمعی تعمیم‌یافته برای مکان، شکل و مقیاس <code>gmrf(district, polys=columb.polys)</code> |              |

در شکل ۱۹.۴ میزان جرم با توجه به طیف رنگی قابل تفسیر می‌باشد. نواحی زرد رنگ نشان دهنده‌ی منطقه‌ای است که بیشترین میزان جرم و جنایت در آن اتفاق افتاده است و نواحی آبی رنگ بیانگر کاهش میزان جرم و جنایت در آن منطقه می‌باشد.



شکل ۱۹.۴: برازش مقادیر برای داده‌های جرم جنایت با استفاده از تابع GMRF

## ۴.۴ نتیجه‌گیری

در این پایان‌نامه، به تفصیل، انواع مختلف اسپلاین‌ها را معرفی کردیم و روش‌های انتخاب پارامتر هموارساز را مطرح کردیم. یکی از جدی‌ترین معایب اسپلاین‌هایی که از توابع پایه ساخته می‌شوند، ناپایداری عددی است که با افزایش درجه چندجمله‌ای، هم‌خطی بین توابع بریده‌شده ایجاد می‌شود. بنابراین در این موارد، از رابطه بازگشتی برای ساختن توابع پایه استفاده می‌شود. در رابطه با اسپلاین‌های هموارساز یکی از معایب آن‌ها وابستگی به تعداد گره‌ها است که با افزایش تعداد گره‌ها با مشکلاتی از قبیل بیش‌برازشی و محاسبات سنگین مواجه می‌شویم. برای رفع این مشکل نیز از اسپلاین‌های جریمه‌ای استفاده می‌شود و همچنین مدل‌های `gamlss` را معرفی کردیم که به منظور تحلیل و مدل‌سازی مورد استفاده قرار می‌گیرند که مدل‌های رگرسیونی و برازش‌های بسیار انعطاف‌پذیری را با کاربرد توابع ناپارامتری از جمله اسپلاین‌ها فراهم می‌کند با معیارهای انتخاب مدل هم‌چون  $AIC, CV, \dots$  بهترین مدل برازشی را انتخاب می‌کنیم که بتواند داده‌ها را به‌خوبی توصیف نماید و دارای خطای کمتری باشد.

# پیوست آ

## ۱. آ ضرب کرونکر

ضرب کرونکر با نماد  $\otimes$  نشان داده می‌شود، دارای خواص جالبی است که باعث شده در زمینه‌های مختلف به‌طور گسترده مورد استفاده قرار گیرد. اولین بار توسط یوهان گئورگ زفوس استفاده شد که به‌عنوان ضرب تانسوری نیز شناخته می‌شود. فرض کنید  $\mathbb{R}$  میدان اعداد مختلط و  $\mathbb{R}^{m \times n}$  مجموعه ماتریس‌های شامل  $m$  سطر و  $n$  ستون با درایه‌های مختلط باشد [۱۵].

**تعریف ۱.۱.۱.** ضرب کرونکر دو ماتریس  $A \in \mathbb{R}^{m_A \times n_A}$  و  $B \in \mathbb{R}^{m_B \times n_B}$  به‌صورت  $A \otimes B$  نشان داده می‌شود و برابر است با

$$A \otimes B = \begin{pmatrix} a_{11}B & \cdots & a_{1n_A}B \\ \vdots & \ddots & \vdots \\ a_{m_A1}B & \cdots & a_{m_An_A}B \end{pmatrix}$$

هر  $a_{ij}B$  یک بلوک در سائز  $m_B \times n_A$  است بنابراین سائز  $A \otimes B$  برابر  $m_A m_B \times n_A n_B$  است.

**تعریف ۲.۱.۱.** جمع کرونکر ماتریس‌های مربعی  $A \in \mathbb{R}^{n_A \times n_A}$  ،  $B \in \mathbb{R}^{n_A \times n_A}$  به صورت  $A \otimes B$  نمایش داده می‌شود و برابر است با

$$A \otimes B = A \otimes I_{n_B} + I_{n_A} \otimes B$$

**خواص پایه:**

شرکت پذیری

$$A \otimes (B \otimes C) = (A \otimes B) \otimes C$$

توزیع پذیری ضرب به جمع

$$(A + B) \otimes (C + D) = A \otimes C + B \otimes C + A \otimes D + B \otimes D$$

سازگاری با معکوس ماتریس

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

سازگاری با ترانهادی ماتریس

$$(A \otimes B)^T = A^T \otimes B^T$$

یکی از مهمترین خواص ضرب کرونکر ارتباط آن با ضرب ماتریس است. برای ماتریس‌های  $n_B = m_D$  ،  $n_A = m_C$  و اگر  $C \in \mathbb{R}^{m_C \times n_C}$  ،  $D \in \mathbb{R}^{m_D \times n_D}$  ،  $B \in \mathbb{R}^{m_B \times n_B}$  ،  $A \in \mathbb{R}^{m_A \times n_A}$  آن‌گاه

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$$

## ۲. مشتق تابع B-اسپلاین

دی بور (۱۹۷۸) فرمولی ساده برای مشتق توابع B-اسپلاین ارائه نمود.

$$h \sum_j \alpha_j \beta_j'(x, k) = \sum \alpha_j \beta_j(x, k-1) - \sum \alpha_{j-1} \beta_{j-1}(x, k-1) = \sum \Delta \alpha_j \beta_j(x, k-1)$$

که  $\Delta \alpha_j = \alpha_j - \alpha_{j-1}$  تفاضل مرتبه اول و  $h$  فاصله‌ی بین گره‌ها است. با توجه به این مقدمات مشتق دوم را به دست می‌آوریم:

$$\begin{aligned} h^2 \sum_j \alpha_j \beta_j''(x, k) &= (h \sum_j \alpha_j \beta_j'(x, k))' = (\sum \alpha_j \beta_j(x, k-1) - \sum \alpha_{j-1} \beta_{j-1}(x, k-1))' = \\ &= (\sum \alpha_j \beta_j(x, k-2) - \sum \alpha_{j-1} \beta_{j-1}(x, k-2)) - (\sum \alpha_{j-1} \beta_{j-1}(x, k-2) - \sum \alpha_{j-2} \beta_{j-2}(x, k-2)) = \\ &= \sum \alpha_j \beta_j(x, k-2) - 2 \sum \alpha_{j-1} \beta_{j-1}(x, k-2) + \sum \alpha_{j-2} \beta_{j-2}(x, k-2) = \\ &= \sum \alpha_j - 2 \alpha_{j-1} + \alpha_{j-2} \beta_j(x, k-2) = \sum \Delta^2 \alpha_j \beta_j(x, k-2) \end{aligned}$$

$$\Delta^2 \alpha_j = \alpha_j - 2\alpha_{j-1} + \alpha_{j-2}$$

تفاضل مرتبه دوم  $\alpha_{j-2} + \alpha_{j-1} - \alpha_j$

## ۳.آ جدول توزیع‌های پیوسته و گسسته موجود در بسته‌ی

### `gamlss.dist`

جدول ۱.آ: خلاصه‌ای از توزیع‌های پیوسته موجود در بسته‌ی `gamlss.dist` با تابع پیوند پیش‌فرض

| توزیع                        | تابع چگالی احتمال  | تابع پیوند  |
|------------------------------|--|---|
|                              |  | $\mu \quad \sigma \quad \nu \quad \tau$                 |
| باکس کاکس کول و گرین (BCCGo) | $f(y \mu, \sigma, \nu) = \left(\frac{1}{\sqrt{(\pi)\sigma}}\right) \left(\frac{y^{\nu-1}}{\mu^\nu}\right) \exp\left(-\frac{z^2}{\nu}\right)$ $z = \begin{cases} \frac{(y\mu)^{(\nu)-1}}{(\nu\sigma)} & \nu \neq 0 \\ z = \frac{\log(\frac{y}{\mu})}{\sigma} & \nu = 0 \end{cases}$   | $\log \quad \text{ident} \quad \log \quad \log$         |
| باکس کاکس توان نمایی (BCPE)  | $f(y \mu, \sigma, \nu, \tau) = \left(\frac{y^{\nu-1}}{\mu^\nu}\right) \left(\frac{\tau}{\sigma}\right) \frac{\exp(-1/\tau   \frac{z}{\sigma}  ^\tau)}{(c^2(1+\frac{1}{\tau}))\Gamma(\frac{1}{\tau})}$ $c = \left[\frac{2^{(\frac{1}{\tau})}\Gamma(\frac{1}{\tau})}{\Gamma(\frac{1}{\tau})^{(\frac{1}{\tau})}}\right]$ $z = \frac{\log(\frac{y}{\mu})}{\sigma}$ | $\log \quad \text{ident} \quad \log \quad \text{ident}$ |
| گاما (GA)                    | $f_z(z \mu, \sigma) = \frac{1}{(\mu\sigma^2)^{1/\sigma^2}} \frac{z^{(\frac{1}{\sigma^2}-1)} e^{-z/(\mu\sigma)^2}}{\Gamma(1/(\sigma)^2)}$   | $\log \quad \log$                                       |
| گاوسی وارون (IG)             | $f(z \mu, \sigma) = \frac{1}{(\sqrt{2\pi\sigma^2 y^2} \exp(-(y-\mu)^2/(2\mu^2\sigma^2 y)))}$   | $\log \quad \log$                                       |

جدول ۲.آ: خلاصه‌ای از توزیع‌های گسسته موجود در بسته‌ی gamlss.dist با تابع پیوند پیش‌فرض

| توزیع                        | تابع جرم احتمال  | تابع پیوند   |
|------------------------------|--|--|
| پواسن (PO)                   | $f(y \mu) = \frac{e^{-\mu}\mu^y}{\Gamma(y+1)}$   | $\tau \quad \nu \quad \sigma \quad \mu$<br>• • • log |
| پواسن دوگانه (DPO)           | $f(y \mu, \sigma) = \left(\frac{1}{\sigma}\right)^{\frac{1}{\sigma}} \left[\frac{e^{-y}y^y}{y!}\right] \left[\frac{(e\mu)}{y}\right]^{\frac{y}{\sigma}} C$   | • • log log  |
| دوجمله‌ای منفی نوع یک (NBI)  | $f(y \mu, \sigma) = \frac{\Gamma(y+\frac{1}{\sigma})}{\Gamma(y+1)\Gamma(\frac{1}{\sigma})} \frac{(\mu\sigma)^y}{(\mu\sigma+1)^{(y+(\frac{1}{\sigma}))}}$   | • • log log  |
| دوجمله‌ای منفی نوع دو (NBII) | $f(y \mu, \sigma) = \frac{\Gamma(y+\frac{\mu}{\sigma})\sigma^y}{\Gamma(\frac{\mu}{\sigma})\Gamma(y+1)(1+\sigma)^{y+\frac{\mu}{\sigma}}}$   | • • log log  |
| پواسن گاوسی وارون (PIG)      | $f(y \mu, \sigma) = \left(\frac{\gamma\alpha}{\pi}\right)^{\frac{1}{\gamma}} \mu^y e^{(\frac{1}{\sigma})} \frac{K(\alpha)}{(\alpha\sigma)^y y!}$<br>$\alpha^{\frac{1}{\gamma}} = \frac{1}{\sigma^{\frac{1}{\gamma}}} + \frac{\gamma\mu}{\sigma}$<br>$K_{\lambda}(t) = \frac{1}{\gamma} \int_0^{\infty} x^{\lambda-1} \exp\{-\frac{1}{\gamma}t(x+x^{-1})\}dx$ | • • log log  |
| دلپورت (DEL)                 | $f(y \mu, \sigma, \nu) = \left(\frac{\exp(-\mu\nu)}{\Gamma(\frac{1}{\sigma})}\right) \{1 + \mu\sigma(1-\nu)\}^{(\frac{1}{\sigma})} S$<br>$S = \sum (P(y, j) ((\mu^y) (\frac{\nu^{y-j}}{y!}) \left\{1 + (\frac{1}{\sigma(1-\nu)})\right\}^j \Gamma((\frac{1}{\sigma})j))$   | • logit log log                                      |
| سیچل (SI)                    | $f(y \mu, \sigma, \nu) = \frac{\mu^y K_{y+n}(\alpha)}{(\alpha\sigma)^{(y+v)} y! K_{\nu}(\frac{1}{\sigma})}$<br>$\alpha^{\frac{1}{\gamma}} = \frac{1}{\sigma^{\frac{1}{\gamma}}} + \frac{\gamma\mu}{\sigma}$<br>$K_{\lambda}(t) = \frac{1}{\gamma} \int_0^{\infty} x^{\lambda-1} \exp\{-\frac{1}{\gamma}t(x+x^{-1})\}dx$                                      | • logit log log                                      |

## ۴.آ دستورات نرم‌افزار R

شکل 1.2

```
rm(list=ls())
plus <- function (x) {
  ifelse( x >= 0, x, 0 )
}
plus <- function (x) {
  ifelse( x >= 0, x, 0 )
}
op <- par(mfrow=c(1,2))
plot(0,type="n",xlim = c(-1, 2), ylim = c(-0.1, 0.5),xlab="x",ylab="pluse(x-a)")
for( a in c(-.5,1,0,.5)){
  curve( plus(x - a),
        lwd = 3,
```

```

        type="l"      , add=T)
    abline(h=0, v=c(-0.5,0,0.5,1), lty=3)
}
col = "grey60"
plot(0,type="n",xlim = c(-1, 2), ylim = c(-0.1, 0.5),xlab="x",ylab="pluse(x-a)^3")
for( a in c(-.5,1,0,.5)){
  curve( plus(x - a)^5,
        lwd = 3,
        col = "grey30", add=T)
  abline(h=0, v=c(-0.5,0,0.5,1), lty=3)
}

```

## شکل 2.2

```

library ( ISLR)
attach (Wage )
agelims = range ( age )
age.grid =seq ( from= agelims [1], to= agelims [2])
library(splines)
plot(age ,wage , col =" gray ")
fit =lm(wage~ bs(age , knots =c(25 ,40 ,60) ),data= Wage)
pred= predict(fit , newdata = list(age =age.grid),se=T)
plot(age ,wage , col =" gray ")
lines (age.grid ,pred$fit ,col ="blue ",lwd =2)
lines (age.grid , pred$fit +2* pred$se,col ="blue" ,lty =2)
lines (age.grid ,pred$fit -2* pred$se ,col ="blue",lty =2)
abline(v=c(25 ,40, 60 ),lty=2)

fit2=lm( wage~ ns(age ,df =4) ,data= Wage)
pred2 = predict (fit2 , newdata =list (age =age.grid),se=T)
lines (age.grid , pred2$fit ,col =" red ",lwd =2)
lines (age.grid , pred2$fit + pred2$se,col =" red ",lty=2)
lines (age.grid ,pred2$fit -2* pred$se ,col =" red ",lty=2)

legend("topright", inset=0.02, legend=c("Natural Cubic Spline ","Cubic Spline"),

```



---

```
lwd=2, col=c("red", "blue"))
```

### شکل 3.2

```
library(splines)
x <- seq(0, 2, length=50)
B=bs(x, knots=1,Boundary.knots=c(0,2), intercept=TRUE, degree=1)
matplot(x, B,type = "l", col =2:5, lwd = 2)
matplot(x, B,type = "l", col =c("black","red","blue"), lwd = 2)
```

### شکل 4.2

```
n=10
xr=seq(0,n,by=0.1)
yr=sin(xr/2)+rnorm(length(xr))/2
db=data.frame(x=xr,y=yr)
plot(db)
par(mfrow=c(3,1))
B=bs(xr,knots=1:9 ,Boundary.knots=c(0,10),degree=1)
matplot(xr,B,type="l")
title(main="Bsplines of Order 1",font.main=1)
B=bs(xr,knots=1:9,Boundary.knots=c(0,10),degree=2)
matplot(xr,B,type="l")
title(main="Bsplines of Order 2",font.main=1)
B=bs(xr,knots=1:9,Boundary.knots=c(0,10),degree=3)
matplot(xr,B,type="l")
title(main="Bsplines of Order 3",font.main=1)
```

### شکل 5.2

```
library(monreg)
op <- par(mfrow=c(1,2))
y <- c(1:3,5,9,7:3,2*(2:5),rep(15,4))
plot(y, col.main = 1)
x <- seq(1,length(y), len=201)
s02 <- smooth.spline(y, spar = 0.1)
```

```
lines(predict(s02, x), col = 2)
title(main="spar=0.1")
y <- c(1:3,5,9,7:3,2*(2:5),rep(15,4))
plot(y, col.main = 2)
x <- seq(1,length(y), len=201)
s02 <-smooth.spline(y, spar = 0.8)
lines(predict(s02, x), col = 2)
title(main="spar=0.8")
```

شکل 6.2

```
library(MASS) # data sets
library(splines) # B-splines
library(mgcv)
knots_eq <- function(x, k, m)
{
  c(min(x) - ((k-1):0) * (max(x)-min(x))/(m+1),
    seq(from=min(x), to=max(x), length.out=m+2)[-c(1,m+2)],
    max(x) + (0:(k-1)) * (max(x)-min(x))/(m+1))
}
# functions for polynomial regression
polynom <- function(x, p)
{
  if (p==0) { return("1") }
  if (p >0) {
    return(paste(polynom(x=x, p=p-1), " + I(", x, "^",
      p, ")", sep="")) }
}
polynom_m <- function(x, p)
{
  if (p==0) { return(as.matrix(rep(1,length(x)))) }
  if (p >0) { return(cbind(polynom_m(x=x,p=p-1),x^p)) }
}
# functions for information criteria
AIC <- function(res, H)
```

---

```

{
log(sum(res^2)/length(res)) + 2*sum(diag(H))/length(res)
}
SIC <- function(res, H)
{
log(sum(res^2)/length(res)) +sum(diag(H))*log(length(res))/length(res)
}
CV <- function(res, H)
{
sum((res/(1-diag(H)))^2)
}
GCV <- function(res, H)
{
sum(res^2)/(1-sum(diag(H))/length(res))^2
}
normed <- function(x) { (x-min(x))/(max(x) - min(x)) }
times <- mcycle$times
accel <- mcycle$accel

k <- 4
m <- 20
D <- matrix(0, nrow=m+k-2, ncol=m+k)
for (j in 1:(m+k-2))
{
d_j <- c(rep(0,j-1),1,-2,1,rep(0,(m+k)-3-(j-1)))
e_j <- c(rep(0,j-1), 1 ,rep(0,(m+k)-3-(j-1)))
D <- D + e_j%*%t(d_j)
}
D
y_star <- c(accel, rep(0,m+k-2))
y_star
X_spl <- splineDesign(x=times, knots=knots_eq(times,k,m), ord=k)
X_spl
ls <-c(0,1,50,10^10)

```

```
cols_1 <- rainbow(length(ls))
cols_1
par(mai=c(0.65,0.6,0.1,0.1), mgp=c(2,1,0))
plot(times, accel, xlab="time", ylab="acceleration", pch=16)
for (ll in 1:length(ls))
{
  l <- ls[ll]
  X_star <- rbind(X_spl, sqrt(l)*D)
  est <- lm(y_star ~ -1 + X_star)
  plot(function(x) splineDesign(x=x, knots=knots_eq(times, k, m), ord=k) %*%
  est$coef,
  from=min(times), to=max(times), add=TRUE, lwd=2, col=cols_1[ll])
}
legend("bottomright", inset=0.02, col=cols_1, lwd=2, ncol=2, cex=0.8,
legend=parse(text=paste("lambda==", ls, sep=""))))
```

شکل 1.3

```
lambdas <- c(10^{-5}, 1:400/100)
AICs <- rep(NA, length(lambdas))
SICs <- AICs
GCVs <- AICs
CVs <- AICs
for (ll in 1:length(lambdas))
{
  l <- lambdas[ll]
  X_star <- rbind(X_spl, sqrt(l)*D)
  est <- lm(y_star ~ -1 + X_star)
  Hat <- X_spl %*% solve(t(X_star)%*%X_star) %*% t(X_spl)
  Res <- accel - X_spl%*% est$coef
  AICs[ll] <- AIC(res=Res, H=Hat)
  SICs[ll] <- SIC(res=Res, H=Hat)
  GCVs[ll] <- GCV(res=Res, H=Hat)
  CVs[ll] <- CV(res=Res, H=Hat)
}
```

---

```

par(mai=c(0.65,0.4,0.1,0.1), mgp=c(2,1,0))
plot(range(lambdas), 0:1, type="n", yaxt="n", xlab=expression(lambda),ylab="")
axis(2, at=0.5, label="(scaled) AIC, SIC, CV, GCV", tick=FALSE, line=-0.5)
SCs <- c("AICs", "SICs", "CVs", "GCVs")
for (SC in 1:length(SCs))
{
points(lambdas, normed(get(SCs[SC])),
type="l", lwd=2, col=rainbow(4)[SC])
abline(v=lambdas[which(get(SCs[SC])==min(get(SCs[SC])))],col=1)
}
legend("topright", inset=0.02,
legend=c(paste("AIC (", lambdas[which(AICs==min(AICs))], ")", sep=""),
paste("SIC (", lambdas[which(SICs==min(SICs))], ")", sep=""),
paste("CV (", lambdas[which(CVs ==min(CVs ))], ")", sep=""),
paste("GCV (", lambdas[which(GCVs==min(GCVs))], ")", sep="")),
col=rainbow(4), lwd=2, bg="white")

```

شکل 2.3 و 3.3

```

library(MASS) # data sets
library(splines) # B-splines
library(mgcv)
knots_eq <- function(x, k, m)
{
c(min(x) - ((k-1):0) * (max(x)-min(x))/(m+1),
seq(from=min(x), to=max(x), length.out=m+2)[-c(1,m+2)],
max(x) + (0:(k-1)) * (max(x)-min(x))/(m+1))
}
# functions for polynomial regression
polynom <- function(x, p)
{
if (p==0) { return("1") }
if (p >0) {
return(paste(polynom(x=x, p=p-1), " + I(", x, "^",
p, ")", sep="")) }
}

```

```

}

polynom_m <- function(x, p)
{
  if (p==0) { return(as.matrix(rep(1,length(x)))) }
  if (p >0) { return(cbind(polynom_m(x=x,p=p-1),x^p)) }
}

# functions for information criteria
AIC <- function(res, H)
{
  log(sum(res^2)/length(res)) + 2*sum(diag(H))/length(res)
}

SIC <- function(res, H)
{
  log(sum(res^2)/length(res)) +sum(diag(H))*log(length(res))/length(res)
}

CV <- function(res, H)
{
  sum((res/(1-diag(H)))^2)
}

GCV <- function(res, H)
{
  sum(res^2)/(1-sum(diag(H))/length(res))^2
}

normed <- function(x) { (x-min(x))/(max(x) - min(x)) }

times <- mcycle$times
accel <- mcycle$accel
m <- 3
ks <-4

#cols_k <- rainbow(length(ks))
par(mai=c(0.65,0.6,0.1,0.1), mgp=c(2,1,0))
plot(times, accel, xlab="time", ylab="acceleration", pch=16)
abline(v=knots_eq(times, 1, m), lty=c(1,rep(2,m),1), col="grey60")
for (kk in 1:length(ks))
{

```

---



---

```

k <- ks[kk]
est <- lm(accel ~ -1 +
splineDesign(x=times, knots=knots_eq(times, k, m), ord=k))
plot(function(x) splineDesign(x=x, knots=knots_eq(times, k, m), ord=k) %*%
est$coef,
from=min(times), to=max(times), n=1001, lwd=2, col=cols_k[kk], add=TRUE)
}
legend("bottomright", inset=0.02, legend=c(paste("k = ",ks,sep=""), "knots=3"),
lwd=2, lty=c(rep(1,length(ks)),2), col=c(cols_k, "grey60"), bg="white")
m <-30
ks <- 4
#cols_k <- rainbow(length(ks))
par(mai=c(0.65,0.6,0.1,0.1), mgp=c(2,1,0))
plot(times, accel, xlab="time", ylab="acceleration", pch=16)
abline(v=knots_eq(times, 1, m), lty=c(1,rep(2,m),1), col="grey60")
for (kk in 1:length(ks))
{
k <- ks[kk]
est <- lm(accel ~ -1 +
splineDesign(x=times, knots=knots_eq(times, k, m), ord=k))
plot(function(x) splineDesign(x=x, knots=knots_eq(times, k, m), ord=k) %*%
est$coef,
from=min(times), to=max(times), n=1001, lwd=2, col=cols_k[kk], add=TRUE)
}
legend("bottomright", inset=0.02, legend=c(paste("k = ",ks,sep=""), "knots=30"),
lwd=2, lty=c(rep(1,length(ks)),2), col=c(cols_k, "grey60"), bg="white")

```

### شكل 4.3

```

library ( ISLR)
attach (Wage )
agelims = range ( age )
age.grid =seq ( from= agelims [1], to= agelims [2])
library(splines)
plot(age ,wage , col =" gray ")

```

```
fit2=lm( wage~ ns(age ,df =4) ,data= Wage)
pred2 = predict (fit2 , newdata =list (age =age.grid),se=T)
lines (age.grid , pred2$fit ,col =" red ",lwd =2)
abline(v=c(33.75 ,42.00, 51.00 ),lty=2)
lines (age.grid , pred$fit +2* pred$se,col =" red ",lty=2)
lines (age.grid ,pred$fit -2* pred$se ,col =" red ",lty=2)
dim (ns(age , knots =c(25 ,40 ,60) ))
dim (ns(age ,df =4))
attr(ns(age,df =4) ,"knots")
```

شکل 5.3

```
set.seed(11)
x <- runif(100,-2,2)
# generate the y response
y <- 2*x^3 + x^2 - 2*x +5 + rnorm(100)
plot(x,y)
xy <- data.frame(x=x, y=y)
library(ggplot2)
# specify the maximum polynomial degree that will be explored
max.poly <- 7
# cretaing data.frame which will store model predictions
# that will be used for the smooth curves in Fig. 1
x.new <- seq(min(x), max(x), by=0.1)
degree <- rep(1:max.poly, each=length(x.new))
predicted <- numeric(length(x.new)*max.poly)
new.dat <- data.frame(x=rep(x.new, times=max.poly),
                      degree,
                      predicted)
# fitting lm() polynomials of increasing complexity
# (up to max.degree) and storing their predictions
# in the new.dat data.frame
for(i in 1:max.poly)
{
  sub.dat <- new.dat[new.dat$degree==i,]
```



---

```

new.dat[new.dat$degree==i,3] <- predict(lm(y~poly(x, i)),
                                         newdata=data.frame(x=x.new))
}
# plotting the data and the fitted models
p <- ggplot()
p + geom_point(aes(x, y), xy, colour="darkgrey">#darkgrey
p + geom_line(aes(x, predicted,colour=as.character(degree)) , new.dat)
p + scale_colour_discrete(name = "Degree")
p

```

### شكل 6.3

```

# creating empty data.frame that will store
# AIC and BIC values of all of the models
AIC.BIC <- data.frame(criterion=c(rep("AIC",max.poly),
                                   rep("BIC",max.poly)),
                     value=numeric(max.poly*2),
                     degree=rep(1:max.poly, times=2))
# calculating AIC and BIC values of each model
for(i in 1:max.poly)
{
  AIC.BIC[i,2] <- AIC(lm(y~poly(x,i)))
  AIC.BIC[i+max.poly,2] <- BIC(lm(y~poly(x,i)))
}
a <- data.frame(cross=numeric(max.poly))
for(i in 1:max.poly)
{
  a[i,1] <- crossvalidate(x, y, degree=i)
}
# plotting AIC and BIC against model complexity
# (which is the polynomial degree)
AIC.plot <- qplot(degree, value, data=AIC.BIC,
                  geom="line", linetype=criterion) +
  xlab("Polynomial degree") +
  ylab("Criterion value") +

```

```
labs(title="Information theory & Bayes")+
geom_segment(aes(x=3, y=400,
                  xend=3, yend=325),
arrow = arrow(length = unit(0.3, "cm"),
angle=20, type="closed")) +
theme(legend.position=c(0.8,0.5))
```

AIC.plot

شکل 7.3

```
# function that will perform the "leave one out"
# crossvalidation for a y~poly(x, degree) polynomial
crossvalidate <- function(x, y, degree)
{
  preds <- numeric(length(x))
  for(i in 1:length(x))
  {
    x.in <- x[-i]
    x.out <- x[i]
    y.in <- y[-i]
    y.out <- x[i]
    m <- lm(y.in ~ poly(x.in, degree=degree) )
    new <- data.frame(x.in = seq(-3, 3, by=0.1))
    preds[i]<- predict(m, newdata=data.frame(x.in=x.out))
  }
  # the squared error:
  return(sum((y-preds)^2))
}

# crossvalidating all of the polynomial models
# and storing their squared errors in
# the "a" object
# plotting crossvalidated squared errors against
# model complexity
cross.plot <- qplot(1:max.poly,cross, data=a, geom=c("line"))+
```

```

      xlab("Polynomial degree") +
      ylab("Squared error") +
      geom_segment(aes(x=3, y=400,
                        xend=3, yend=200),
      arrow = arrow(length = unit(0.3, "cm"),
      angle=20, type="closed")) +
      labs(title="Crossvalidation")

cross.plot

```

کد مربوط به مثال داده‌های اجاره

```

library(gamlss)
PPP <- par(mfrow=c(2,2))
plot(R~F1, data=rent, col=gray(0.7), pch=15, cex=0.5)
plot(R~A, data=rent, col=gray(0.7), pch=15, cex=0.5)
plot(R~H, data=rent, col=gray(0.7), pch=15, cex=0.5)
plot(R~loc, data=rent, col=gray(0.7), pch=15, cex=0.5)
par(PPP)
r1 <- gamlss(R ~ F1+A+H+loc, family=N0, data=rent, trace=FALSE)
plot(r1)
r2 <- gamlss(R ~ F1+A+H+loc, family=GA, data=rent)
plot(r2)
r3 <- gamlss(R ~ pb(F1)+pb(A)+H+loc, family=GA, data=rent, trace=FALSE)
drop1(r3)
wp(r3, ylim.all=.6)
r4 <- gamlss(R ~ pb(F1)+pb(A)+H+loc, sigma.fo=~pb(F1)+pb(A)+H+loc,
family=GA, data=rent, trace=FALSE)
wp(r4, ylim.all=.6)
r6 <- gamlss(R ~ pb(F1)+pb(A)+H+loc, sigma.fo=~pb(F1)+pb(A)+H+loc,
nu.fo=~1, family=BCCGo, data=rent, trace=FALSE)
r7 <- gamlss(R ~ pb(F1)+pb(A)+H+loc, sigma.fo=~pb(F1)+pb(A)+
H+loc, nu.fo=~pb(F1)+pb(A)+H+loc, family=BCCGo, data=rent, trace=FALSE)
wp(r6, ylim.all=.6) ; title("r6: BCCG(mu, sigma)")
wp(r7, ylim.all=.6) ; title("r7: BCCG(mu, sigma, nu)")

```

کد مربوط به مثال داده‌های گونه‌های ماهی

```
library(gamlss)
data(species)
# creating the log(lake)
species <- transform(species, x=log(lake))
plot(fish~x,data=species)
mSI<-gamlss(fish~poly(x,2),data=species, sigma.fo=~1, nu.fo=~x,
family=SICHEL, n.cyc=60, trace=FALSE)
plot(fish~log(lake), data=species)
lines(species$x[order(species$lake)], fitted(mSI)[order(
species$lake)], col="red")
m9 <- gamlss(fish~poly(x,2), nu.fo=~x, data=species, family=SICHEL,
trace=FALSE)
wp(m9)
```

کد مربوط به داده‌های مدت بستری

```
data(aep)
prop<-with(aep, noinap/los)
par(mfrow = c(2, 2))
plot(prop~age, data=aep, cex=los/30)
plot(prop~sex,data=aep)
plot(prop~ward,data=aep)
plot(prop~year,data=aep)
```

کد مربوط به داده‌های فیلم

```
library(gamlss)
data(film90)
#attach(film90)
names(film90)
## [1] "lnosc" "lboopen" "lborev1" "dist"
par(mfcol = c(1,2))
with(film90, plot(lnosc,lborev1,pch=c(21,24)[unclass(dist)],
bg=c("red","lightgray")[unclass(dist)],
xlab="log no of screens", ylab="log extra revenue", main="(a)"))
legend("bottomright",legend=c("Independent","Major"),pch=c(21,24),
```

---

```

pt.bg=c("red","lightgray"),cex=0.7)
with(film90, plot(lboopen,lborev1,pch=c(21,24)[unclass(dist)],
bg=c("red","lightgray")[unclass(dist)],
xlab="log opening revenue", ylab="log extra revenue", main="(b)")
legend("bottomright",legend=c("Independent","Major"),pch=c(21,24),
pt.bg=c("red","lightgray"),cex=0.7)
install.packages("rgl")
library(rgl)
with(film90, plot3d(lboopen, lnosc, lborev1,
col=c("red","green3")[unclass(dist)]))
install.packages("gamlss.add")
library(gamlss.add)
m6 <- gamlss(lborev1~ga(~te(lboopen,lnosc))+dist, data=film90,trace=FALSE)
wp(m6, ylim.all=1.1)
library(mgcv)
plot(getSmo(m6))
vis.gam(getSmo(m6),theta = 0, phi = 30)
m7<- gamlss(lborev1~ga(~te(lboopen,lnosc))+dist,
sigma.fo=~ga(~te(lboopen,lnosc))+dist,
data=film90, trace=FALSE)
wp(m7, ylim.all=1.1)
m8 <- gamlss(lborev1 ~ pb(lboopen)+pb(lnosc) + dist,
sigma.fo = ~ pb(lboopen)+pb(lnosc) + dist,
nu.fo = ~ pb(lboopen)+pb(lnosc) + dist,
tau.fo = ~ pb(lboopen)+pb(lnosc) + dist,
family = BCPE, data = film90, trace=FALSE)
m9 <- gamlss(lborev1 ~ ga(~te(lboopen,lnosc)) + dist,
sigma.fo = ~ ga(~te(lboopen,lnosc)) + dist,
nu.fo = ~ ga(~te(lboopen,lnosc)) + dist,
tau.fo = ~ ga(~te(lboopen,lnosc)) + dist,
family = BCPE, data = film90, n.cyc=20, trace=FALSE)
par(mfrow=c(1,2))
wp(m8, ylim.all=0.5)
wp(m9, ylim.all=0.5)

```

```
wp(m9, xvar=~lboopen+lnosc, ylim.worm=1)
```

کد مربوط به داده‌های جرم و جنایت

```
library(gamlss.spatial)
data(columb)
data(columb.polys)
m1 <- gamlss(crime~ gmrf(district, polys=columb.polys), data=columb)
draw.polys(columb.polys, getSmo(m1), scheme="topo")
```



# مراجع

- [۱] بازرگان لاری ع.ر، (۱۳۸۵)، ”رگرسیون خطی کاربردی“ چاپ دوم، انتشارات دانشگاه شیراز،
- [۲] نیرومند ح.ع، (۱۳۸۴)، ”الگوهای خطی تعمیم یافته: با کاربردهای آن در علوم و مهندسی“ انتشارات دانشگاه فردوسی مشهد،
- [3] Agresti, A. (2015), ” **Foundations of linear and generalized linear models**”, John Wiley, Sons.
- [4] Burnham, K. P and Anderson, D. R. (2002), ”Model selection and multimodel inference: A practical information-theoretic approach”, **Springer**.
- [5] De Boor, C. (1977), ”Package for calculating with B-splines”, **SIAM Journal on Numerical Analysis**, 14(3), 441-472.
- [6] De Boor, C. (1978), ”**A Practical Guide to Splines**”, Springer, Berlin.
- [7] Dierckx, P. (1993), ”**Curve and Surface Fitting with Splines**”, Oxford University Press.
- [8] Eilers, P.H and Marx, B.D. (1996), ”Flexible smoothing with B-splines and penalties”, **Statistical science**, 11, 89-102.
- [9] Eilers, P.H., Marx, B.D. and Durbán, M. (2015), ”Twenty years of P-splines”, **statistics and operations research transactions**, 39(2), 149-186
- [10] Ghosh, A. and Kindermann, D.D.S. (2010), ”**Efficient thin plate spline interpolation and its application to adaptive optics**”.
- [11] Green, P.J. and Silverman, B.W. (1994), ”**Nonparametric regression and generalized linear models**” , CRC Press.
- [12] Hastie, T.J and Tibshirani, R.J. (1990), ”**Generalized additive models**”. Chapman Hall, London.



- 
- [13] James, G., Witten, D., Hastie, T. and Tibshirani, R. (2013), "**An introduction to statistical learning**", springer.
- [14] Kagerer, K. (2013), "**A short introduction to splines in least squares regression analysis**", University of Regensburg, Faculty of Business, Economics and Management Information Systems.
- [15] Langville, A.N. and Stewart, W.J. ( 2004), "The Kronecker product and stochastic automata networks", **computational and applied mathematics**, 167(2), 429-447.
- [16] Lambaert, H. (2006), "**Manual registration with thin plates**".
- [17] Mayr, A., Fenske, N., Hofner, B., Kneib, T. and Schmid, M. (2012), "Generalized additive models for location, scale and shape for high dimensional data—a flexible approach based on boosting", **Journal of the Royal Statistical Society: Series C (Applied Statistics)**, 61(3), 403-427.
- [18] McCullagh, P. and Nelder, J.A. (1989), "**Generalized linear models**", Chapman Hall, London.
- [19] McKinley, S. and Levine, M. (1998), "Cubic spline interpolation", **College of the Redwoods**, 45(1), 1049-1060.
- [20] Marra, G. and Radice, R. (2010), "Penalised regression splines: theory and application to medical research", **Statistical Methods in Medical Research**, 19(2), 107-125.
- [21] Nelder, J.A. (1972), "Generalized linear models", **Statistical Society Series A** , 135, 370-384.
- [22] O'Sullivan, F. (1986), "A statistical perspective on ill-posed inverse problems", **Statistical science**, 1, 502-518.
- [23] O'Sullivan, F. (1988), "Fast computation of fully automated log-density and log-hazard estimators", **SIAM Journal on scientific and statistical computing**, 9(2), 363-379.
- [24] Rigby, R.A. and Stasinopoulos, D.M. (2005), "Generalized additive models for location, scale and shape", **Journal of the Royal Statistical Society: Series C (Applied Statistics)**, 54(3), 507-554.
- [25] Ruppert, D., Wand, M.P. and Carroll, R.J. (2009), "Semiparametric regression during 2003–2007", **Electronic journal of statistics**, 3, 1193-1256.

- [26] Ruppert, D. (2002), "Selecting the number of knots for penalized splines", **Journal of computational and graphical statistics**, 11(4), 735-757.
- [27] Ruppert, D. and Carroll, R.J. (2000), "Theory Methods: Spatially adaptive Penalties for Spline Fitting", **Australian New Zealand Journal of Statistics**, 42(2), 205-223.
- [28] Stasinopoulos, M.D., Rigby, R.A., Heller, G.Z., Voudouris, V. and De Bastiani, F. (2017), **"Flexible regression and smoothing: using gamlss in r"**, CRC Press.
- [29] Stein, G.Z. and Juritz, J.M. (1988), "Linear models with an inverse gaussian poisson error distribution", **Communications in Statistics-Theory and Methods**, 17(2), 557-571.
- [30] Strang, G. and Aarikka, K. (1986), **"Introduction to applied mathematics"**, Wellesley-Cambridge Press Wellesley, MA.
- [31] Stoer, J. and Bulirsch, R. (2013), **"Introduction to numerical analysis"**, Springer Science Business Media.
- [32] Wahba, G. (1990), **Spline models for observational data**, SIAM, Philadelphia.
- [33] Wand, M.P. and Ormerod, J.T. (2008), "On semiparametric regression with O'Sullivan penalized splines", **Australian New Zealand Journal of Statistics**, 50(2), 179-198.
- [34] Wood, S.N. (2006), **Generalized additive models: an introduction with R**, CRC Press.
- [35] Wood, S.N. (2003) "Thin plate regression splines", **Journal of the Royal Statistical Society: Series B (Statistical Methodology)**, 65(1), 95-114.
- [36] Wassermann, L. (2006), "All of nonparametric statistics", **New York**.
- [37] Voudouris, V., Gilchrist, R., Rigby, R., Sedgwick, J. and Stasinopoulos, D. (2012), "Modelling skewness and kurtosis with the BCPE density in GAMLSS", **Journal of Applied Statistics**, 39(6), 1279-1293.

## **Aabstract**

Regression analysis is, in fact, the most widely used method amongst statistical techniques and a statistical tool for identifying the relationship between one or more covariates with a response variable. When a linear regression model is not efficient for analyzing a data set, we can use non-parametric regression methods. Splines as one of the interpolation tools are a powerful non-parametric method for modeling nonparametric regression. Splines construct curves that include polynomials of the same degree on substructures of a specified interval and join together with defined contiguous conditions, and pass through common knots between the two infrastructures. During the years since the introduction of Splines, their theoretical foundations have been developed, and different versions have been introduced. Penalized Splines (P-Splines) are among the most widely used tools for nonparametric modeling and smoothing problems. The appearance of P-Splines with many changes in smoothing problems has become a running and active field of research. Due to the broad applications and importance of Splines, we are going to introduce them and their growth as well as the development process and some of their features.

**Keywords:** Penalized Splines, Smoothing, Generalized Linear Models, Generalized Additive Models, Akaike Information Criterion.



**Shahrood University of Technology**

**Faculty Of Mathematical Sciences**

**MSc Thesis in: Statistics**

**the spline penalized and new application**

**By: Akram Ghaemi zadeh**

**Supervisor**

**Dr. Negar Eghbal**

**Advisor**

**Dr. Hossein Baghishani**

**July 2018**