

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



دانشکده علوم ریاضی

رشته آمار، گرایش آمار ریاضی

پایان نامه کارشناسی ارشد

# استنباط مبتنی بر درست‌نمایی در رگرسیون ریج برای داده‌های طولی

نگارنده: محمد رحمانی

استاد راهنما

دکتر محمد آرشی

شهریور ۱۳۹۷

شماره:

تاریخ:

باسمه تعالی



مدیریت تحصیلات تکمیلی

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

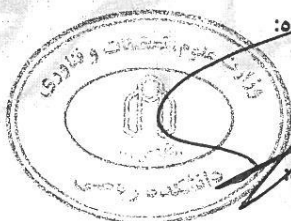
با نام و یاد خداوند متعال، ارزیابی جلسه دفاع از پایان نامه کارشناسی ارشد خانم / آقای محمد رحمانی با شماره دانشجویی ۹۴۰۷۷۷۴ رشته آمار گرایش آمار ریاضی تحت عنوان استنباط مبتنی بر درستمایی در رگرسیون ریبج برای داده‌های طولی که در تاریخ ۹۷/۶/۱۲ با حضور هیأت محترم داوران در دانشگاه صنعتی شاهرود برگزار گردید به شرح ذیل اعلام می‌گردد:

<input type="checkbox"/> مردود <input checked="" type="checkbox"/> قبول (با درجه: عالی.....)			
نوع تحقیق: <input checked="" type="checkbox"/> نظری <input type="checkbox"/> عملی			
عضو هیأت داوران	نام و نام خانوادگی	مرتبه علمی	امضاء
۱- استاد راهنمای اول	دکتر محمد آرشی	دانشیار	
۲- استاد راهنمای دوم	-----	-----	-----
۳- استاد مشاور	-----	-----	-----
۴- نماینده تحصیلات تکمیلی	دکتر نگار اقبال	استادیار	
۵- استاد ممتحن اول	دکتر حسین باغیشنی	استادیار	
۶- استاد ممتحن دوم	دکتر محمدرضا ربیعی	استادیار	

نام و نام خانوادگی رئیس دانشکده:

دکتر ابراهیم هاشمی

تاریخ و امضاء و مهر دانشکده:



تبصره: در صورتی که کسی مردود شود حداکثر یکبار دیگر (در مدت مجاز تحصیل) می‌تواند از پایان نامه خود دفاع نماید (دفاع مجدد نباید زودتر از ۴ ماه برگزار شود).

تقدیم بہ

پدر بزرگوار و مادر مہربانم

آن دو فرشتہ ای کہ از برایم عاشقانہ و خالصانہ ز خویش تن گذشتند...

پس ہستی من ز ہستی اوست

تا ہستم و ہست دارمش دوست

ایرج میرزا

تقدیم بہ

استاد کرامت قدم و کترا آرش

کہ از چشمہ جوشان علم و اخلاق شایستہ ایشان بہرہ بردم.  
زیر سایہ لطف الہی، پایدار باشند و برقرار...

# سپاس گزارمی...

منت خدای را عزوجل که طاعتش موجب قرمت و به شکر اندرش مزید نعمت، هر نفسی که فرومی رود مدتی است و چون برمی آید مفرح ذات، پس در هر نفسی دو نعمت موجود است و بر هر نعمت شکر می واجب...

از دست و زبان که برآید  
کز عهده شکرش به درآید

سعدی

پس از حمد خداوند باری تعالی بر خود لازم می دانم تا زحمات اساتید گرانقدرم را ارج نهم. بسیار خرسندم که افتخار شاگردی استاد فرزانه جناب آقای دکتر محمد آرشی را دارم، بیش از همه سپاس گزار زحمات، محبت ها و راهنمایی های ایشان هستم و از خداوند رحمان بهترین ها را برای ایشان آرزو دارم.

از اساتید بزرگوار و دلسوزم آقایان دکتر حسین باغیثنی و دکتر محمدرضا ربیعی که خالصانه از علم خویش به من آموختند و زحمت داوری این پایان نامه را نیز متقبل شدند، کمال تقدیر و تشکر را دارم.

بر خود لازم می دانم از زحمات اساتید ارجمندم آقایان دکتر شاهسونی، دکتر نزاکتی و سرکارخانم دکتر اقبال که همواره راهنما و راه گشای اینجانب بودند و تمامی عزیزانی که به هر نحوی در طول دوران تحصیل و تهیه این مجموعه مرا یاری دادند صمیمانه تشکر و قدردانی نمایم و از خداوند متعال برای همه ایشان موفقیت و سربلندی را خواستارم.

سعی نابوده در این راه به جایی نرسی  
مزد اگر میطلبی طاعت استاد بیر

حافظ

محمد رحمانی  
شهریور ۱۳۹۷

## تعهد نامه

اینجانب محمد رحمانی دانشجوی کارشناسی ارشد رشته آمار دانشکده علوم ریاضی دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان **استنباط مبتنی بر درست‌نمایی در رگرسیون ریج برای داده‌های طولی**، تحت راهنمایی محمد آرشی متعهد می‌شوم:

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

محمد رحمانی

شهریور ۱۳۹۷

### مالکیت نتایج و حق نشر

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

## چکیده

داده‌های طولی به داده‌هایی اطلاق می‌شود که در جمع‌آوری آن‌ها، اندازه‌گیری ویژگی‌های مختلف برای هر کدام از واحدهای آزمایشی در طول زمان تکرار می‌شود. یکی از ابزارهای تحلیل داده‌های طولی استفاده از مدل‌های آمیخته خطی است، اما وجود مشکل همخطی در این نوع داده‌ها استفاده از روش‌های مختلفی همچون روش‌های انقباضی را موجب می‌شود. در این پایان‌نامه به تحلیل این نوع داده‌ها با استفاده از مدل‌های آمیخته خطی و ترکیب آن‌ها با برخی از روش‌هایی همچون رگرسیون ریج پرداخته شده و برآوردهای ریج و لیو با خطای بیضی‌گون ارائه شده‌اند. همچنین رفتار برآوردهای حاصل از این روش‌ها به کمک مطالعات شبیه‌سازی و تحلیل داده‌های واقعی مورد بررسی قرار گرفته‌است.

کلمات کلیدی: الگوریتم EM، برآوردهای انقباضی، برآوردهای لیو، توزیع بیضی‌گون، رگرسیون ریج، مدل‌های آمیخته خطی.

## پیش‌گفتار

پیشینه روش‌های تحلیل داده‌های طولی به حدود نیم قرن پیش باز می‌گردد. پایه این روش‌ها توسط یک ستاره‌شناس انگلیسی بنا شد که زمینه‌ای برای فعالیت‌های پژوهشگران در این موضوع، همچون لیرد و ویر (۱۹۸۲) با معرفی مدل اثرات تصادفی برای داده‌های طولی فراهم آمد. با گسترش این روش‌ها تاکنون مدل‌های مختلفی مبنی بر تحلیل داده‌های طولی از جمله مدل‌های حاشیه‌ای، مدل‌های انتقال و مدل‌های آمیخته توسط پژوهشگران معرفی شده‌اند. در شرایطی که در داده‌های متغیرهای توضیحی همخطی وجود دارد، روش‌های معمول برآورد کارا نبوده و باید از روش‌های جایگزینی همچون ریج (هورل و کنارد، ۱۹۷۰) استفاده کرد. تمرکز این پایان‌نامه بر تحلیل داده‌های طولی با استفاده از مدل‌های آمیخته خطی است. الیوت و همکاران (۲۰۱۱) روش رگرسیون ریج را در مواجهه با داده‌های همخط در مدل آمیخته خطی، به منظور افزایش کارایی در تحلیل داده‌های طولی به کار بردند. در این مجموعه رهیافت آن‌ها را بررسی کرده و روشی جایگزین را ارائه می‌کنیم. همچنین تعمیم‌های دیگری از جمله تغییر توزیع جامعه هدف به بیضی‌گون و بررسی مدل‌های خطی تعمیم‌یافته را در نظر می‌گیریم. این مجموعه مشتمل بر ۳ فصل و یک پیوست است که محتویات آن‌ها به‌طور خلاصه در زیر آمده‌است.

- در فصل اول پس از معرفی برخی از مدل‌های ذکرشده مبنی بر تحلیل داده‌های طولی، مدل‌های آمیخته خطی را با نگاه دقیق‌تری مورد مطالعه قرار داده‌ایم.
- در فصل دوم با هدف دستیابی به نتایج مطمئن‌تر و دقیق‌تر در تحلیل داده‌های طولی، پس از معرفی و بررسی مدل آمیخته خطی ریج به منظور مقابله با وجود مشکل همخطی در داده‌ها، مدل آمیخته خطی ریج با خطای بیضی‌گون و برآوردگرهای انقباضی در مدل آمیخته خطی ریج به جهت بهبود برآوردگرها ارائه شده‌اند.
- در فصل سوم به سبب ویژگی خطی بودن برآوردگر لیمو نسبت به پارامتر آن، مدل آمیخته خطی لیمو، مدل آمیخته خطی لیمو با خطای بیضی‌گون و برآوردگر انقباضی در مدل آمیخته خطی لیمو ارائه شده و با بهره‌گیری از روش‌های شبیه‌سازی مورد بررسی قرار گرفته‌اند.
- همچنین گزیده‌ای از برنامه‌های کامپیوتری که تحت محیط R نوشته شده‌اند، در پیوست آمده‌است.

لازم به ذکر است محتویات بخش‌های ۳.۲، ۴.۲ و فصل ۳ کاملاً جدید بوده و از یافته‌های نویسنده این مجموعه است که حاصل آن تهیه ۴ مقاله است که در فهرست مقالات مستخرج از پایان‌نامه آمده‌اند.



# لیست مقالات مستخرج از پایان نامه

1. Rahmani, M. and Arashi, M. (2017), Longitudinal Data Analysis Using Liu Regression, *11th Seminar on Probability and Stochastic Processes*, Imam Khomeini International University, Qazvin, Iran.
2. Rahmani, M. and Arashi, M. (2017), Ridge Regression for Longitudinal Data with Elliptically Contoured Errors, *48th Annual Iranian Mathematics Conference*, Bu-Ali Sina University, Hamedan, Iran.
3. Rahmani, M., Arashi, M. and Jenatabadi, H.S. (2018), Longitudinal data analysis using Liu regression, *Biostatistics and Biometric Open Access Journal*, 7(5).
4. Rahmani, M., Arashi, M., Mamode Khan, N. and Sunecher, Y. (2018), Improved mixed model for longitudinal data analysis using shrinkage method, Submitted.

# فهرست مطالب

ق	فهرست تصاویر
ش	فهرست جداول
۱	۱ مدل‌های خطی برای داده‌های طولی
۱	۱.۱ مقدمه
۱	۲.۱ داده‌های طولی
۲	۱.۲.۱ نمادگذاری داده‌های طولی
۳	۲.۲.۱ اهداف و ضرورت تحلیل داده‌های طولی
۵	۳.۱ انواع مدل‌بندی داده‌های طولی
۵	۱.۳.۱ مدل‌های حاشیه‌ای
۱۰	۲.۳.۱ مدل‌های انتقال
۱۰	۳.۳.۱ مدل‌های آمیخته خطی
۱۲	۴.۱ مدل آمیخته خطی
۱۳	۱.۴.۱ برآورد کلاسیک پارامترهای مدل
۱۵	۲.۴.۱ مثال کاربردی
۱۷	۵.۱ مدل آمیخته خطی تعمیم‌یافته
۱۸	۱.۵.۱ مؤلفه‌های مدل خطی تعمیم‌یافته
۲۰	۲.۵.۱ برآورد پارامترهای مدل
۲۵	۲ رگرسیون ریبج برای داده‌های طولی
۲۵	۱.۲ مقدمه
۲۵	۲.۲ مدل آمیخته خطی ریبج
۳۱	۱.۲.۲ مطالعه شبیه‌سازی
۳۳	۲.۲.۲ مثال واقعی
۳۵	۳.۲ مدل آمیخته خطی ریبج با خطای بیضی‌گون

۳۶	.....	توزیع بیضی گون	۱.۳.۲
۴۰	.....	مدل MR بیضی گون	۲.۳.۲
۴۲	.....	مطالعه شبیه سازی	۳.۳.۲
۴۴	.....	برآوردگرهای انقباضی در مدل آمیخته خطی ریج	۴.۲
۴۷	.....	مطالعه شبیه سازی	۱.۴.۲
۴۹	.....	مثال واقعی	۲.۴.۲
۵۲	.....	مدل آمیخته خطی تعمیم یافته ریج	۵.۲
۵۷		<b>رگرسیون لیو برای داده های طولی</b>	<b>۳</b>
۵۷	.....	مقدمه	۱.۳
۵۸	.....	مدل آمیخته خطی لیو	۲.۳
۶۱	.....	مطالعه شبیه سازی	۱.۲.۳
۶۳	.....	مثال واقعی	۲.۲.۳
۶۵	.....	مدل آمیخته خطی لیو با خطای بیضی گون	۳.۳
۶۶	.....	مطالعه شبیه سازی	۱.۳.۳
۶۸	.....	برآوردگرهای انقباضی در مدل آمیخته خطی لیو	۴.۳
۶۹	.....	مطالعه شبیه سازی	۱.۴.۳
۷۰	.....	نتیجه گیری و آینده تحقیق	۵.۳
۷۰	.....	خلاصه و نتیجه گیری	۱.۵.۳
۷۱	.....	پیشنهادات برای پژوهش های آینده	۲.۵.۳
۷۳		<b>مراجع</b>	
۷۷		<b>آ گزیده های از برنامه های کامپیوتری</b>	
۹۱		واژه نامه فارسی به انگلیسی	
۹۳		واژه نامه انگلیسی به فارسی	

# فهرست تصاویر

۲	.....	نمودار ساختگی به منظور توجیه ساختار داده‌های طولی	۱.۱
۴	.....	نمودار قدرت خواندن در کودکان	۲.۱
۳۷	.....	نمودارهای منحنی‌های تراز و تابع چگالی توزیع نرمال دومتغیره	۱.۲
۳۸	.....	نمودارهای منحنی‌های تراز و تابع چگالی توزیع $t$ دومتغیره	۲.۲
۳۸	.....	نمودارهای منحنی‌های تراز و تابع چگالی توزیع کاتز دومتغیره	۳.۲
۳۹	.....	نمودارهای منحنی‌های تراز و تابع چگالی توزیع لاپلاس دومتغیره	۴.۲
۵۲	.....	نمودار جعبه‌ای مربوط به مقدار MPE برآوردها	۵.۲

# فهرست جداول

۳	.....	۱.۱	ساختار داده‌های طولی
۱۷	.....	۲.۱	برآورد ضرایب با استفاده از مدل خطی و آمیخته خطی
۳۲	.....	۱.۲	نتایج مطالعه شبیه‌سازی مدل آمیخته خطی ریج
۳۴	.....	۲.۲	همبستگی بین متغیرهای پیش‌بین در داده‌های GENE
۳۵	.....	۳.۲	برآورد ضرایب در داده‌های GENE
۴۳	.....	۴.۲	نتایج مطالعه شبیه‌سازی مدل آمیخته خطی ریج با خطای بیضی‌گون
۴۷	.....	۵.۲	نتایج مطالعه شبیه‌سازی برآوردگر نوع استاین در مدل آمیخته خطی ریج
۴۸	.....	۶.۲	نتایج مطالعه شبیه‌سازی برآوردگر آزمون اولیه در مدل آمیخته خطی ریج
۵۰	.....	۷.۲	توصیف متغیرهای مجموعه داده بخش ۲.۴.۲
۵۰	.....	۸.۲	اطلاعات خلاصه شده از متغیرهای مجموعه داده بخش ۲.۴.۲
۵۱	.....	۹.۲	برآورد ضرایب در داده‌های بخش ۲.۴.۲
۵۲	.....	۱۰.۲	انحراف استاندارد برآوردهای جدول ۹.۲
۶۲	.....	۱.۳	نتایج مطالعه شبیه‌سازی مدل آمیخته خطی لیو
۶۴	.....	۲.۳	توصیف متغیرهای مجموعه داده بخش ۱.۲.۳
۶۵	.....	۳.۳	برآورد ضرایب در مجموعه داده بخش ۱.۲.۳
۶۷	.....	۴.۳	نتایج مطالعه شبیه‌سازی مدل آمیخته خطی لیو با خطای بیضی‌گون
۷۰	.....	۵.۳	نتایج مطالعه شبیه‌سازی برآوردگر نوع استاین در مدل آمیخته خطی لیو

# فصل ۱

## مدل‌های خطی برای داده‌های طولی

### ۱.۱ مقدمه

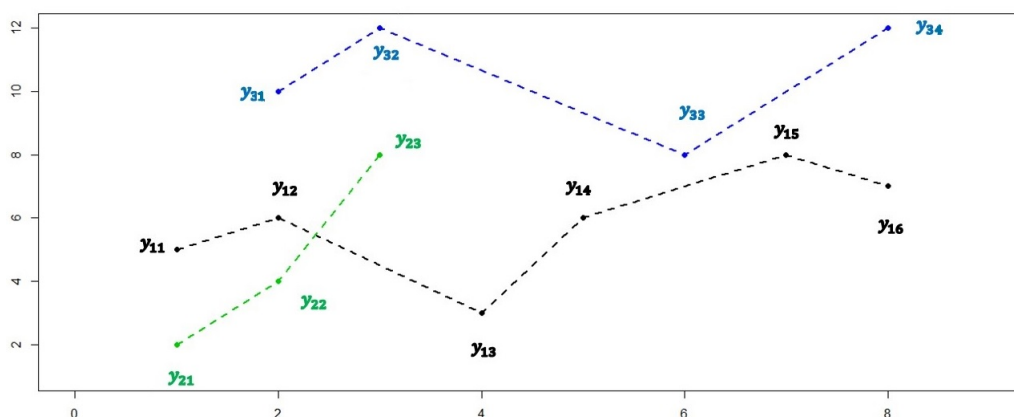
از آنجایی که آشنایی با یک موضوع جدید، مستلزم شناخت مقدمات مربوط به آن موضوع می‌باشد، در این فصل به معرفی داده‌های طولی<sup>۱</sup> و همچنین شرح چگونگی مدل‌بندی این داده‌ها می‌پردازیم. مطالب این فصل در قسمت تعریف داده‌های طولی و انواع مدل‌های آن عمدتاً برگرفته از تعاونی (۱۳۹۲) و در قسمت مدل‌های خطی (تعمیم‌یافته) از مک کولاک و سیرل (۲۰۰۱) است.

### ۲.۱ داده‌های طولی

نوع مطالعه و روش گردآوری داده‌ها یکی از مسائل مهم در یک تحلیل آماری می‌باشد. تحلیل‌های رگرسیونی و سری‌های زمانی دو روش متداول در تحلیل داده‌ها به شمار می‌آیند. تحلیل داده‌های طولی، ترکیبی از تحلیل‌های رگرسیونی و سری‌های زمانی که در آن ارتباط بین چندین ویژگی در طول زمان مد نظر است را فراهم می‌آورد. بنابراین داده‌های طولی به مجموعه وسیعی از داده‌های گسسته و پیوسته اطلاق می‌شود که در آن ویژگی‌های چندین واحد آزمایشی (مانند انسان، حیوان یا نمونه‌های آزمایشگاهی) به‌عنوان نمونه‌ای از اعضای

<sup>۱</sup>Longitudinal data

جامعه، در طول زمان اندازه‌گیری می‌شوند. لازم به ذکر است همان‌طور که در شکل ۱.۱ نشان داده شده‌است، در داده‌های طولی تعداد مشاهدات همه واحدها و همچنین فاصله نقاط زمانی آن‌ها لزوماً برابر نیستند. جهت آشنایی بیشتر با ساختار داده‌های طولی در ادامه به معرفی نمادگذاری متداول این داده‌ها می‌پردازیم.



شکل ۱.۱: نمودار ساختگی به منظور توجیه ساختار داده‌های طولی

### ۱.۲.۱ نمادگذاری داده‌های طولی

جدول ۱.۱ را در نظر بگیرید که در آن مقادیر متغیر پاسخ و مقدار  $x_{ijk}$  مقدار  $k$  امین متغیر پیش‌بین مربوط به واحد  $i$  ام در زمان  $j$  ام می‌باشند، به طوری که  $i = 1, \dots, n$ ،  $j = 1, \dots, t_i$  و  $k = 1, \dots, p$ ، که در آن  $n$  بیان‌گر تعداد واحدهای تحت مطالعه،  $t_i$  تعداد دوره‌های زمانی برای واحد  $i$  ام و  $p$  تعداد متغیرهای پیش‌بین است. لازم به ذکر است که هر دو متغیر پاسخ و پیش‌بین ممکن است گسسته یا پیوسته باشند. همچنین متغیرهای پیش‌بین می‌توانند در طول زمان ثابت (مانند نژاد) یا متغیر (مانند سن) باشند.

جدول ۱.۱: ساختار داده‌های طولی

$x_{ijk}$	$y_{ij}$	زمان	واحد
$x_{111}, \dots, x_{11p}$	$y_{11}$	۱	۱
⋮	⋮	⋮	⋮
$x_{1t_11}, \dots, x_{1t_1p}$	$y_{1t_1}$	$t_1$	۱
$x_{211}, \dots, x_{21p}$	$y_{21}$	۱	۲
⋮	⋮	⋮	⋮
$x_{1t_21}, \dots, x_{1t_2p}$	$y_{2t_2}$	$t_2$	۲
⋮	⋮	⋮	⋮
$x_{n11}, \dots, x_{n1p}$	$y_{n1}$	۱	n
⋮	⋮	⋮	⋮
$x_{nt_n1}, \dots, x_{nt_np}$	$y_{nt_n}$	$t_n$	n

طبق آن چه که بیان شد، این داده‌ها را می‌توان از دو بعد زمانی و مقطعی<sup>۲</sup> مورد توجه قرار داد. این ماهیت دو بعدی، اهداف و لزوم ارائه ترکیبی از ابزارهای آماری را در تحلیل داده‌های طولی آشکار می‌سازد.

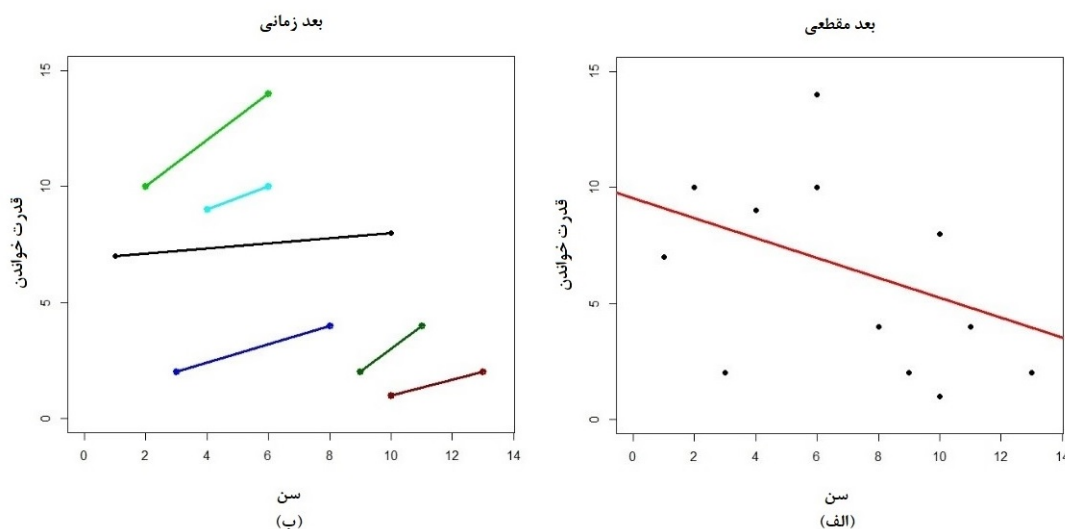
### ۲.۲.۱ اهداف و ضرورت تحلیل داده‌های طولی

با جمع‌آوری مشاهدات مربوط به واحدها در طول زمان، علاوه بر دستیابی به رابطه بین متغیرها می‌توان تغییرات واحدها در طول زمان را نیز مورد مطالعه قرار داد. به‌عنوان مثال اگر داده‌هایی که در زمان زام جمع‌آوری شده‌اند را در نظر بگیریم، در واقع واحدهای اول تا  $n$ ام را تنها در یک مقطع از زمان مورد بررسی قرار داده‌ایم. بنابراین برای هر دوره زمانی، داده‌ها مقطعی هستند. اما اگر داده‌های مربوط به واحد  $i$ ام را در نظر بگیریم، تنها یک واحد را در زمان‌های متوالی مورد بررسی قرار داده‌ایم که در این صورت، داده‌ها اثر زمان را نشان می‌دهند. لذا در تحلیل داده‌های طولی سه هدف اصلی می‌تواند مد نظر هستند: (۱) مطالعه تغییرات واحدها در طول زمان، (۲) مطالعه پیرامون رابطه بین متغیرهای پاسخ و پیش‌بین، (۳) مطالعه اثر زمان بر رابطه بین متغیرها.

برای درک بهتر این موضوع به مثالی از دیگل و همکاران (۲۰۰۲) اشاره می‌کنیم.

<sup>2</sup>Sectional





شکل ۲.۱: نمودار قدرت خواندن در کودکان

در شکل ۲.۱، داده‌ها حاصل یک مطالعه طولی مربوط به قدرت خواندن کودکان هستند که در آن شاخص توانایی مطالعه هر کودک، دو بار در طول زمان اندازه‌گیری شده‌است. بُعد مقطعی داده‌ها در قسمت (الف)، کاهش قدرت خواندن کودکان با گذشت زمان را نشان می‌دهد، درحالی‌که بُعد زمانی داده‌ها در قسمت (ب)، افزایش توانایی هر کودک با گذشت زمان را نشان می‌دهد. به‌طور کلی در مورد این داده‌ها می‌توان گفت توانایی خواندن در کودکان کم‌سن‌تر سطح بالاتری دارد و همچنین تمام کودکان مورد مطالعه، توانایی خود را با گذشت زمان بهبود می‌بخشند. بنابراین تحلیل داده‌های طولی که با هدف مشاهده تغییر رفتار پدیده‌ها در طول زمان انجام می‌گیرد، نسبت به داده‌های صرفاً مقطعی یا صرفاً زمانی از اهمیت ویژه‌ای برخوردار است.

علاوه بر حجیم بودن اطلاعات حاصل از داده‌های طولی و ماهیت دو بعدی آن‌ها که پیچیدگی داده‌ها را سبب می‌شوند، دلیل دیگری نیز وجود دارد که تحلیل این داده‌ها را با چالش رو برو می‌سازد. برخلاف روش‌های آماری متداول که بر پایه پذیره استقلال مشاهدات طراحی شده‌اند، در این داده‌ها مجموعه مشاهدات مربوط به هر واحد، دارای همبستگی درونی بوده و از مشاهدات سایر واحدها مستقل می‌باشند. واضح است که برای دستیابی به استنباط آماری و معتبر، در نظر گرفتن این همبستگی امری ضروری است.

دلایل فوق از جمله مهمترین مواردی است که روش‌های آماری کلاسیک برای تحلیل داده‌های طولی از قبیل رگرسیون و سری‌های زمانی به تنهایی قادر به تحلیل دقیق این داده‌ها نمی‌باشند. با هدف تحلیل معتبر این داده‌ها آماردانان روش‌های متعددی معرفی نموده‌اند که در ادامه به آن‌ها خواهیم پرداخت.

## ۳.۱ انواع مدل‌بندی داده‌های طولی

در این بخش به معرفی و بررسی مختصر برخی از مدل‌های رگرسیونی همچون مدل‌های حاشیه‌ای<sup>۳</sup>، مدل‌های انتقال<sup>۴</sup> و مدل‌های خطی آمیخته<sup>۵</sup> که در تحلیل داده‌های طولی به کار می‌روند، می‌پردازیم.

### ۱.۳.۱ مدل‌های حاشیه‌ای

مدل‌های حاشیه‌ای، مدل‌هایی هستند که در آن‌ها اثر متغیرهای پیش‌بین بر روی متغیر پاسخ به طور مجزا از همبستگی بین پاسخ‌ها، مدل‌بندی می‌شوند. به عبارت دیگر، در این روش مدل‌های جداگانه‌ای برای میانگین حاشیه‌ای، واریانس حاشیه‌ای و همبستگی حاشیه‌ای در نظر گرفته می‌شود. در صورت پیوسته بودن متغیر پاسخ، مدل حاشیه‌ای به صورت یک مدل خطی کلی با در نظر گرفتن همبستگی داده‌ها ارائه می‌شود. مدل خطی زیر یک مدل حاشیه‌ای به شمار می‌آید:

$$\mathbf{y}_i = \mathbf{X}_i^T \boldsymbol{\beta} + \epsilon_i, i = 1, \dots, n, \quad (1.1)$$

که در آن  $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{in_i})^T$  و  $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{it_i})^T$ . بردار متغیرهای پیش‌بین مربوط به پارامتر  $\beta_i$  از  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  و بردار خطای تصادفی بوده، به طوری که  $\epsilon_i$ ‌ها از هم مستقل و هم‌توزیع (iid) با توزیع زیر هستند

$$\epsilon_i \stackrel{iid}{\sim} \mathcal{N}_{t_i}(\mathbf{0}, \mathbf{V}_i).$$

همچنین  $\mathbf{V}_i > \circ, i = 1, \dots, n$ ، ماتریس معین مثبت است. وجود پذیره استقلال بین واحدهای مختلف و همبستگی درون واحدها (همبستگی حاشیه‌ای) که در مدل به صورت ماتریس  $\mathbf{V}_i$  نشان داده شده‌است، تمایز مدل‌های حاشیه‌ای را با مدل‌های خطی کلاسیک، آشکار می‌سازد.

### برآورد پارامترهای مدل

روش‌های متعددی برای برآورد پارامترهای مدل خطی (۱.۱) وجود دارند، که از جمله متداول‌ترین آن‌ها می‌توان به روش ماکزیمم درست‌نمایی<sup>۶</sup> (ML) و کمترین توان‌های دوم تعمیم‌یافته<sup>۷</sup> (GLS) اشاره کرد. در این بخش به برآورد پارامترهای مدل با استفاده از روش ML خواهیم پرداخت.

<sup>3</sup>Marginal

<sup>4</sup>Transition

<sup>5</sup>Linear mixed models

<sup>6</sup>Maximum likelihood

<sup>7</sup>Generalized least squares

با توجه به مدل (۱.۱)، می‌توان نوشت

$$E(\mathbf{y}|\mathbf{X}) = \mathbf{X}\boldsymbol{\beta} \quad (۲.۱)$$

که در آن  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$  و  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^T$ . تحت مفروضات مدل (۱.۱) داریم

$$\mathbf{y}_i \stackrel{iid}{\sim} N_{t_i}(\mathbf{X}_i^T \boldsymbol{\beta}, \mathbf{V}_i).$$

فرض کنید ماتریس کواریانس مدل به یک پارامتر  $\alpha$  وابسته باشد که آن را با  $\mathbf{V}_i(\alpha)$  نشان می‌دهیم. دلیل این تخصیص بررسی انواع ساختارهای همبستگی است که جلوتر به آن اشاره می‌کنیم.

تابع درستنمایی مدل (۱.۱) و لگاریتم آن به صورت زیر حاصل می‌شوند

$$L(\boldsymbol{\beta}, \alpha) = \prod_{i=1}^n (\pi)^{\frac{-t_i}{\nu}} |\mathbf{V}_i(\alpha)|^{-\frac{1}{\nu}} \exp \left\{ -\frac{1}{\nu} (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta})^T \mathbf{V}_i(\alpha)^{-1} (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta}) \right\} \quad (۳.۱)$$

$$l(\boldsymbol{\beta}, \alpha) = \sum_{i=1}^n \left\{ -\frac{t_i}{\nu} \log(\pi) - \frac{1}{\nu} \log |\mathbf{V}_i(\alpha)| \right\} - \frac{1}{\nu} \sum_{i=1}^n \left\{ (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta})^T \mathbf{V}_i(\alpha)^{-1} (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta}) \right\}.$$

فرض کنید  $U_i(\alpha) = \frac{d}{d\alpha} \mathbf{V}_i(\alpha)$ ، در این صورت با استفاده از دو قانون مشتق‌گیری

$$\frac{d \log |\mathbf{A}|}{dx} = \text{tr} \left( \mathbf{A}^{-1} \frac{d\mathbf{A}}{dx} \right), \quad \frac{d\mathbf{A}^{-1}}{dx} = -\mathbf{A}^{-1} \frac{d\mathbf{A}}{dx} \mathbf{A}^{-1}$$

معادلات نرمال به صورت زیر به دست می‌آیند

$$\frac{dl(\boldsymbol{\beta}, \alpha)}{d\boldsymbol{\beta}} = \sum_{i=1}^n \mathbf{X}_i^T \mathbf{V}_i(\alpha)^{-1} (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta}) = 0 \quad (۴.۱)$$

$$\begin{aligned} \frac{dl(\boldsymbol{\beta}, \alpha)}{d\alpha} &= \sum_{i=1}^n -\frac{1}{\nu} \text{tr} \left( \mathbf{V}_i^{-1}(\alpha) U_i(\alpha) \right) \\ &+ \frac{1}{\nu} (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta})^T \left( \mathbf{V}_i^{-1}(\alpha) U_i(\alpha) \mathbf{V}_i^{-1}(\alpha) \right) (\mathbf{y}_i - \mathbf{X}_i^T \boldsymbol{\beta}) = 0 \end{aligned} \quad (۵.۱)$$

با حل رابطه (۴.۱)، برآورد پارامتر  $\boldsymbol{\beta}$  به صورت زیر به دست می‌آید

$$\hat{\boldsymbol{\beta}} = \left( \sum_{i=1}^n \mathbf{X}_i^T \hat{\mathbf{V}}_i^{-1}(\alpha) \mathbf{X}_i \right)^{-1} \sum_{i=1}^n \mathbf{X}_i^T \hat{\mathbf{V}}_i^{-1}(\alpha) \mathbf{y}_i \quad (۶.۱)$$

که در آن  $\hat{\mathbf{V}}_i(\alpha) = \mathbf{V}_i(\hat{\alpha})$  که از رابطه (۵.۱) حاصل می‌شود. به سبب وجود  $\hat{\mathbf{V}}_i^{-1}(\alpha)$  در ساختار برآوردگر  $\hat{\boldsymbol{\beta}}$ ، یافتن برآورد آن از طریق روش‌های عددی همچون نیوتن-رافسون امکان‌پذیر است. در ادامه فرض می‌کنیم ماتریس کواریانس  $\mathbf{V}_i(\alpha)$ ، به صورت زیر است

$$\mathbf{V}_i(\alpha) = \phi \mathbf{A}_i^{\frac{1}{\nu}} \mathbf{R}_i(\alpha) \mathbf{A}_i^{\frac{1}{\nu}}. \quad (۷.۱)$$

که در آن،  $A_i$  یک ماتریس قطری در اندازه  $t_i \times t_i$  با عناصر روی قطر  $\phi = var(y_{ij})$  بوده و  $R_i(\alpha)$  به عنوان ماتریس همبستگی عملی شناخته می‌شود که برای آن ساختارهای متعددی در نظر گرفته شده است. نظر به این که برآورد پارامترهای رگرسیونی به انتخاب این ساختارها وابسته است، می‌توان در انواع ساختارها به موارد زیر اشاره نمود.

**ساختار استقلال<sup>۸</sup>:** که در آن پاسخ‌های مکرر هر واحد آزمایشی مستقل از هم فرض

می‌شود.

برای ساده نویسی از نماد  $\alpha$  در آرگومان  $R_i(\alpha)$  صرف نظر می‌کنیم.

$$R_i = \begin{pmatrix} 1 & \circ & \dots & \circ \\ \circ & 1 & \dots & \circ \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \dots & 1 \end{pmatrix}_{t_i \times t_i} \quad (۸.۱)$$

بدیهی است در این حالت  $R_i = R_i(\circ)$ .

**ساختار همبستگی متبادل پذیر<sup>۹</sup>:** در آن مشاهدات مربوط به هر واحد آزمایشی، دارای

همبستگی یکسان می‌باشند.

$$R_i = \begin{pmatrix} 1 & \alpha & \dots & \alpha \\ \alpha & 1 & \dots & \alpha \\ \vdots & \vdots & \ddots & \vdots \\ \alpha & \alpha & \dots & 1 \end{pmatrix}_{t_i \times t_i} \quad (۹.۱)$$

**همبستگی غیر ساختاری<sup>۱۰</sup>:** نوع دیگری از ساختارهای همبستگی که به صورت زیر است

$$R_i = \begin{pmatrix} 1 & \alpha_{12} & \dots & \alpha_{1t_i} \\ \alpha_{21} & 1 & \dots & \alpha_{2t_i} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{t_i 1} & \alpha_{t_i 2} & \dots & 1 \end{pmatrix}_{t_i \times t_i} \quad (۱۰.۱)$$

**همبستگی  $m$ <sup>۱۱</sup>:** در آن تنها مشاهداتی دارای همبستگی می‌باشند که اختلاف طول زمانی

مشاهده آن‌ها کمتر از  $m$  است. بنابراین ساختار همبستگی به صورت زیر است

$$Corr(y_{ij}, y_{ij'}) = \begin{cases} 1 & s = 0 \\ \alpha_s & s = 1, 2, \dots, m, \quad s = |i - j'| \\ 0 & s > m \end{cases}$$

<sup>۸</sup>Independence

<sup>۹</sup>Exchangeable

<sup>۱۰</sup>Unstructured

<sup>۱۱</sup>m-dependent

صورت ماتریسی آن برای حالت  $t_i = 4$  و  $m = 2$  دارای ساختار زیر است

$$\mathbf{R}_i = \begin{pmatrix} 1 & \alpha_1 & \alpha_2 & 0 \\ \alpha_1 & 1 & \alpha_1 & \alpha_2 \\ \alpha_2 & \alpha_1 & 1 & \alpha_1 \\ 0 & \alpha_2 & \alpha_1 & 1 \end{pmatrix}_{4 \times 4} \quad (11.1)$$

**همبستگی اتورگرسیو<sup>۱۲</sup>:** متداول‌ترین آن اتورگرسیو مرتبه اول  $AR(1)$  است و به صورت

زیر تعیین می‌شود

$$\text{Corr}(y_{ij}, y_{ij'}) = \alpha^s, \quad s = 0, 1, \dots, t_i - j \quad (12.1)$$

صورت ماتریسی آن برای حالت  $t_i = 4$  دارای ساختار زیر است

$$\mathbf{R}_i = \begin{pmatrix} 1 & \alpha & \alpha^2 & \alpha^3 \\ \alpha & 1 & \alpha & \alpha^2 \\ \alpha^2 & \alpha & 1 & \alpha \\ \alpha^3 & \alpha^2 & \alpha & 1 \end{pmatrix}_{4 \times 4} \quad (13.1)$$

در صورتی که طبق اطلاعات قبلی ساختار ماتریس همبستگی معلوم باشد، در این صورت ماتریس همبستگی به صورت یک ماتریس ثابت و معلوم در روابط محاسبه برآورد پارامتر  $\beta$ ، جایگزین می‌شود.

پس از آشنایی با ساختارهای مختلف ماتریس همبستگی، به روش‌های عددی در برآورد پارامترهای مدل می‌پردازیم. در رابطه (۴.۱) و (۷.۱) دیدیم که معادلات نرمال علاوه بر پارامتر  $\beta$ ، به پارامترهای همبستگی  $\alpha$  و مقیاس<sup>۱۳</sup>  $\phi$  نیز وابسته است. لذا به جهت این که معادلات نرمال به صورت تابعی از  $\beta$  به دست آیند، لازم است که  $\alpha$  و  $\phi$  برآورد شوند. یکی از راه‌حل‌ها، محاسبه برآورد گشتاوری آن‌ها به صورت زیر است

$$\hat{\phi} = \frac{1}{N-p} \sum_{i=1}^n \sum_{j=1}^{t_i} R_{ij}^2 \quad (14.1)$$

که در آن،  $R_{ij} = \frac{y_{ij} - E(y_{ij})}{\sqrt{\text{var}(y_{ij})}}$  باقیمانده‌های پیرسون و  $N = \sum_{i=1}^n t_i$  تعداد کل مشاهدات است.

برآورد  $\alpha$  وابسته به انتخاب ساختارهای همبستگی است. به‌عنوان مثال برای همبستگی‌های متبادله‌پذیر، غیر ساختاری و اتورگرسیو مرتبه اول، این برآوردگرها به ترتیب به صورت زیر

<sup>12</sup>Auto regressive

<sup>13</sup>Scale parameter

$$\hat{\alpha} = \hat{\phi}^{-1} \sum_{i=1}^n \frac{1}{n_i(n_i - 1)} \sum_{j \neq j'} R_{ij} R_{ij'},$$

$$\hat{\alpha} = \hat{\phi}^{-1} \frac{1}{n} \sum_{i=1}^n R_{ij} R_{ij'},$$

$$\hat{\alpha} = \hat{\phi}^{-1} \sum_{i=1}^n \frac{1}{n_i - 1} \sum_{j \leq n_i - 1} R_{ij} R_{ij+1},$$

که در آن ها  $R_{ij}$  مؤلفه سطر  $i$ ام و ستون  $j$ ام ماتریس همبستگی می باشد. با به دست آمدن برآوردهای فوق در نهایت معادلات نرمال تابعی بر حسب  $\beta$  است. برای به دست آوردن  $\hat{\beta}$  از روش نیوتن-رافسون<sup>۱۴</sup> استفاده می کنیم. در زیر الگوریتم محاسبه  $\hat{\beta}$  آورده شده است.

### الگوریتم ۱ الگوریتم نیوتن-رافسون در برآورد پارامتر مدل حاشیه ای

**گام ۱)** قرار می دهیم  $r = 0$  و مقدار اولیه  $\beta^0$  را برآورد LS یا ML مدل خطی کلاسیک تحت فرض استقلال محاسبه می کنیم.

**گام ۲)** پارامتر مقیاس  $\phi$  و ماتریس همبستگی عملی را براساس مانده های پیرسن و انتخاب یکی از ساختارهای معرفی شده، به ازای  $\beta^0$  محاسبه می کنیم.

**گام ۳)** برآورد ماتریس کواریانس را از رابطه زیر محاسبه می کنیم

$$\hat{V}_i(\alpha) = \hat{A}_i^{\forall} R_i(\hat{\alpha}) \hat{A}_i^{\forall}.$$

که در آن

$$\hat{A}_i = A_i(\hat{\phi})$$

**گام ۴)** مقدار  $\beta^{r+1}$  را با استفاده از رابطه زیر محاسبه می کنیم

$$\beta^{r+1} = \beta^r + \left( \sum_{i=1}^n X_i \hat{V}_i(\alpha) X_i^T \right)^{-1} \left( \sum_{i=1}^n X_i \hat{V}_i(\alpha) (y_i - X_i^T \beta^r) \right)$$

**گام ۵)** گام های (۲) تا (۴) را تا رسیدن به همگرایی ( $|\beta^{r+1} - \beta^r| < \gamma$ ) به ازای یک مقدار  $\gamma$  کوچک و از پیش تعیین شده، تکرار می کنیم.

<sup>14</sup>Newton-Raphson method

### ۲.۳.۱ مدل‌های انتقال

مدل انتقال، بسط مدل‌های خطی تعمیم‌یافته<sup>۱۵</sup> (که در بخش ۵.۱ در مورد آن‌ها صحبت خواهیم کرد) برای توصیف توزیع شرطی بر روی هر پاسخ  $y_{ij}$  به‌عنوان تابعی معلوم از متغیرهای پاسخ قبلی  $y_{i1}, y_{i2}, \dots, y_{i(j-1)}$  و متغیر پیش‌بین می‌باشد. مدل‌های انتقال یا مدل‌های مارکوف، مسئله همبستگی موجود هر واحد نمونه را با مدل‌بندی پاسخ در زمان حال به شرط پاسخ در زمان‌های گذشته مورد توجه قرار می‌دهند. به عبارت دیگر در مدل‌های خطی تعمیم‌یافته، اگر مشاهدات گذشته برای واحد  $i$ ام را با

$$H_{ij} = \{y_{i1}, y_{i2}, \dots, y_{i(j-1)}\}$$

و امید شرطی مشاهده  $i$ ام را با

$$\mu_{ij} = E(y_{ij}|H_{ij})$$

و واریانس شرطی آن را با

$$\nu = \text{var}(y_{ij}|H_{ij})$$

نمایش دهیم، یک مدل انتقال با سه مشخصه مشخص می‌شود به‌طوری که  $y_{ij}|H_{ij}$  از توزیعی از خانواده توزیع‌های نمایی پیروی می‌کند

$$f(y_{ij}|H_{ij}) = \exp \left\{ \frac{y_{ij}\theta_{ij} - b(\theta_{ij})}{\varphi} + c(y_{ij}, \varphi) \right\}$$

$$\mu_{ij} = \frac{db(\theta_{ij})}{d\theta_{ij}}$$

$$h(\mu_{ij}) = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \sum_{r=1}^s f_r(\mathbf{H}_{ij}; \boldsymbol{\alpha})$$

$$\nu_{ij} = \varphi \nu(\mu_{ij})$$

که در آن  $h$ ،  $\nu$  و  $f_r$  توابعی مشخص هستند و  $\boldsymbol{\alpha}$  برداری از پارامترها است. رایج‌ترین مدل مورد استفاده در مدل‌های انتقال استفاده از زنجیره مارکوف<sup>۱۶</sup> است که در آن پاسخ فعلی  $y_{ij}$  به شرط  $H_{ij}$  تنها به  $q$  پاسخ پیش از خود وابسته است که از آن به‌عنوان مرتبه مدل یاد می‌شود. به‌عنوان مثال در یک مدل انتقال با مرتبه  $q = 1$  داریم

$$h(\mu_{ij}) = h(E(y_{ij}|y_{i(j-1)})) = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \alpha(\mathbf{y}_{i(j-1)} - \mathbf{x}_{i(j-1)}^T \boldsymbol{\beta}).$$

### ۳.۳.۱ مدل‌های آمیخته خطی

در بسیاری از تحقیقات و مطالعات پزشکی، پزشک بیماران مورد نظر را در طول زمان به دفعات مورد معاینه قرار می‌دهد. داده‌ها و اطلاعات حاصل از معاینه که مکرر ثبت خواهند شد، مثالی

<sup>15</sup>Generalized linear models

<sup>16</sup>Markov chain

از داده‌های طولی است. یکی از رویکردهای تحلیل این نوع داده‌ها استفاده از مدل‌های آمیخته است. مدل‌های آمیخته ترکیبی از مدل با اثرات تصادفی<sup>۱۷</sup> و مدل با اثرات ثابت<sup>۱۸</sup> هستند. برای درک بهتر این مفهوم به مثال زیر از مک کولاک و همکاران (۲۰۰۸) توجه کنید. فرض کنید در یک آزمایش بالینی<sup>۱۹</sup> در یک کلینیک نوعی دارو در ۴ دُز مختلف تهیه شده‌است که اثر دُز مصرفی بیماران تشنجی بر تعداد تشنج‌های آن‌ها مورد مطالعه است. در مدل زیر  $y_{ij}$  تعداد تشنج‌های  $j$ امین شخصی که  $i$ امین دُز را دریافت کرده،  $\mu$  میانگین کل و  $\alpha_i$  اثر دُز  $i$ ام روی متوسط تعداد تشنج‌ها است.

$$E[y_{ij}] = \mu + \alpha_i, \quad i = 1, \dots, 4. \quad (15.1)$$

واضح است که  $\alpha_i$ ها اثرات ثابت هستند و هدف تفاوت بین دزهای مختلف و اثر آن‌ها روی کاهش تعداد تشنج‌های بیمار است، اما آنچه حائز اهمیت است این است که اگر مدلی ترکیب خطی از عامل‌هایی با اثرات ثابت باشد، یک مدل خطی با اثرات ثابت است. حال فرض کنید همین آزمایش را در ۲۰ مرکز کلینیک یک شهر انجام دهیم و بیمار تنها دُز شماره ۱ را دریافت کند. در این صورت  $y_{ij}$  نشان دهنده تعداد تشنج‌های  $j$ امین شخصی است که در  $i$ امین کلینیک دارو را دریافت کرده و  $\alpha_i$  اثر کلینیک  $i$ ام روی متوسط تعداد تشنج‌ها است، لذا مدل مورد بررسی به صورت زیر است

$$E[y_{ij}] = \mu + a_i, \quad i = 1, \dots, 20. \quad (16.1)$$

کلینیک  $i$ ام کلینیکی است که شخص بیمار به‌طور تصادفی از بین تمامی ۲۰ کلینیک موجود در شهر برای مراجعه به آن انتخاب می‌کند؛ لذا مراجعه هر شخص به هر مرکز کلینیک به‌طور تصادفی رخ خواهد داد، پس  $a_i$ ها اثرات تصادفی هستند. در حالتی دیگر فرض کنید در هر کلینیک برخلاف حالت قبل هر ۴ دُز دارو موجود باشد و شخصی که به‌طور تصادفی به یکی از ۲۰ کلینیک مراجعه می‌کند از بین ۴ دُز مختلف دارو آنچه را که مربوط به خود است دریافت کند. در این صورت خواهیم داشت

$$E[y_{ijk}] = \mu + \alpha_i + \beta_j + c_{ij} \quad (17.1)$$

به‌طوری که  $y_{ijk}$  تعداد تشنج‌های  $k$ امین شخصی است که دُز  $j$ ام را در کلینیک  $i$ ام دریافت کرده‌است و  $\alpha_i$  اثر کلینیک  $i$ ام روی متوسط تعداد تشنج‌های مربوط به کلینیک  $i$ ام و  $\beta_j$  اثر دُز  $j$ ام روی متوسط تعداد تشنج‌های مربوط به دُز  $j$ ام و  $c_{ij}$  اثر متقابل<sup>۲۰</sup> بین کلینیک  $i$ ام و دُز  $j$ ام است. اکنون با توجه به توضیحات قبل واضح است که  $\alpha_i$ ها اثرات تصادفی و  $\beta_j$ ها اثرات ثابت هستند و چون  $c_{ij}$  تقابلی بین اثر تصادفی و ثابت است، پس  $c_{ij}$  نیز یک اثر تصادفی است.

<sup>17</sup>Random effects

<sup>18</sup>Fixed effects

<sup>19</sup>Clinical trial

<sup>20</sup>Interaction



بنابراین مدل (۱۷.۱) آمیزه‌ای از اثرات تصادفی و ثابت است و به آن مدل آمیخته می‌گویند. همچنین چون ترکیبی خطی از این اثرات است به آن مدل آمیخته خطی گویند. از آن جایی که تمرکز این پایان‌نامه بر روی مدل‌های آمیخته خطی است لذا در بخش بعدی به طور جداگانه به معرفی این مدل خواهیم پرداخت.

## ۴.۱ مدل آمیخته خطی

مدل خطی  $E[y|X] = X\beta$  را در نظر بگیرید که در آن بردار مؤلفه‌های پاسخ،  $X = (X_1^T, \dots, X_n^T)^T$  ماتریس طرح در اندازه  $N \times p$  و  $\beta = (\beta_1, \dots, \beta_p)^T$  بردار اثرات ثابت است. تنها تفاوت مدل آمیخته خطی با مدل خطی افزودن جمله  $Zb$  است که در آن  $Z$  همانند  $X$  یک ماتریس معلوم از متغیرهای پیش‌بین و  $b = (b_1^T, \dots, b_n^T)^T$  بردار اثرات تصادفی در رخ داد متغیر پاسخ  $y$  است که در آن  $b_i$ ها بردارهای  $q \times 1$  هستند، لذا مدل آمیخته خطی را می‌توان به صورت زیر نوشت:

$$E[y|X, Z] = X\beta + Zb. \quad (18.1)$$

در مطالعات طولی که هر مشاهده برای هر فرد به دفعات اندازه‌گیری می‌شود، شکل ماتریسی مدل آمیخته خطی با در نظر گرفتن جمله خطاها به صورت زیر است

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{N \times 1} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix}_{N \times p} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}_{p \times 1} + \begin{pmatrix} Z_1^T & \circ & \dots & \circ \\ \circ & Z_2^T & \dots & \circ \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \dots & Z_n^T \end{pmatrix}_{N \times nq} \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}_{nq \times 1} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}_{N \times 1}$$

که می‌توان به صورت زیر نمایش داد

$$y_i = X_i\beta + Z_i^T b_i + \epsilon_i, \quad i = 1, \dots, n, \quad (19.1)$$

که در آن  $y_i = (y_{i1}, y_{i2}, \dots, y_{in_i})^T$  برداری از  $n_i$  مشاهده برای فرد  $i$ ام،  $X_i$  ماتریس طرح  $n_i \times p$  از اثرات ثابت و  $\beta = (\beta_1, \dots, \beta_p)^T$  بردار  $p \times 1$  ضرایب ثابت و نامعلوم، و  $Z_i$  ماتریس طرح  $n_i \times q$  از اثرات تصادفی،  $b_i$  بردار  $q \times 1$  ضرایب تصادفی و  $\epsilon_i$  بردار  $n \times 1$  خطاهای تصادفی است. همچنین فرض می‌شود به ازای هر  $i = 1, \dots, n$ ،  $b_i$  از  $\epsilon_i$  مستقل بوده و

$$\begin{cases} E(b_i) = \circ, & E(b_i b_i^T) = D \\ E(\epsilon_i) = \circ, & E(\epsilon_i \epsilon_i^T) = \sigma^2 I_{n_i} \end{cases}$$

با فرض نرمال بودن توزیع  $b_i$ ها و  $\epsilon_i$ ها، توزیع متغیر پاسخ به صورت زیر نتیجه می‌شود

$$y \sim \mathcal{N}_N(X\beta, V) \quad (20.1)$$

که در آن  $\mathbf{X} = (\mathbf{X}_1^T, \dots, \mathbf{X}_n^T)^T$  و  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$

$$\begin{aligned} \mathbf{V} &= \text{var}(\mathbf{y}) = \text{var}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}^T \mathbf{b} + \boldsymbol{\epsilon}) \\ &= \text{var}(\mathbf{Z}^T \mathbf{b}) + \text{var}(\boldsymbol{\epsilon}) \\ &= \mathbf{ZDZ}^T + \sigma^2 \mathbf{I}. \end{aligned} \quad (21.1)$$

### ۱.۴.۱ برآورد کلاسیک پارامترهای مدل

در روش ML، برآورد پارامترها با ماکزیمم کردن تابع درستنمایی مدل محاسبه می‌شود. در این بخش برآورد اثرات ثابت مدل (۱۹.۱) را به این روش محاسبه می‌کنیم. لازم به ذکر است در طول این پایان‌نامه تمرکز بر روی برآورد اثرات ثابت  $\boldsymbol{\beta}$  می‌باشد. برای برآورد اثرات تصادفی  $\mathbf{b}$ ، مک کولاک و همکاران (۲۰۰۸) را ببینید. تحت مفروضات مدل (۲۰.۱) داریم

$$f(\mathbf{y}; \boldsymbol{\beta}; \mathbf{V}) = \frac{|\mathbf{V}|^{-\frac{1}{2}}}{(2\pi)^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}.$$

بنابراین لگاریتم تابع درستنمایی عبارتست از

$$l(\mathbf{y}) = -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \log |\mathbf{V}| - \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (22.1)$$

واضح است که ماکزیمم کردن تابع فوق برحسب  $\boldsymbol{\beta}$  معادل با مینیمم کردن عبارت وابسته به  $\boldsymbol{\beta}$  است، در نتیجه

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\text{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}. \quad (23.1)$$

با مساوی صفر قرار دادن مشتق رابطه (۲۳.۱) بر حسب  $\boldsymbol{\beta}$  خواهیم داشت

$$\begin{aligned} \circ &= \frac{d}{d\boldsymbol{\beta}} \left[ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \\ &= \frac{d}{d\boldsymbol{\beta}} \left[ \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} - \mathbf{y}^T \mathbf{V}^{-1} \mathbf{X}\boldsymbol{\beta} - (\mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} \mathbf{y} + (\mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{X}\boldsymbol{\beta}) \right] \\ &= \frac{d}{d\boldsymbol{\beta}} \left[ \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} - 2\mathbf{y}^T \mathbf{V}^{-1} \mathbf{X}\boldsymbol{\beta} + (\mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{X}\boldsymbol{\beta}) \right] \\ &= -2\mathbf{y}^T \mathbf{V}^{-1} \mathbf{X} + 2\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}. \end{aligned}$$

با شرط معلوم بودن  $\mathbf{V}$  معکوس پذیر بودن  $(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})$  داریم

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} \quad (24.1)$$

به همین ترتیب نتیجه می‌شود

$$E(\hat{\beta}) = \beta \quad (25.1)$$

و

$$\begin{aligned} \text{var}(\hat{\beta}) &= \text{var}((\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{Y}) \\ &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}. \end{aligned} \quad (26.1)$$

واضح است که در رابطه (۲۴.۱) محاسبه مقدار  $\mathbf{V}$ ، مستلزم برآورد مؤلفه‌های واریانس یعنی  $\theta = (\sigma, \mathbf{D})$  است. بدین منظور لیرد و ویر (۱۹۸۲) الگوریتم EM را برای برآورد پارامتر  $\beta$  پیشنهاد داد که در ادامه بیان خواهیم کرد.

الگوریتم ۲ الگوریتم EM در برآورد پارامتر مدل آمیخته خطی

گام (۱) قرار می‌دهیم  $\hat{\theta}^{(t)} = \theta$ ، سپس مقدار آماره‌های  $\hat{t}_\psi^{(t)}$  و  $\hat{t}_\gamma^{(t)}$  را برای مرحله  $t$ ام به صورت زیر محاسبه می‌کنیم:

$$\begin{aligned} \hat{t}_\gamma^{(t)} &= E \left( \sum_{i=1}^n \epsilon_i^T \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \epsilon_i^T(\hat{\theta}^{(t)}) \epsilon_i(\hat{\theta}^{(t)}) + \text{tr} \text{Var}(\epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)}) \right], \end{aligned} \quad (27.1)$$

$$\begin{aligned} \hat{t}_\psi^{(t)} &= E \left( \sum_{i=1}^n \mathbf{b}_i \mathbf{b}_i^T \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \mathbf{b}_i(\hat{\theta}^{(t)}) \mathbf{b}_i^T(\hat{\theta}^{(t)}) + \text{Var}(\mathbf{b}_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)}) \right]. \end{aligned} \quad (28.1)$$

گام (۲) برای به دست آوردن  $\hat{\theta}^{(t+1)}$ ، مقدار  $\hat{\sigma}^2(t+1)$  و  $\hat{\mathbf{D}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم:

$$\hat{\sigma}^2(t+1) = \hat{t}_\gamma^{(t)} / N,$$

$$\hat{\mathbf{D}}^{(t+1)} = \hat{t}_\psi^{(t)} / n.$$

که در آن  $N = \sum_{i=1}^n n_i$  و  $n$  برابر تعداد واحدهای آزمایشی است. حال برای محاسبه  $\hat{\mathbf{V}}^{(t+1)}$  داریم

$$\hat{\mathbf{V}}^{(t+1)} = \mathbf{Z} \hat{\mathbf{D}}^{(t+1)} \mathbf{Z}^T + \hat{\sigma}^2(t+1) \mathbf{I}. \quad (29.1)$$

گام (۳) در نهایت مقدار  $\hat{\beta}$  از رابطه زیر بدست می‌آید

$$\hat{\beta}^{(t+1)} = (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{y}. \quad (30.1)$$

گام (۴) گام‌های (۱) تا (۳) را تا رسیدن به همگرایی ( $|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}| < \gamma$ ) به ازای یک مقدار  $\gamma$  کوچک و از پیش تعیین شده، تکرار می‌کنیم.

در ادامه چگونگی محاسبه برآوردگر  $\beta$  را در قالب یک مثال شهودی با استفاده از نرم‌افزار R بررسی می‌کنیم.

### ۲.۴.۱ مثال کاربردی

در این بخش یک نمونه از داده‌های طولی موجود در بسته nlme در نرم‌افزار R به نام Orthodont را مورد بررسی قرار می‌دهیم. این داده‌ها توسط محققان دانشکده دندان پزشکی دانشگاه

کارولینای شمالی<sup>۲۱</sup> جمع‌آوری شده‌اند و شامل اطلاعات رشد ۲۷ کودک است (۱۶ پسر و ۱۱ دختر) که از سن ۸ تا ۱۴ سالگی هر دو سال یک‌بار فاصله بین مرکز غده هیپوفیز<sup>۲۲</sup> تا شکاف پریگوماکسیلار<sup>۲۳</sup> واقع در جمجمه انسان اندازه‌گیری شده‌است. ۸ سطر ابتدایی داده‌ها به‌صورت است

	distance	age	Subject	Sex
1	26.0	8	M01	Male
2	25.0	10	M01	Male
3	29.0	12	M01	Male
4	31.0	14	M01	Male
5	21.5	8	M02	Male
6	22.5	10	M02	Male
7	23.0	12	M02	Male
8	26.5	14	M02	Male

متغیر Subject در این مجموعه داده یک عامل برای شناسایی افراد است (مثلاً M01 یعنی پسر شماره ۱ و F01 یعنی دختر شماره ۱). لذا به‌عنوان اثر تصادفی در مدل در نظر گرفته می‌شود، زیرا هر کودک به‌طور تصادفی برای انجام آزمایش انتخاب شده‌است. لذا مدل آمیخته خطی را به داده‌ها برازش می‌دهیم. اما قبل از هر چیز ممکن است یک سوال مطرح شود که چرا باید اثرات تصادفی را در مدل وارد کنیم؟

در ادامه پاسخ این سوال و دلیل این امر را با برازش یک مدل خطی به داده‌ها بدون در نظر گرفتن اثر تصادفی و مقایسه نتایج آن با مدل آمیخته خطی بیان می‌کنیم. برای این منظور ابتدا مدل (۲.۱) و سپس مدل (۱۸.۱) را با استفاده از بسته nlme در R به داده‌ها برازش می‌دهیم. کدهای مربوطه در ضمیمه آمده‌است.

نتایج خروجی این مدل در جدول ۲.۱ خلاصه شده‌است.

<sup>21</sup>North Carolina

<sup>22</sup>Pituitary

<sup>23</sup>Pterygomaxillary fissure

جدول ۲.۱: برآورد ضرایب با استفاده از مدل خطی و آمیخته خطی

مدل خطی			
متغیر	برآورد	انحراف استاندارد	P-value
عرض از مبدأ	۱۶/۳۴۰۶	۱/۴۱۶۲	< ۰/۰۰۰۱
سن	۰/۷۸۴۴	۰/۱۲۶۲	< ۰/۰۰۰۱
جنسیت	۱/۰۳۲۱	۲/۲۱۸۸	۰/۶۴۳۰
اثر متقابل سن و جنسیت	-۰/۳۰۴۸	۰/۱۹۷۷	۰/۱۲۶۰
مدل آمیخته خطی			
متغیر	برآورد	انحراف استاندارد	P-value
عرض از مبدأ	۱۶/۳۴۰۶	۰/۹۸۱۳	< ۰/۰۰۰۱
سن	۰/۷۸۴۳	۰/۰۷۷۵	< ۰/۰۰۰۱
جنسیت	۱/۰۳۲۱	۱/۵۳۷۴	۰/۵۰۸۲
اثر متقابل سن و جنسیت	-۰/۳۰۴۸	۰/۱۲۱۴	۰/۰۱۴۱

با توجه به نتایج جدول ۲.۱ در خروجی دو مدل دیده می‌شود که برآورد پارامترهای هر دو مدل بسیار بهم نزدیک هستند اما آنچه حائز اهمیت است میزان خطای استاندارد برآورد پارامترها می‌باشد که در مدل آمیخته خطی کاهش قابل ملاحظه‌ای داشته‌است و چنانچه اثر تصادفی در مدل نادیده گرفته شود، اثر متقابل دو متغیر سن و جنسیت معنا دار نخواهد بود.

## ۵.۱ مدل آمیخته خطی تعمیم یافته

در بخش قبل دیدیم که در استنباط مبتنی بر درست‌نمایی در مدل خطی فرض بر این است که متغیر پاسخ از توزیع نرمال پیروی می‌کند، اما در برخی موارد که متغیر پاسخ دو سطحی<sup>۲۴</sup> یا شمارشی<sup>۲۵</sup> است استفاده از مدل‌های خطی امکان‌پذیر نیست چرا که متغیر پاسخ از توزیع نرمال تبعیت نمی‌کند. در این صورت اگر متغیر پاسخ دارای توزیعی غیر از توزیع نرمال باشد از توزیع‌های متعلق به خانواده توزیع‌های نمایی در مدل‌های خطی تعمیم یافته، استفاده می‌شود. همچنین در مواردی از جمله توزیع پواسن که اغلب برای داده‌های شمارشی به کار می‌رود، میانگین توزیع پاسخ باید نامنفی باشد. در این حالت استفاده از مدل‌های خطی ممکن است برآورد منفی را برای میانگین پاسخ در پی داشته باشد. لذا در چنین مواردی با استفاده از

<sup>24</sup>Binary

<sup>25</sup>Count

مدل‌های خطی تعمیم‌یافته تابعی از میانگین پاسخ در مدل‌سازی استفاده می‌شود. بنابراین به سبب تحلیل این دسته از داده‌ها، مدل‌های خطی تعمیم‌یافته به عنوان تعمیمی از مدل‌های خطی به کار گرفته شده که با استفاده از یک تابع پیوند، بین میانگین پاسخ‌ها و متغیرهای بیش‌بین ارتباط برقرار می‌شود.

### ۱.۵.۱ مؤلفه‌های مدل خطی تعمیم‌یافته

یک مدل خطی تعمیم‌یافته دارای سه مؤلفه به صورت زیر است:

#### مؤلفه تصادفی

این مؤلفه توزیع متغیر پاسخ را مشخص می‌کند که باید از خانواده توزیع‌های نمایی باشد. خانواده توزیع‌های نمایی یک مؤلفه مهم در مدل‌های خطی تعمیم‌یافته است که تابع چگالی این خانواده را می‌توان به صورت زیر بیان کرد

$$f(y; \gamma, \tau) = \exp \left\{ (y\gamma - b(\gamma)) / \tau^2 - c(y, \tau) \right\} \quad (31.1)$$

که در آن  $b(\cdot)$  و  $c(\cdot)$  توابعی معلوم هستند،  $b(\cdot)$  یک تابع انباشته<sup>۲۶</sup>، پارامتر  $\gamma$  یک پارامتر مکانی<sup>۲۷</sup> و  $\tau$  را اغلب پارامتر پراکندگی<sup>۲۸</sup> می‌نامند. همچنین میانگین و واریانس خانواده توزیع‌های نمایی به صورت زیر محاسبه می‌شوند:

$$E(y) = \mu = \frac{db(\gamma)}{d\gamma} \quad (32.1)$$

$$Var(y) = \frac{d^2 b(\gamma)}{d\gamma^2} \tau^2. \quad (33.1)$$

از جمله توزیع‌های مربوط به خانواده توزیع‌های نمایی می‌توان به نرمال، پواسن، هندسی، نمایی، گاما، نرمال وارون<sup>۲۹</sup> و دو جمله‌ای منفی اشاره کرد. به عنوان مثال، تابع احتمال توزیع پواسن را در نظر بگیرید

$$f(y; \mu) = \frac{e^{-\mu} \mu^y}{y!}.$$

این تابع احتمال را می‌توان به صورت خانواده توزیع‌های نمایی نوشت:

$$f(y; \mu) = \exp \{ y \ln \mu - \mu - \ln y! \}$$

<sup>26</sup>Cumulant

<sup>27</sup>Location parameter

<sup>28</sup>Dispersion parameter

<sup>29</sup>Inverse normal

که در آن  $b(\gamma) = e^\gamma$ ،  $\gamma = \ln \mu$ ،  $\tau^2 = 1$  و  $c(y; \tau) = \exp\{-\ln(y!)\}$  است. برای محاسبه میانگین و واریانس آن نیز با استفاده از روابط (۳۲.۱) و (۳۳.۱) داریم

$$E(y) = \mu = \frac{db(\gamma)}{d\gamma} = \frac{de^\gamma}{d\gamma} = e^\gamma = e^{\ln \mu} = \mu$$

$$Var(y) = \frac{d^2 b(\gamma)}{d\gamma^2} \tau^2 = \frac{d^2 e^\gamma}{d\gamma^2} = e^\gamma = e^{\ln \mu} = \mu.$$

## مؤلفه اصلی

این مؤلفه دارای یک پیش‌گوی خطی<sup>۳۰</sup> است که متغیرهای پیش‌بین را شامل شده و به صورت زیر بیان می‌شود

$$\eta_i = \mathbf{x}_i^T \boldsymbol{\beta} = \beta_0 + \sum_{i=1}^p \beta_i x_i.$$

که در آن  $\mathbf{x}_i = (x_1, \dots, x_p)^T$

## تابع پیوند

این مؤلفه پیش‌گوی خطی را به میانگین متغیر پاسخ مربوط می‌کند. به طوری که

$$\eta_i = g(E(\mathbf{y}_i | \mathbf{X}_i)) = \mathbf{x}_i^T \boldsymbol{\beta}$$

که در آن  $g(\cdot)$  را تابع پیوند<sup>۳۱</sup> گویند که یک تابع یکنوای مشتق‌پذیر و همچنین معکوس‌پذیر است. لذا می‌توان آن را به صورت زیر بیان کرد:

$$\mu_i = E(\mathbf{y}_i | \mathbf{x}_i) = g^{-1}(\eta_i) = g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}). \quad (34.1)$$

تابع پیوند انواع گوناگونی دارد که از جمله آن‌ها می‌توان به تابع پیوند همانی<sup>۳۲</sup>،  $g(p) = p$ ، تابع پیوند لگاریتمی،  $g(p) = \log(p)$ ، تابع پیوند لجیت<sup>۳۳</sup>،  $g(p) = \log\left(\frac{p}{1-p}\right)$ ، تابع پیوند پرابیت<sup>۳۴</sup>،  $g(p) = \Phi^{-1}(p)$  و تابع پیوند Cloglog<sup>۳۵</sup>،  $g(p) = \log(-\log(1-p))$ ، اشاره کرد.

با انتخاب تابع پیوند  $g$ ، الگوهای رگرسیون متفاوتی همچون رگرسیون لجستیک، رگرسیون پرابیت و رگرسیون پواسن به عنوان یکی از اعضای رده مدل‌های خطی تعمیم‌یافته در مدل‌بندی متغیرهای پاسخ استفاده می‌شود. واضح است که با در نظر گرفتن تابع پیوند همانی، الگوی رگرسیون خطی استاندارد نیز یک مدل خطی تعمیم‌یافته است.

<sup>30</sup>Linear predictor

<sup>31</sup>Link function

<sup>32</sup>Identity

<sup>33</sup>Logit

<sup>34</sup>Probit

<sup>35</sup>Complementary Log-Log



یکی از فرضیات مدل‌های خطی تعمیم‌یافته استقلال بردار پاسخ است. چنان‌چه این فرض برقرار نباشد از مدل‌های خطی آمیخته تعمیم‌یافته استفاده می‌شود. به عبارت دقیق‌تر، می‌توان گفت که در مدل (۱۸.۱)، رفتار متغیر پاسخ، استفاده از مدل‌های تعمیم‌یافته را به جهت موارد گفته‌شده، ایجاب می‌کند. لذا یکی از روش‌های تحلیل آن دسته از مجموعه داده‌های طولی که بردار پاسخ آن‌ها دارای چنین رفتاری است، مدل آمیخته خطی تعمیم‌یافته خواهد بود. زیرا همبستگی بین مشاهدات هر واحد در بردار پاسخ داده‌های طولی وجود دارد. با این توضیحات واضح است که تفاوت مدل‌های خطی تعمیم‌یافته با مدل‌های آمیخته خطی تعمیم‌یافته در مؤلفه اصلی خواهد بود، به طوری که داریم

$$E(\mathbf{y}_i | \mathbf{b}_i) = \mu_i$$

$$\eta_i = g(\mu_i) = \mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{Z}_i^T \mathbf{b}_i \quad (۳۵.۱)$$

$$\begin{aligned} E(\mathbf{y}_i) &= E(E(\mathbf{y}_i | \mathbf{b}_i)) \\ &= E(\mu_i) \\ &= E(g^{-1}(\mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{Z}_i^T \mathbf{b}_i)). \end{aligned}$$

## ۲.۵.۱ برآورد پارامترهای مدل

### مدل خطی تعمیم‌یافته

مدل خطی زیر را در نظر بگیرید

$$E[\mathbf{y} | \mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$$

که در آن بردار مؤلفه‌های پاسخ،  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$  بردار مؤلفه‌های پاسخ،  $\mathbf{X} = (\mathbf{X}_1^T, \dots, \mathbf{X}_n^T)^T$  ماتریس طرح در اندازه  $N \times p$  و  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  بردار اثرات ثابت است. با فرض این که توزیع متغیر پاسخ از خانواده توزیع‌های نمایی باشد، داریم

$$\mathbf{y}_i \stackrel{iid}{\sim} f_{\mathbf{y}_i}(\mathbf{y}_i)$$

که در آن

$$f_{\mathbf{y}_i}(\mathbf{y}_i) = \exp \left\{ (\mathbf{y}_i \boldsymbol{\gamma}_i - b(\boldsymbol{\gamma}_i)) / \tau^2 - c(\mathbf{y}_i, \tau) \right\} \quad (۳۶.۱)$$

لذا لگاریتم تابع درست‌نمایی به صورت زیر است

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^n [\mathbf{y}_i \boldsymbol{\gamma}_i - b(\boldsymbol{\gamma}_i)] / \tau^2 - \sum_{i=1}^n c(\mathbf{y}_i, \tau). \quad (۳۷.۱)$$

قبل از ماکزیمم کردن تابع (۳۷.۱)، توجه به موارد زیر که در ادامه از آن‌ها بهره می‌گیریم ضروری است.

$$E \left[ \frac{\partial \log f_{\mathbf{y}_i}(\mathbf{y}_i)}{\partial \gamma_i} \right] = 0 \quad (38.1)$$

و

$$Var \left( \frac{\partial \log f_{\mathbf{y}_i}(\mathbf{y}_i)}{\partial \gamma_i} \right) = -E \left[ \frac{\partial^2 \log f_{\mathbf{y}_i}(\mathbf{y}_i)}{\partial \gamma_i^2} \right]. \quad (39.1)$$

تحت شرایط نظم (کسلا و برگر (۱۹۹۰))، با استفاده از رابطه‌های (۳۶.۱) و (۳۸.۱) داریم

$$E \left[ \left\{ \mathbf{y}_i - \frac{\partial b(\gamma_i)}{\partial \gamma_i} \right\} / \tau^2 \right] = 0, \quad (40.1)$$

$$E[\mathbf{y}_i] = \mu_i = \frac{\partial b(\gamma_i)}{\partial \gamma_i}. \quad (41.1)$$

همچنین با استفاده از رابطه‌های (۳۶.۱) و (۳۹.۱) داریم

$$Var \left( \left\{ \mathbf{y}_i - \frac{\partial b(\gamma_i)}{\partial \gamma_i} \right\} / \tau^2 \right) = -E \left[ -\frac{1}{\tau^2} \frac{\partial^2 b(\gamma_i)}{\partial \gamma_i^2} \right], \quad (42.1)$$

حال با استفاده از رابطه (۴۱.۱) می‌توان نوشت

$$\begin{aligned} Var \left( \frac{\mathbf{y}_i - \mu_i}{\tau^2} \right) &= \frac{1}{\tau^2} \frac{\partial^2 b(\gamma_i)}{\partial \gamma_i^2} \\ Var(\mathbf{y}_i) &= \tau^2 \frac{\partial^2 b(\gamma_i)}{\partial \gamma_i^2} \end{aligned} \quad (43.1)$$

با قرار دادن  $\nu(\mu_i) = \partial^2 b(\gamma_i) / \partial \gamma_i^2$  داریم

$$Var(\mathbf{y}_i) = \tau^2 \nu(\mu_i).$$

همچنین

$$\frac{\partial \gamma_i}{\partial \mu_i} = \left( \frac{\partial \mu_i}{\partial \gamma_i} \right)^{-1} = \left( \frac{\partial^2 b(\gamma_i)}{\partial \gamma_i^2} \right)^{-1} = \frac{1}{\nu(\mu_i)} \quad (44.1)$$

و طبق رابطه (۳۴.۱) داریم

$$\begin{aligned} \frac{\partial \mu_i}{\partial \beta} &= \frac{\partial \mu_i}{\partial g(\mu_i)} \frac{\partial g(\mu_i)}{\partial \beta} = \left( \frac{\partial g(\mu_i)}{\partial \mu_i} \right)^{-1} \frac{\partial \mathbf{X}_i^T \beta}{\partial \beta} \\ &= \left( \frac{\partial g(\mu_i)}{\partial \mu_i} \right)^{-1} \mathbf{X}_i^T. \end{aligned} \quad (45.1)$$

حال برای ماکزیمم کردن تابع (۳۷.۱) برحسب  $\beta$  داریم

$$\begin{aligned} \frac{\partial \ell(\beta)}{\partial \beta} &= \frac{1}{\tau^2} \sum \left[ \mathbf{y}_i \frac{\partial \gamma_i}{\partial \beta} - \frac{\partial b(\gamma_i)}{\partial \gamma_i} \frac{\partial \gamma_i}{\partial \beta} \right] \\ &= \frac{1}{\tau^2} \sum (y_i - \mu_i) \frac{\partial \gamma_i}{\partial \beta} \\ &= \frac{1}{\tau^2} \sum (y_i - \mu_i) \frac{\partial \gamma_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \beta} \\ &= \frac{1}{\tau^2} \sum \frac{(y_i - \mu_i)}{\nu(\mu_i) g_\mu(\mu_i)} \mathbf{X}_i^T \\ &= \sum (y_i - \mu_i) w_i g_\mu(\mu_i) \mathbf{X}_i^T, \end{aligned} \quad (۴۶.۱)$$

که در آن  $w_i = [\tau^2 \nu(\mu_i) g_\mu^2(\mu_i)]^{-1}$ .

شکل ماتریسی رابطه (۴۶.۱) را می‌توان به صورت زیر نوشت

$$\frac{\partial \ell(\beta)}{\partial \beta} = \mathbf{X}^T \mathbf{W} \Delta (\mathbf{y} - \boldsymbol{\mu}). \quad (۴۷.۱)$$

با مساوی صفر قرار دادن رابطه (۴۷.۱) نتیجه می‌شود

$$\mathbf{X}^T \mathbf{W} \Delta \mathbf{y} = \mathbf{X}^T \mathbf{W} \Delta \boldsymbol{\mu}, \quad (۴۸.۱)$$

که در آن  $\mathbf{W} = \{n w_i\}$  و  $\Delta = \{n g_\mu(\mu_i)\}$  و  $\mu$  به پارامتر  $\beta$  وابسته می‌باشند.

برای حل معادله (۴۸.۱) بر حسب  $\beta$ ، از روش کمترین توان‌های دوم موزون تکراری<sup>۳۶</sup> بهره می‌گیریم. در این روش برای بردار پارامتر مورد نظر  $\theta$  داریم

$$\theta^{(m+1)} = \theta^{(m)} + \mathbf{I}(\theta^{(m)})^{-1} \frac{\partial \ell}{\partial \theta} \Big|_{\theta=\theta^{(m)}}, \quad (۴۹.۱)$$

که در آن  $(m)$  نشان دهنده تکرار  $m$ ام و  $\mathbf{I}(\theta)$  ماتریس اطلاع فیشر<sup>۳۷</sup> است که به صورت زیر محاسبه می‌شود

$$\mathbf{I}(\theta) = -E \left[ \frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^T} \right].$$

لذا با استفاده از رابطه (۴۷.۱) داریم

$$\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = -\mathbf{X}^T \mathbf{W} \Delta \frac{\partial \boldsymbol{\mu}^T}{\partial \beta} + \mathbf{X}^T \frac{\partial \mathbf{W} \Delta^T}{\partial \beta} (\mathbf{y} - \boldsymbol{\mu}) \quad (۵۰.۱)$$

در نتیجه

$$\begin{aligned} -E \left[ \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right] &= \mathbf{X}^T \mathbf{W} \Delta \frac{\partial \boldsymbol{\mu}^T}{\partial \beta} + \circ \\ &= \mathbf{X}^T \mathbf{W} \Delta \Delta^{-1} \mathbf{X} \\ &= \mathbf{X}^T \mathbf{W} \mathbf{X}. \end{aligned} \quad (۵۱.۱)$$

<sup>۳۶</sup>Iterative weighted least squares method

<sup>۳۷</sup>Fisher information matrix

اکنون طبق رابطه (۴۹.۱) و (۵۱.۱) برآوردگر خطی تعمیم یافته (GL) پارامتر  $\beta$  به صورت زیر حاصل می شود

$$\hat{\beta}_{GL}^{(m+1)} = \hat{\beta}_{GL}^{(m)} + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \Delta (\mathbf{y} - \mu). \quad (52.1)$$

### مدل آمیخته خطی تعمیم یافته

مدل آمیخته خطی زیر را در نظر بگیرید

$$E[\mathbf{y} | \mathbf{X}, \mathbf{Z}] = \mathbf{X}\beta + \mathbf{Z}\mathbf{b}$$

که در آن  $\mathbf{y} = (y_1^T, \dots, y_n^T)^T$  بردار مؤلفه های پاسخ،  $\mathbf{X} = (\mathbf{X}_1^T, \dots, \mathbf{X}_n^T)^T$  ماتریس طرح در اندازه  $N \times p$  و  $\beta = (\beta_1, \dots, \beta_p)^T$  بردار اثرات ثابت و  $\mathbf{Z}$  همانند  $\mathbf{X}$  یک ماتریس معلوم از متغیرهای پیش بین و  $\mathbf{b} = (b_1^T, \dots, b_n^T)^T$  بردار اثرات تصادفی در رخ داد متغیر پاسخ  $\mathbf{y}$  است. با فرض این که توزیع متغیر پاسخ از خانواده توزیع های نمایی باشد، داریم

$$y_i | b_i \stackrel{iid}{\sim} f_{y_i | b_i}(y_i | b_i), \quad (53.1)$$

که در آن

$$f_{y_i | b_i}(y_i | b_i) = \exp \left\{ (y_i \gamma_i - b(\gamma_i)) / \tau^2 - c(y_i, \tau) \right\}, \quad (54.1)$$

به طوری که

$$E[y_i | \mathbf{b}] = \mu_i$$

$$g(\mu_i) = \mathbf{X}_i^T \beta + \mathbf{Z}_i^T \mathbf{b}_i$$

$$\mathbf{b} \sim f_{\mathbf{b}}(\mathbf{b} | \mathbf{D})$$

لذا تابع درستنمایی و لگاریتم آن، به ترتیب به صورت زیر است

$$L(\beta) = \int \left( \prod_i f_{y_i | b_i}(y_i | b_i) f_{\mathbf{b}}(\mathbf{b}_i) \right) d\mathbf{b}, \quad (55.1)$$

$$\ell(\beta) = \log \int f_{\mathbf{y} | \mathbf{b}}(\mathbf{y} | \mathbf{b}) f_{\mathbf{b}}(\mathbf{b}) d\mathbf{b} = \log f_{\mathbf{y}}(\mathbf{y}). \quad (56.1)$$

با مشتق گیری از رابطه (۵۶.۱) نسبت به بردار اثرات ثابت  $(\beta)$  داریم

$$\begin{aligned} \frac{\partial \ell(\beta)}{\partial \beta} &= \frac{\partial}{\partial \beta} \int f_{\mathbf{y} | \mathbf{b}}(\mathbf{y} | \mathbf{b}) f_{\mathbf{b}}(\mathbf{b}) d\mathbf{b} / f_{\mathbf{y}}(\mathbf{y}) \\ &= \int \left[ \frac{\partial}{\partial \beta} f_{\mathbf{y} | \mathbf{b}}(\mathbf{y} | \mathbf{b}) \right] f_{\mathbf{b}}(\mathbf{b}) d\mathbf{b} / f_{\mathbf{y}}(\mathbf{y}). \end{aligned} \quad (57.1)$$

از آن جایی که  $f_b(\mathbf{b})$  به پارامتر  $\beta$  وابسته نیست، می‌توان به صورت زیر نوشت

$$\begin{aligned} \frac{\partial}{\partial \beta} f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b}) &= \left( \frac{1}{f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b})} \frac{\partial f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b})}{\partial \beta} \right) f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b}) \\ &= \frac{\partial \log f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b})}{\partial \beta} f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b}), \end{aligned} \quad (58.1)$$

حال رابطه (57.1) را به صورت زیر بازنویسی می‌کنیم

$$\begin{aligned} \frac{\partial \ell(\beta)}{\partial \beta} &= \int \frac{\partial \log f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b})}{\partial \beta} f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b}) f_{\mathbf{b}}(\mathbf{b}) d\mathbf{b} / f_{\mathbf{y}}(\mathbf{y}) \\ &= \int \frac{\partial \log f_{\mathbf{y}|\mathbf{b}}(\mathbf{y}|\mathbf{b})}{\partial \beta} f_{\mathbf{b}|\mathbf{y}}(\mathbf{b}|\mathbf{y}) d\mathbf{b}. \end{aligned} \quad (59.1)$$

با استفاده از رابطه (47.1) داریم

$$\begin{aligned} \frac{\partial \ell(\beta)}{\partial \beta} &= \int \mathbf{X}^T \mathbf{W} \Delta (\mathbf{y} - \boldsymbol{\mu}) f_{\mathbf{b}|\mathbf{y}}(\mathbf{b}|\mathbf{y}) d\mathbf{b} \\ &= \mathbf{X}^T E_{\mathbf{b}|\mathbf{y}} [\mathbf{W} \Delta | \mathbf{y}] \mathbf{y} - \mathbf{X}^T E_{\mathbf{b}|\mathbf{y}} [\mathbf{W} \Delta \boldsymbol{\mu} | \mathbf{y}] \end{aligned} \quad (60.1)$$

با مساوی صفر قرار دادن رابطه (60.1) داریم

$$\mathbf{X}^T E_{\mathbf{b}|\mathbf{y}} [\mathbf{W} \Delta | \mathbf{y}] \mathbf{y} = \mathbf{X}^T E_{\mathbf{b}|\mathbf{y}} [\mathbf{W} \Delta \boldsymbol{\mu} | \mathbf{y}] \quad (61.1)$$

در نتیجه طبق رابطه (49.1) برآوردگر آمیخته خطی تعمیم‌یافته (GLM) پارامتر  $\beta$  به صورت زیر حاصل می‌شود

$$\hat{\beta}_{GLM}^{(m+1)} = \hat{\beta}_{GLM}^{(m)} + (\mathbf{X}^T E_{\mathbf{b}|\mathbf{y}} [\mathbf{W} | \mathbf{y}] \mathbf{X})^{-1} \mathbf{X}^T E_{\mathbf{b}|\mathbf{y}} [\mathbf{W} \Delta (\mathbf{y} - \boldsymbol{\mu}) | \mathbf{y}]. \quad (62.1)$$

برای آگاهی بیشتر در مورد برآورد پارامترهای مدل آمیخته خطی تعمیم‌یافته به مک کولاک و همکاران (2008) مراجعه کنید.

## فصل ۲

# رگرسیون ريج برای داده‌های طولی

### ۱.۲ مقدمه

همان‌طور که در فصل قبل گفته شد، یکی از روش‌های مدل‌بندی داده‌های طولی استفاده از مدل‌های آمیخته خطی است. الیوت و همکاران (۲۰۱۱) برای بهبود بخشیدن به برآورد پارامترهای مدل آمیخته خطی، رگرسیون آمیخته ريج را مطرح کردند. در این فصل علاوه بر شرح مدل آن‌ها، برآورد ضرایب رگرسیونی این مدل را با خطاهای بیضی‌گون بررسی کرده‌ایم. همچنین به جهت بهبود خطای برآوردگر مدل رگرسیون آمیخته ريج، برآوردگرهای انقباضی نوع استاین را در این مدل ارائه کرده و نتایج شبیه‌سازی و تحلیل مثال واقعی را گزارش می‌کنیم. انتهای این فصل در بر گیرنده تعمیم نتایج در مدل‌های خطی تعمیم‌یافته است.

### ۲.۲ مدل آمیخته خطی ريج

مدل رگرسیون خطی زیر را در نظر بگیرید

$$y = X\beta + \epsilon \quad (1.2)$$

که در آن  $\mathbf{y} = (y_1, \dots, y_n)^T$  یک بردار  $n$  مؤلفه‌ای شامل متغیر پاسخ،  $y_i$  نشان دهنده  $i$ امین پاسخ مشاهده شده،  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$  ماتریس غیر تصادفی در اندازه  $n \times p$  از متغیرهای پیش‌بین که در آن نشان دهنده  $i$ امین مشاهده در  $i$ امین ستون  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T \in \mathbb{R}^p$  و  $x_{ij}$  نشان دهنده  $i$ امین مشاهده در  $j$ امین ستون متغیر پیش‌بین است. فرض می‌کنیم ماتریس  $\mathbf{X}$  دارای رتبه کامل ستونی است تا ماتریس  $\mathbf{X}^T \mathbf{X}$  که بعد آن  $p$  بوده، معکوس پذیر باشد. لازم به ذکر است در این حالت  $\mathbf{X}^T \mathbf{X}$  دارای  $p$  ستون مستقل خطی است و لذا معکوس پذیر است. بردار ضرایب رگرسیونی  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  پارامترهای مدل و  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$  بردار  $n$  مؤلفه‌ای خطای تصادفی است. در ابتدا فرض می‌کنیم

$$E(\boldsymbol{\epsilon}) = \mathbf{0}, \quad E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T) = \sigma^2 \mathbf{I}$$

که  $\sigma^2 \in \mathbb{R}^+$  مجهول است.

یکی از مسائلی که می‌تواند استفاده از برآوردگر کمترین توان‌های دوم در مدل رگرسیونی (۱.۲) را با مشکل مواجه کند، همبستگی چندگانه بین متغیرهای رگرسیونی است. چراکه این امر سبب پیدایش مشکل هم‌خطی در داده‌ها می‌شود.

اگر یکی از متغیرهای پیش‌بین یک تابع دقیق خطی از یک یا چند متغیر پیش‌بین دیگر باشد، گوییم رگرسیون دارای هم‌خطی کامل<sup>۱</sup> است و اگر این رابطه‌ی خطی، به‌طور تقریبی باشد هم‌خطی ناقص<sup>۲</sup> اتفاق می‌افتد. در صورت وجود هم‌خطی کامل دترمینان ماتریس  $\mathbf{X}^T \mathbf{X}$  صفر است و لذا معکوس آن وجود ندارد. اگرچه می‌توان از معکوس تعمیم‌یافته  $\mathbf{X}^T \mathbf{X}$  استفاده کرد. لازم به ذکر است که معکوس تعمیم‌یافته یکتا نیست و لذا کارایی ندارد. از طرفی اگر معکوس تعمیم‌یافته مور-پنروز<sup>۳</sup> که یکتا است نیز استفاده کنیم استنباط در خصوص برآوردگرهای حاصل به راحتی امکان پذیر نخواهد بود. وجود هم‌خطی ناقص باعث می‌شود دترمینان ماتریس  $\mathbf{X}^T \mathbf{X}$  بسیار کوچک شود و در نتیجه  $(\mathbf{X}^T \mathbf{X})^{-1}$  دارای مولفه‌هایی با قدر مطلق بسیار بزرگ است که در این صورت دقت برآورد ضرایب رگرسیونی کم می‌شود و رفتار ناپایداری دارد. جهت برخورد با این مشکل استفاده از رگرسیون ریدج<sup>۴</sup> که توسط هورل و کنارد (۱۹۷۰) معرفی شد، می‌تواند راهکار مناسبی باشد. برای آگاهی بیشتر در خصوص مشکلات هم‌خطی و چگونگی مواجهه با آن به روزبه (۱۳۹۰) مراجعه کنید.

در رگرسیون ریدج، پارامتر جریمه<sup>۵</sup>  $\lambda$  در یک تابع جریمه درجه دوم از  $\boldsymbol{\beta}$  ضرب شده و برآوردگر ریدج با مینیمم کردن مجموع توان‌های دوم جریمه شده به صورت زیر محاسبه می‌شود

$$\hat{\boldsymbol{\beta}}_{Ridge} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}\} \quad (۲.۲)$$

<sup>۱</sup>Exact collinearity

<sup>۲</sup>Partial collinearity

<sup>۳</sup>Moore–Penrose

<sup>۴</sup>Ridge regression

<sup>۵</sup>Penalty

با مشتق‌گیری از مجموع توان‌های دوم جریمه شده نسبت به پارامتر  $\beta$  داریم

$$\begin{aligned} & \frac{d}{d\beta} [(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T\beta] \\ &= \frac{d}{d\beta} [\mathbf{y}^T\mathbf{y} - \mathbf{y}^T\mathbf{X}\beta - (\mathbf{X}\beta)^T\mathbf{y} + (\mathbf{X}\beta)^T(\mathbf{X}\beta) + \lambda\beta^T\beta] \\ &= \frac{d}{d\beta} [\mathbf{y}^T\mathbf{y} - 2\mathbf{y}^T\mathbf{X}\beta + (\beta^T\mathbf{X}^T\mathbf{X}\beta) + \lambda\beta^T\beta] \\ &= -2\mathbf{y}^T\mathbf{X} + 2\beta^T\mathbf{X}^T\mathbf{X} + 2\lambda\beta^T, \end{aligned}$$

که با مساوی صفر قرار دادن رابطه فوق، برآوردگر ریج به صورت زیر حاصل می‌شود

$$\hat{\beta}_{Ridge} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1} \cdot \mathbf{X}^T\mathbf{y}$$

امید ریاضی برآوردگر ریج به صورت زیر است

$$\begin{aligned} E(\hat{\beta}_{Ridge}) &= E((\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^T\mathbf{y}) \\ &= (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}^T\mathbf{X}\beta. \end{aligned} \quad (3.2)$$

لذا  $\hat{\beta}_{Ridge}$  برآوردگری اریب برای  $\beta$  است. به همین ترتیب نتیجه می‌شود

$$\text{Var}(\hat{\beta}_{Ridge}) = \sigma^2 (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1} \mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p)^{-1}.$$

همچنین میانگین توان‌های دوم خطای برآوردگر ریج عبارت است از

$$\begin{aligned} \text{MSE}(\hat{\beta}_{Ridge}) &= E\left(\left(\hat{\beta}_{Ridge} - \beta\right)^T \left(\hat{\beta}_{Ridge} - \beta\right)\right) \\ &= \text{tr}\left(\text{Var}\left(\hat{\beta}_{Ridge}\right)\right) + \text{Bias}^T\left(\hat{\beta}_{Ridge}\right) \text{Bias}\left(\hat{\beta}_{Ridge}\right) \\ &= \sigma^2 \text{tr}\left(\left(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p\right)^{-1} \mathbf{X}^T\mathbf{X}\left(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}_p\right)^{-1}\right) + \lambda^2 \beta^T (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-2} \beta. \end{aligned} \quad (4.2)$$

اولین جمله سمت راست (۴.۲) مجموع واریانس‌های مولفه‌های پارامتر  $\hat{\beta}_{Ridge}$  و دومین عبارت، توان دوم اریبی (ترانهاده اریبی در خودش) است. بنابراین اریبی  $\hat{\beta}_{Ridge}$  با افزایش  $\lambda$ ، افزایش و واریانس آن با افزایش  $\lambda$ ، کاهش می‌یابد و برعکس. در رگرسیون ریج، پارامتر  $\lambda$  طوری انتخاب می‌شود که کاهش در واریانس بیشتر از افزایش توان دوم اریبی باشد. مقدار بهینه‌ی  $\lambda$  به روش اعتبارسنجی متقابل<sup>۶</sup> توسط کراون و واهبا (۱۹۷۹) از رابطه زیر به دست می‌آید

$$\hat{\lambda} = \arg \min_{\lambda} \left\{ n^{-1}(\mathbf{y} - \hat{\mathbf{y}})^T(\mathbf{y} - \hat{\mathbf{y}}) / (1 - \text{tr}(\mathbf{S})/n) \right\}^2, \quad (5.2)$$

<sup>۶</sup>Cross validation



که در آن  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{Ridge}$  و  $\mathbf{S} = \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T$  همچنین جنتل و همکاران (۲۰۰۴) برای مقادیر کوچک  $\text{tr}(\mathbf{S})/n$  محاسبه مقدار تقریبی  $\lambda$  ی بهینه را از رابطه زیر پیشنهاد کردند

$$\hat{\lambda} \approx \arg \min_{\lambda} \left\{ n^{-1}(\mathbf{y} - \hat{\mathbf{y}})^T(\mathbf{y} - \hat{\mathbf{y}}) + 2n^{-2} \text{tr}(\mathbf{S})(\mathbf{y} - \hat{\mathbf{y}})^T(\mathbf{y} - \hat{\mathbf{y}}) \right\}.$$

لازم به ذکر است که برآوردهای  $\lambda$  محدود به دو حالت فوق نشده و می‌توان انواع دیگر آن را در روزه و همکاران (۲۰۱۵) دید. از طرفی توجه به این نکته لازم است که برآورد ریج یک تابع غیر خطی نسبت به  $\lambda$  است و این خود موجب می‌شود که نتوان به راحتی آن را برآورد کرد. روش رگرسیون ریج در برخورد با داده‌های دارای مشکل همخطی در مطالعات طولی نیز می‌تواند کارساز باشد. الیوت و همکاران (۲۰۱۱) روش رگرسیون آمیخته ریج<sup>۷</sup> (MR) را برای مرتفع کردن این مشکل در داده‌های طولی مطرح کردند. مدل آمیخته خطی زیر را در نظر بگیرید

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i^T\mathbf{b}_i + \varepsilon_i, \quad (۶.۲)$$

که در آن برداری از  $n_i$  مشاهده،  $\mathbf{X}_i$  ماتریس طرح  $n_i \times p$  از اثرات ثابت و  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  بردار  $p \times 1$  ضرایب ثابت و نامعلوم، و  $\mathbf{Z}_i$  ماتریس طرح  $q \times n_i$  از اثرات تصادفی،  $\mathbf{b}_i \sim \mathcal{N}_q(\mathbf{0}, \mathbf{D})$  بردار  $q \times 1$  ضرایب تصادفی و  $\varepsilon_i \sim \mathcal{N}_n(\mathbf{0}, \sigma^2\mathbf{I})$  بردار  $n \times 1$  خطاهای تصادفی است. برای اعمال رگرسیون آمیخته ریج، تابع درستنمایی جریمه شده را می‌توان به صورت زیر نوشت

$$f(\boldsymbol{\beta}, \mathbf{V}; \mathbf{y}) = \frac{|\mathbf{V}|^{-\frac{1}{2}}}{(2\pi)^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \lambda\boldsymbol{\beta}^T\boldsymbol{\beta} \right\}.$$

در نتیجه لگاریتم تابع درستنمایی عبارتست از

$$\ell^{pen}(\boldsymbol{\beta}, \mathbf{V}; \mathbf{y}) = -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \log |\mathbf{V}| - \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \lambda\boldsymbol{\beta}^T\boldsymbol{\beta}.$$

برای دستیابی به برآوردگر ماکزیمم درستنمایی عبارت وابسته به  $\boldsymbol{\beta}$  را در تابع  $\ell^{pen}(\boldsymbol{\beta}, \mathbf{V}; \mathbf{y})$  می‌نیم می‌کنیم

$$\hat{\boldsymbol{\beta}}_{MR} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda\boldsymbol{\beta}^T\boldsymbol{\beta}\}. \quad (۷.۲)$$

با مساوی صفر قرار دادن مشتق عبارت داخل آکولاد در رابطه‌ی (۷.۲) بر حسب  $\boldsymbol{\beta}$  خواهیم

<sup>۷</sup>Mixed ridge

$$\begin{aligned}
 \circ &= \frac{d}{d\beta} \left[ (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta \right] \\
 &= \frac{d}{d\beta} \left[ \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} - \mathbf{y}^T \mathbf{V}^{-1} \mathbf{X}\beta - (\mathbf{X}\beta)^T \mathbf{V}^{-1} \mathbf{y} + (\mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{X}\beta) + \lambda \beta^T \beta \right] \\
 &= \frac{d}{d\beta} \left[ \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + (\mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{X}\beta) + \lambda \beta^T \beta \right] \\
 &= -2\mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + 2\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}\beta + 2\lambda\beta.
 \end{aligned} \tag{۸.۲}$$

در نتیجه برآوردگر MR به شرط معلوم بودن  $\mathbf{V}$  به صورت

$$\hat{\beta}_{MR} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} \tag{۹.۲}$$

حاصل می‌شود.

امیدریاضی برآوردگر MR عبارت است از

$$\begin{aligned}
 E(\hat{\beta}_{MR}) &= E\left((\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}\right) \\
 &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}\beta
 \end{aligned}$$

همچنین واریانس آن برابر

$$\text{Var}(\hat{\beta}_{MR}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I})^{-1}$$

است. مقدار بهینه  $\lambda$  در برآوردگر MR نیز از معادله (۵.۲) با تعریف

$$\mathbf{S} = \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{V}^{-1} + \mathbf{ZDZ}^T \left( \mathbf{I} - \mathbf{X} \left[ (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \right] \right)$$

محاسبه می‌شود.

همان‌طور که در فصل اول گفته شد محاسبه مقدار  $\mathbf{V}$  در رابطه‌ی (۹.۲) مستلزم برآورد پارامتر  $\theta = (\sigma, \mathbf{D})$  است. بدین منظور الیوت و همکاران (۲۰۱۱) الگوریتم مطرح شده توسط لیرد و ویر (۱۹۸۲) را برای  $\hat{\beta}_{MR}$  تعمیم دادند که در ادامه بیان خواهیم کرد.

**الگوریتم ۳ الگوریتم EM در برآورد پارامتر مدل آمیخته خطی ریبج**

**گام ۱)** قرار می‌دهیم  $\hat{\theta}^{(t)} = \theta$  و  $\hat{\lambda}^{(t)} = \lambda$ . سپس مقدار آماره‌های  $\hat{t}_\lambda^{(t)}$  و  $\hat{t}_\theta^{(t)}$  را برای مرحله  $t$ ام به صورت زیر محاسبه می‌کنیم:

$$\begin{aligned} \hat{t}_\lambda^{(t)} &= E \left( \sum_{i=1}^n \epsilon_i^T \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \epsilon_i^T(\hat{\theta}^{(t)}) \epsilon_i(\hat{\theta}^{(t)}) + \text{tr} \text{Var} \left( \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \right], \end{aligned} \quad (10.2)$$

$$\begin{aligned} \hat{t}_\theta^{(t)} &= E \left( \sum_{i=1}^n \mathbf{b}_i \mathbf{b}_i^T \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \mathbf{b}_i(\hat{\theta}^{(t)}) \mathbf{b}_i^T(\hat{\theta}^{(t)}) + \text{Var} \left( \mathbf{b}_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \right]. \end{aligned} \quad (11.2)$$

**گام ۲)** برای به دست آوردن  $\hat{\theta}^{(t+1)}$  مقدار  $\hat{\sigma}^2(t+1)$  و  $\hat{\mathbf{D}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم:

$$\hat{\sigma}^2(t+1) = \hat{t}_\lambda^{(t)} / N,$$

$$\hat{\mathbf{D}}^{(t+1)} = \hat{t}_\theta^{(t)} / n.$$

که در آن  $N = \sum_{i=1}^n n_i$  و  $n$  برابر تعداد واحدهای آزمایشی است. حال برای محاسبه  $\hat{\mathbf{V}}^{(t+1)}$  داریم

$$\hat{\mathbf{V}}^{(t+1)} = \mathbf{Z} \hat{\mathbf{D}}^{(t+1)} \mathbf{Z}^T + \hat{\sigma}^2(t+1) \mathbf{I}. \quad (12.2)$$

**گام ۳)** ماتریس  $\hat{\mathbf{S}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم

$$\begin{aligned} \hat{\mathbf{S}}^{(t+1)} &= \mathbf{X}(\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{\lambda}^{(t)} \mathbf{I})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \\ &+ \mathbf{Z} \hat{\mathbf{D}}^{(t+1)} \mathbf{Z}^T \left( \mathbf{I} - \mathbf{X} \left[ (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{\lambda}^{(t)} \mathbf{I})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \right] \right), \end{aligned} \quad (13.2)$$

سپس برای محاسبه  $\hat{\lambda}^{(t+1)}$  داریم

$$\hat{\lambda}^{(t+1)} = \arg \min_{\lambda} \left\{ n^{-1} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) / (1 - \text{tr}(\hat{\mathbf{S}}^{(t+1)})/n) \right\} \quad (14.2)$$

**گام ۴)** در نهایت مقدار  $\hat{\beta}_{MR}^{(t+1)}$  از رابطه زیر بدست می‌آید

$$\hat{\beta}_{MR}^{(t+1)} = (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{\lambda}^{(t+1)})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{y}. \quad (15.2)$$

**گام ۵)** گام‌های (۱) تا (۴) را تا رسیدن به همگرایی  $\left( |\hat{\beta}_{MR}^{(t+1)} - \hat{\beta}_{MR}^{(t)}| < \gamma \right)$  به ازای یک مقدار  $\gamma$  کوچک و از پیش تعیین شده، تکرار می‌کنیم.

در ادامه چگونگی محاسبه  $\hat{\beta}_{MR}$  را در قالب یک مطالعه شبیه‌سازی و مثال واقعی با استفاده از نرم‌افزار R بررسی می‌کنیم.

## ۱.۲.۲ مطالعه شبیه‌سازی

در این بخش به بررسی عملکرد برآوردگر آمیخته ریج و مقایسه رفتار آن با برآوردگر آمیخته خطی در مدل آمیخته خطی خواهیم پرداخت. از آنجایی که روش ریج یکی از ابزارهای مقابله با همخطی چندگانه است، لذا مطالعه شبیه‌سازی معطوف به این امر است. برای داشتن داده‌هایی با میزان همخطی دلخواه، مک دونالد و همکاران (۱۹۷۵) و گیونز (۱۹۸۱) پیشنهاد دادند که داده‌ها براساس رابطه‌ی

$$x_{ij} = (1 - \rho^2)^{\frac{1}{2}} \cdot Z_{ij} + \rho Z_{ip} \quad i = 1, 2, \dots, n_i, j = 1, 2, \dots, p$$

تولید شوند که در آن اعداد شبه تصادفی مستقل از توزیع نرمال استاندارد و  $\rho^2$  همبستگی بین هر دو متغیر توضیحی است. در ادامه ابتدا ویژگی داده‌های شبیه‌سازی شده را بیان و سپس نتایج را ارائه می‌دهیم.

در این مطالعه برای شبیه‌سازی مقادیر متغیر پاسخ در رابطه‌ی (۶.۲) فرض می‌کنیم

$$\beta = (0, 0/4, 1/0, 1/6, 2/0)$$

همچنین  $\varepsilon_{ijk} \sim \mathcal{N}(0, 1)$  و  $b_{ijk} \sim \mathcal{N}(0, 0/6)$  به طوری که برای هر واحد آزمایشی  $i$  ام  $n_i = 4$  و همچنین هر متغیر پیش‌بین از توزیع نرمال با میانگین ۵ و واریانس ۱ در نظر گرفته شده است. برای ایجاد همبستگی بین متغیرهای پیش‌بین در چند حالت مختلف، مقدار  $\rho$  از ۰ تا ۰/۹ فرض شده است. همچنین برای محاسبه‌ی خطای استاندارد برآوردگر روند شبیه‌سازی مونت کارلویی به تعداد  $B = 200$  تکرار شده است. لازم به ذکر است که در این مطالعه تعداد نمونه‌ی آزمایشی  $n = 40$  در نظر گرفته شده است. نتایج شبیه‌سازی در جدول ۱.۲ آمده است.

جدول ۱.۲: نتایج مطالعه شبیه‌سازی مدل آمیخته خطی ریبج

$\rho$	$\beta$	مدل آمیخته خطی		مدل آمیخته خطی ریبج	
		برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
۰/۰	۰/۰	۰/۰۰۴۰۱۸۶۳	۰/۰۳۴۴۱۳۴۲	۰/۰۰۴۷۱۶۳۷	۰/۰۳۰۷۸۲۸۸
	۰/۴	۰/۴۰۸۸۷۲۳۱	۰/۰۳۰۹۵۰۹۹	۰/۴۲۸۲۲۴۶۲	۰/۰۲۹۷۴۹۷۷
	۱/۰	۱/۰۱۷۶۴۹۹۵	۰/۰۲۹۳۳۴۶۹	۱/۰۱۹۱۸۷۶۶	۰/۰۲۷۸۱۳۴۲
	۱/۶	۱/۶۱۸۰۹۶۸۱	۰/۰۳۳۴۰۱۲۴	۱/۶۳۲۱۳۹۲۱	۰/۰۳۲۷۵۳۹۲
	۲/۰	۱/۹۶۲۷۴۴۶۸	۰/۰۳۵۹۴۵۲۲	۱/۹۴۴۲۸۷۳۰	۰/۰۳۳۹۲۸۱۷
MSE		۰/۸۶۷۷۵۸		۰/۷۷۸۱۵	
۰/۳	۰/۰	-۰/۰۰۲۱۳۷۳۸	۰/۰۴۳۸۶۴۱۷	۰/۰۰۵۵۲۴۲۹	۰/۰۳۸۶۳۲۳۲
	۰/۴	۰/۴۰۹۰۴۹۹۵	۰/۰۳۶۰۳۶۹۱	۰/۳۹۱۰۳۴۴۶	۰/۰۳۵۲۱۹۰۰
	۱/۰	۰/۹۹۹۸۰۲۰۱	۰/۰۳۳۲۳۹۸۱	۰/۹۸۵۹۰۳۸۹	۰/۰۲۹۹۸۸۶۸
	۱/۶	۱/۶۱۴۲۷۰۳۲	۰/۰۴۰۷۷۵۳۵	۱/۶۶۸۸۲۷۱۴	۰/۰۳۹۸۸۶۱۱
	۲/۰	۱/۹۸۹۴۲۵۷۱	۰/۰۴۲۹۸۱۱۴	۱/۹۳۴۷۷۱۳۷	۰/۰۳۹۸۹۱۲۳
MSE		۱/۲۵۴۴۲۱		۱/۰۹۹۶۰	
۰/۶	۰/۰	۰/۰۰۴۱۶۶۲۰	۰/۰۵۴۳۴۴۱۰	۰/۴۱۰۸۰۶۵۳	۰/۰۵۰۹۷۳۹۳
	۰/۴	۰/۴۰۴۱۶۶۲۰	۰/۰۵۲۲۱۳۵۵	۰/۴۱۳۲۱۲۸	۰/۰۴۸۲۳۰۷۷
	۱/۰	۰/۹۸۶۴۴۱۴۲	۰/۰۴۲۴۹۵۰۵	۱/۰۱۷۸۴۲۸۲	۰/۰۳۷۵۸۸۴۸
	۱/۶	۱/۶۰۱۸۴۵۹۴	۰/۰۴۷۰۶۸۸۱	۱/۶۷۸۰۲۹۹۸	۰/۰۴۶۵۶۰۵۰
	۲/۰	۲/۰۷۰۲۹۹۸۲	۰/۰۵۲۲۲۵۸۷	۲/۱۹۵۵۰۶۱۶	۰/۰۵۰۳۷۹۷۶
MSE		۱/۹۹۳۷۰۵		۱/۳۰۳۸۶۵	
۰/۹	۰/۰	۰/۰۲۷۶۱۴۱۱	۰/۰۹۹۳۷۳۴۷	۰/۰۴۷۴۴۹۶۱	۰/۰۸۵۰۸۶۳۳
	۰/۴	۰/۴۰۵۴۹۸۱۰	۰/۰۸۴۸۷۲۵۲	۰/۴۳۰۵۷۳۶۲	۰/۰۷۰۲۲۱۶۲
	۱/۰	۱/۰۵۱۷۲۱۷۰	۰/۰۷۸۲۲۹۶۰	۱/۱۴۶۳۲۹۹۷	۰/۰۶۰۹۸۹۳۶
	۱/۶	۱/۶۰۱۳۶۷۸۷	۰/۰۸۶۸۶۹۳۴	۱/۵۳۲۶۷۱۳۸	۰/۰۷۸۴۸۶۳۶
	۲/۰	۲/۰۶۵۲۷۳۱۳	۰/۰۷۷۱۵۰۲۳	۲/۱۷۱۸۱۹۸۷	۰/۰۷۰۶۸۱۶۶
MSE		۵/۸۷۹۲۰۹		۴/۳۸۶۰۸۹	

همان طور که در جدول ۱.۲ ملاحظه می‌شود اگرچه برآوردهای حاصل از روش آمیخته ریج دارای اریبی  $(|\hat{\beta} - \beta|)$  بیشتری نسبت به برآوردهای حاصل از آمیخته خطی است اما به مراتب دارای انحراف استاندارد ( $sd$ ) کمتری است. این امر را مقادیر میانگین توان‌های دوم خطای برآورد (MSE) نیز تایید می‌کنند. آنچه که در تحلیل نتایج خروجی این شبیه‌سازی حائز اهمیت است این است که هر چقدر مقدار  $\rho$  به سمت یک نزدیک‌تر باشد یعنی میزان همخطی قوی‌تری در داده حاکم باشد برآورد حاصل از مدل آمیخته خطی ریج به‌طور قابل توجهی دارای انحراف استاندارد کمتری نسبت به برآورد حاصل از مدل آمیخته خطی دارا است. همین امر نتیجه استفاده از رگرسیون ریج در مقابله با مشکل همخطی است. اگرچه که این میزان همبستگی در حالت  $\rho = 0$  نیز سبب برتری برآوردگر حاصل از مدل آمیخته ریج نسبت به برآورد بدست آمده از روش مدل آمیخته خطی شده است.

در ادامه، یک مثال واقعی به‌کار گرفته شده توسط الیوت و همکاران (۲۰۱۱) را شرح می‌دهیم.

## ۲.۲.۲ مثال واقعی

اندوتوکسین<sup>۸</sup> نوعی سم موجود در باکتری است که توسط باکتری تولید یا آزاد می‌شود و ارتباط مستقیم با سیستم ایمنی بدن دارد، کمبود این ماده در خون عملکرد سلول‌ها را مختل و همچنین باعث انعقاد خون، اختلال بدن در دریافت اکسیژن، افسردگی، تب و بسیاری از اثرات نامطلوب بر بدن می‌شود. از طرفی وجود بیش از حد مجاز آن نیز خطر آفرین است که از جمله عوارض آن می‌توان به شوک بدن به همراه خونریزی‌های داخلی اشاره کرد. از این رو مطالعه بر روی واکنش بدن در برابر میزان اندوتوکسین از اهمیت ویژه‌ای برخوردار است.

مجموعه داده مورد بررسی در این تحقیق به مطالعه GENE<sup>۹</sup> معروف است که حاصل یک آزمایش مداوم برای تشخیص اثرات عوامل ژنتیکی به میزان نیاسین<sup>۱۰</sup> و اندوتوکسین می‌باشد. در این مطالعه ۱۸۹ نفر به جهت تزریق انسداد اندوتوکسین با دُز کم (LPS)، مورد آزمایش قرار گرفته‌اند. نتیجه این تزریق یک التهاب خفیف پس از حدود ۶ الی ۸ ساعت است که طی آن آپولیپوپروتئین A1<sup>۱۱</sup> (Apo-A1)، آپولیپوپروتئین B (Apo-B)، کلسترول کل<sup>۱۲</sup> (Chol)، کلسترول لیپوپروتئین با چگالی بالا<sup>۱۳</sup> (HDL)، کلسترول لیپوپروتئین با چگالی کم<sup>۱۴</sup> (LDL)، فسفولیپید<sup>۱۵</sup>

<sup>8</sup>Endotoxin

<sup>9</sup>Genetics of Evoked-Responses to Niacin and Endotoxemia

<sup>10</sup>Niacin

<sup>11</sup>Apolipoprotein

<sup>12</sup>Total cholesterol

<sup>13</sup>High-density lipoprotein cholesterol

<sup>14</sup>Low-density lipoprotein cholesterol

<sup>15</sup>Phospholipids

Phos، تری گلیسیرید<sup>۱۶</sup> Trig و فاکتور نکروز تومور آلفا<sup>۱۷</sup> TNF در ۶، ۱۲ و ۲۴ ساعت پس از التهاب ناشی شده، اندازه‌گیری شده‌اند. لازم به ذکر است که در این مطالعه متغیر LPS به عنوان متغیر پاسخ و متغیر جنسیت به عنوان اثر تصادفی در مدل در نظر گرفته شده است.

توجه به یک نکته اساسی مهم است. همبستگی با همخطی تفاوت دارد و باید وجود همخطی با استفاده از ابزارهای مختلف از جمله محاسبه عدد شرطی  $X^T X$  و عامل تورم واریانس<sup>۱۸</sup>، بررسی شود. اما همان‌طور که مونت گومری و همکاران (۲۰۱۲) به آن اشاره کردند، می‌توان از روی مقادیر همبستگی بالا نیز تا حدی وجود همخطی را تشخیص داد. برای آگاهی بیشتر در خصوص معیارهای تشخیص همخطی پایان‌نامه آرست (۱۳۹۵) را ببینید.

همبستگی بین تمام جفت متغیرهای پیش‌بین در جدول ۲.۲، استفاده از مدل آمیخته خطی ریبج را به مدل آمیخته خطی ترجیح می‌دهد. نتایج استفاده از این دو مدل در جدول ۳.۲ خلاصه شده‌است.

جدول ۲.۲: همبستگی بین متغیرهای پیش‌بین در داده‌های GENE

	TNF	Apo-A1	Apo-B	Chol	HDL	Trig	Phos	LDL
TNF	۱/۰۰	۰/۱۳	۰/۰۴	۰/۰۴	۰/۰۷	-۰/۰۵	۰/۰۲	۰/۰۳
Apo-A1	۰/۱۳	۱/۰۰	۰/۰۷	۰/۴۰	۰/۸۷	۰/۰۱	۰/۶۶	۰/۰۱
Apo-B	۰/۰۴	۰/۰۷	۱/۰۰	۰/۸۸	-۰/۱۲	۰/۴۰	۰/۵۶	۰/۹۵
Chol	۰/۰۴	۰/۴۰	۰/۸۸	۱/۰۰	۰/۲۹	۰/۳۱	۰/۷۹	۰/۹۰
HDL	۰/۰۷	۰/۸۷	-۰/۱۲	۰/۲۹	۱/۰۰	-۰/۳۲	۰/۴۶	-۰/۱۱
Trig	-۰/۰۵	۰/۰۱	۰/۴۰	۰/۳۱	-۰/۳۲	۱/۰۰	۰/۴۲	۰/۲۷
Phos	۰/۰۲	۰/۶۶	۰/۵۶	۰/۷۹	۰/۴۶	۰/۴۲	۱/۰۰	۰/۵۵
LDL	۰/۰۳	۰/۰۱	۰/۹۵	۰/۹۰	-۰/۱۱	۰/۲۷	۰/۵۵	۱/۰۰

<sup>16</sup>Triglycerides

<sup>17</sup>Tumor necrosis factor alpha

<sup>18</sup>Variance inflation factor

جدول ۳.۲: برآورد ضرایب در داده‌های GENE

	مدل آمیخته خطی		مدل آمیخته خطی ریح	
	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
Time	۰/۴۶	۰/۱۶	۰/۴۷	۰/۱۲
Time1	-۰/۴۰	۰/۲۶	-۰/۴۰	۰/۲۰
Time2	-۰/۱۷	۰/۱۸	-۰/۱۸	۰/۱۵
Apo-A1	۰/۰۸	۰/۰۵	۰/۰۷	۰/۰۴
Apo-B	-۰/۱۸	۰/۰۹	-۰/۲۴	۰/۰۷
Chol	۰/۴۹	۱/۲۰	۱/۱۶	۰/۹۴
HDL	-۰/۵۵	۱/۲۰	-۱/۱۶	۰/۹۴
LDL	-۰/۳۵	۱/۲۰	-۰/۹۸	۰/۹۴
Phos	-۰/۰۵	۰/۰۳	-۰/۰۶	۰/۰۲
Trig	-۰/۰۸	۰/۲۴	-۰/۲۱	۰/۱۹
TNF	۰/۱۵	۰/۲۰	۰/۱۱	۰/۱۳

همان‌طور که در جدول ۳.۲ ملاحظه می‌شود مقادیر مربوط به انحراف استاندارد برآورد ضرایب در مدل آمیخته خطی ریح به مراتب کمتر از انحراف استاندارد برآورد ضرایب در مدل آمیخته خطی است. لذا می‌توان از این حیث مدل آمیخته خطی ریح را بر مدل آمیخته خطی برتر دانست.

## ۳.۲ مدل آمیخته خطی ریح با خطای بیضی گون

در این بخش مدل آمیخته خطی ریح را با فرض پایه‌ای این که بردار خطای تصادفی دارای توزیع بیضی گون<sup>۱۹</sup> (توزیع چندمتغیره با منحنی‌های تراز بیضی شکل) می‌باشند، مورد بررسی قرار می‌دهیم و نتایج مطالعه شبیه‌سازی را در انتهای بخش در مقایسه با مدل آمیخته خطی ارائه می‌دهیم.

توضیحات مربوط به توزیع بیضی گون عمدتاً برگرفته از آرشی (۱۳۸۷) و حسن‌زاده بشتیان (۱۳۸۸) است و سایر نتایج آن جدید هستند.

<sup>19</sup>Elliptically contoured distribution



## ۱.۳.۲ توزیع بیضی گون

بردار تصادفی  $n$  مولفه‌ای  $\mathbf{X}$  دارای توزیع بیضی گون با پارامترهای  $\mu$ ،  $\Sigma$  و تابع مولد مشخصه<sup>۲۰</sup> است و با  $\mathbf{X} \sim \mathcal{E}_n(\mu, \Sigma, \psi)$  نشان می‌دهیم اگر تابع مشخصه آن به صورت زیر باشد

$$\phi_{\mathbf{X}}(t) = \exp(it^T \mu) \psi \left( \frac{t^T \Sigma t}{2} \right)$$

که در آن  $\psi(\cdot)$  متعلق به رده توابعی به صورت  $\psi(t) : [0, \infty] \rightarrow \mathbb{R}$  است به طوری که  $\psi \left( \sum_{i=1}^n t_i^2 \right)$  یک تابع مشخصه  $n$ -بعدی باشد.

توزیع‌های بیضی گون ( $ECD_s$ )، توزیع‌هایی هستند که منحنی‌های تراز آن‌ها بیضی شکل است. از این که  $\mathbf{X} \sim \mathcal{E}_n(\mu, \Sigma, \psi)$  نمی‌توان نتیجه گرفت که متغیر تصادفی  $\mathbf{X}$  دارای تابع چگالی است، ولی در صورتی که تابع چگالی احتمال وجود داشته باشد، به صورت زیر است:

$$f_{\mathbf{X}}(x) = c_n |\Sigma|^{1/2} g \left[ \frac{(x - \mu)^T \Sigma^{-1} (x - \mu)}{2} \right]$$

که در آن تابع  $g(\cdot)$  تابع مولد چگالی<sup>۲۱</sup> است و دارای شرط زیر می‌باشد

$$\int_0^{\infty} x^{n/2-1} g(x) dx < \infty$$

و همچنین  $c_n$  ضریب نرمال سازی می‌باشد که به صورت زیر تعیین می‌گردد

$$c_n = \frac{\Gamma(n/2)}{(\sqrt{\pi})^{n/2}} \left[ \int_0^{\infty} x^{n/2-1} g(x) dx \right]^{-1}$$

در این صورت می‌نویسیم  $\mathbf{X} \sim \mathcal{E}_n(\mu, \Sigma, g)$ . این نکته حائز اهمیت است، در صورتی که تابع چگالی وجود داشته باشد، تابع  $g$ ، تابع  $\psi$  را مشخص می‌کند و برعکس. برای آگاهی بیشتر در این زمینه فنگ و همکاران (۱۹۹۰)، میرهد (۱۹۸۲) و کوپتا و وارگا (۱۹۹۳) را ببینید.

در این توزیع،  $E(\mathbf{X}) = \mu$  و  $\text{Cov}(\mathbf{X}) = -2\psi^T(0)\Sigma$  وجود ندارد در صورتی که  $|\psi^T(0)| < \infty$ . توزیع‌های زیادی متعلق به خانواده توزیع‌های بیضی گون می‌باشند که از آن جمله می‌توان به توزیع‌های نرمال چندمتغیره،  $t$  چندمتغیره، کوشی چندمتغیره<sup>۲۲</sup>، پیرسن چندمتغیره نوع II<sup>۲۳</sup>، پیرسن چندمتغیره نوع VII<sup>۲۴</sup>، بسل چندمتغیره<sup>۲۵</sup>، نمایی توانی چندمتغیره<sup>۲۶</sup>، لاپلاس چندمتغیره<sup>۲۷</sup>، اسلش تعمیم یافته<sup>۲۸</sup> و کاتز چندمتغیره<sup>۲۹</sup> اشاره کرد.

<sup>20</sup>Characteristic generator

<sup>21</sup>Density generator

<sup>22</sup>Multivariate Cauchy

<sup>23</sup>Multivariate Pearson type II

<sup>24</sup>Multivariate Pearson type VII

<sup>25</sup>Multivariate Bessel

<sup>26</sup>Multivariate Exponential power

<sup>27</sup>Multivariate Laplace

<sup>28</sup>Generalized Slash

<sup>29</sup>Multivariate Kotz type

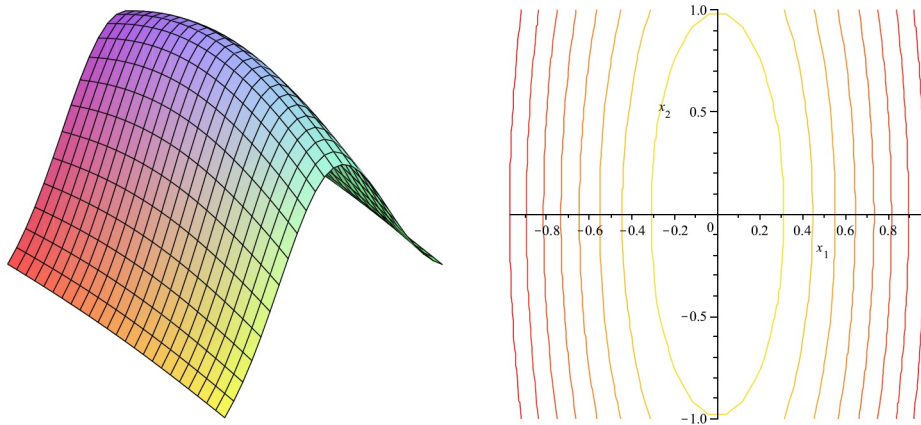
حال به عنوان نمونه، تابع چگالی احتمال و تابع مولد چگالی را برای چند توزیع مهم بیضی گون، همراه با نمودار آن‌ها ارائه می‌دهیم. در نمودارها  $n = 2$  و  $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  فرض شده است.

۱. توزیع نرمال  $n$  متغیره  $(\mathbf{X} \sim \mathcal{N}_n(\mu, \Sigma))$ : توابع مولد چگالی و چگالی احتمال به ترتیب عبارتند از

$$g(u) = e^{-u},$$

$$f_{\mathbf{X}}(x) = \frac{|\Sigma|^{-1/2}}{(\sqrt{2\pi})^n} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}.$$

نمودارهای منحنی‌های تراز و تابع چگالی در حالت  $n = 2$  در شکل ۱.۲ داده شده است.



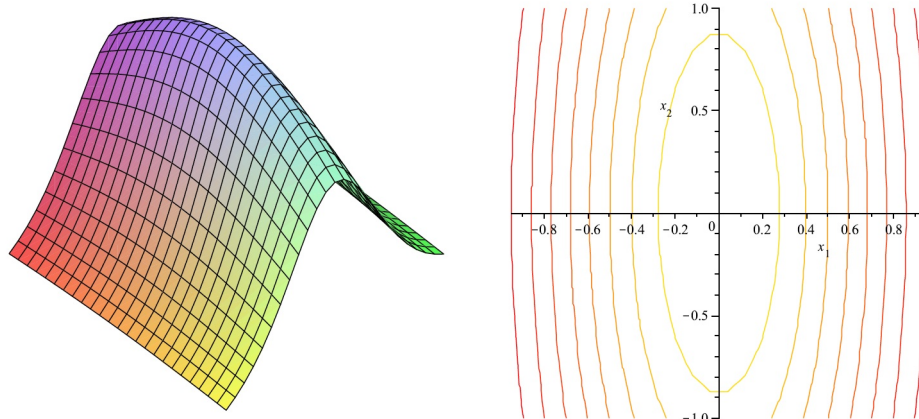
شکل ۱.۲: نمودارهای منحنی‌های تراز و تابع چگالی توزیع نرمال دومتغیره

۲. توزیع  $t$  چندمتغیره  $(\mathbf{X} \sim \mathcal{T}_n(\mu, \Sigma, \mathbf{V}))$ : توابع مولد چگالی و چگالی احتمال به ترتیب عبارتند از

$$g_n(u) = \left( 1 + \frac{2u}{\nu} \right)^{-(n+\nu)/2},$$

$$f_{\mathbf{X}}(x) = \frac{\Gamma\left(\frac{n+\nu}{2}\right) |\Sigma|^{-1/2}}{(\nu\pi)^{n/2} \Gamma\left(\frac{\nu}{2}\right)} \left\{ 1 + \frac{(x - \mu)^T \Sigma^{-1} (x - \mu)}{\nu} \right\}^{-\frac{(n+\nu)}{2}}.$$

نمودارهای منحنی‌های تراز و تابع چگالی به ازای  $n = 2$  و  $\nu = 3$  در شکل ۲.۲ داده شده است.



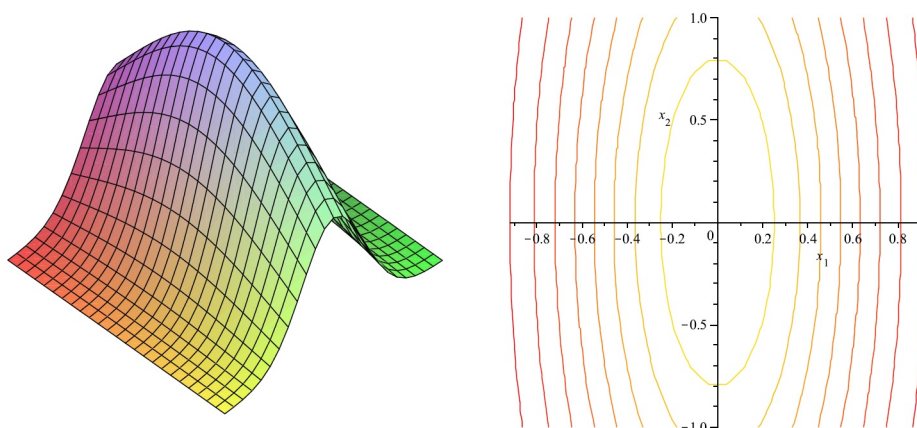
شکل ۲.۲: نمودارهای منحنی‌های تراز و تابع چگالی توزیع  $t$  دومتغیره

۳. توزیع کاتز  $n$  متغیره  $(X \sim \mathcal{K}_n(\mu, \Sigma, r))$ : توابع مولد چگالی و چگالی احتمال به ترتیب عبارتند از

$$g_n(u) = e^{-ru}, r > 0,$$

$$f_{\mathbf{X}}(x) = \frac{\Gamma\left(\frac{n}{\nu}\right) |\Sigma|^{-1/\nu} r^{n/\nu}}{(\nu\pi)^{n/\nu} \Gamma\left(\frac{n}{\nu}\right)} \exp\left\{-\frac{r}{\nu} \left[(x - \mu)^T \Sigma^{-1} (x - \mu)\right]\right\}.$$

نمودارهای منحنی‌های تراز و تابع چگالی به ازای  $r = 3$  در شکل ۳.۲ داده شده است.



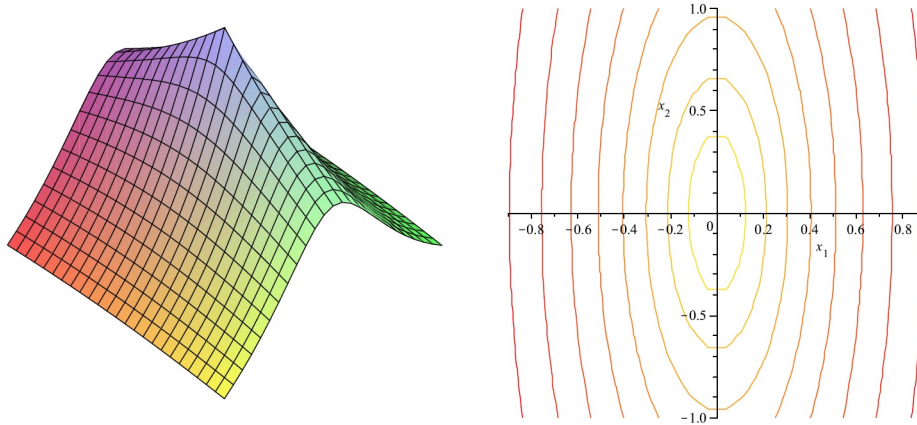
شکل ۳.۲: نمودارهای منحنی‌های تراز و تابع چگالی توزیع کاتز دومتغیره

۴. توزیع لاپلاس چندمتغیره  $(X \sim \mathcal{L}_n(\mu, \Sigma, r))$ : توابع مولد چگالی و چگالی احتمال به

ترتیب عبارتند از

$$g_n(u) = e^{-\sqrt{2}u^{1/2}},$$

$$f_{\mathbf{X}}(x) = \frac{\left(\frac{1}{\sqrt{2}}\right) \Gamma\left(\frac{n}{\sqrt{2}}\right) |\Sigma|^{-1/2} \sqrt{2}^n}{(\sqrt{2}\pi)^{n/\sqrt{2}} \Gamma(n)} \exp\left\{-\frac{\sqrt{2}}{2} \left[(x - \mu)^T \Sigma^{-1} (x - \mu)\right]^{1/2}\right\}.$$



شکل ۴.۲: نمودارهای منحنی‌های تراز و تابع چگالی توزیع لاپلاس دومتغیره

در ادامه چند قضیه مهم را در خصوص توزیع بیضی گون ارائه می‌کنیم.

لم ۱.۳.۲. (چو، ۱۹۷۳). فرض کنید  $\mathbf{X} \sim \mathcal{E}_n(\mu, \Sigma, \psi)$ ، در این صورت تابع چگالی احتمال  $\mathbf{X}$  را می‌توان به صورت زیر نوشت

$$f_{\mathbf{X}}(\mathbf{x}) = \int_0^{+\infty} w(t)h(x)dt$$

که در آن  $h(\cdot)$  نشان‌دهنده تابع چگالی احتمال  $\mathcal{N}_n(\mu, t^{-1}\Sigma)$  است. بنابراین

$$f_{\mathbf{X}}(\mathbf{x}) = \int_0^{+\infty} w(t) \frac{1}{(\sqrt{2}\pi)^{n/\sqrt{2}}} |t^{-1}\Sigma|^{-1/2} \exp\left\{-\frac{1}{\sqrt{2}}(\mathbf{x} - \mu)^T (t^{-1}\Sigma)^{-1} (\mathbf{x} - \mu)\right\} dt$$

که در آن

$$w(t) = (\sqrt{2}\pi)^{n/\sqrt{2}} |t^{-1}\Sigma|^{1/2} L^{-1}[f(s)]$$

و  $L^{-1}[f(s)]$  نشان‌دهنده تبدیل معکوس لاپلاس  $h(s)$  به ازای

$$s = (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)$$

می‌باشد.

قضیه ۱.۳.۲. اگر  $\mathbf{X} \sim \mathcal{E}_p(\mu, \Sigma, \psi)$ ،  $\mathbf{A}_{q \times p}$  و  $\mathbf{C}_{q \times 1}$  ماتریس‌های ثابت باشند، آنگاه

$$\mathbf{A}\mathbf{X} + \mathbf{C} \sim \mathcal{E}_q(\mathbf{A}\mu + \mathbf{C}, \mathbf{A}\Sigma\mathbf{A}^T, \psi).$$

**قضیه ۲.۳.۲.** فرض کنید  $\mathbf{X} \sim \mathcal{E}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \psi)$  و  $\Omega$  مجموعه‌ای در فضای  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  باشد، به طوری که اگر  $(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \in \Omega$  آنگاه  $(\boldsymbol{\mu}, c\boldsymbol{\Sigma}) \in \Omega, \forall c > 0$ . همچنین فرض کنید  $g$  به گونه‌ای باشد که  $g(\mathbf{x}^T \mathbf{x})$  یک تابع چگالی در  $\mathbb{R}^n$  و  $y^{N/2} g(y)$  دارای یک مقدار ماکزیمم مثبت متناهی مانند  $y_g$  باشد و براساس تک مشاهده  $\mathbf{x}$  از  $|\boldsymbol{\Sigma}|^{-1/2} g[(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})]$  برآوردگر درست‌نمایی ماکزیمم  $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \in \Omega$  تحت نرمال بودن، موجود و منحصر بفرد باشد، آنگاه  $MLE$  برای  $g$  عبارت است از

$$\hat{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}} \quad \hat{\boldsymbol{\Sigma}} = \frac{N}{y_g} \tilde{\boldsymbol{\Sigma}} \quad (۱۶.۲)$$

و تابع درست‌نمایی ماکزیمم به صورت زیر است

$$|\hat{\boldsymbol{\Sigma}}|^{-1/2} g(y_g). \quad (۱۷.۲)$$

اثبات: فرض کنید  $\mathbf{B} = |\boldsymbol{\Sigma}|^{-1/N} \boldsymbol{\Sigma}$  و

$$d = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \frac{(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{B}^{-1} (\mathbf{x} - \boldsymbol{\mu})}{|\boldsymbol{\Sigma}|^{1/N}} \quad (۱۸.۲)$$

آنگاه  $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \in \Omega$  و  $|\mathbf{B}| = 1$ .

تابع درست‌نمایی عبارت است از

$$\left[ (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{B}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]^{-N/2} d^{N/2} g(d) \quad (۱۹.۲)$$

تحت نرمال بودن  $g(d) = (2\pi)^{-n/2} e^{-d/2}$ ، ماکزیمم رابطه (۱۹.۲) در  $\boldsymbol{\mu} = \tilde{\boldsymbol{\mu}}$ ،

$$\mathbf{B} = \tilde{\mathbf{B}} = |\tilde{\boldsymbol{\Sigma}}|^{-1/N} \tilde{\boldsymbol{\Sigma}}$$

و  $d = N$  رخ می‌دهد. به طور کلی ماکزیمم رابطه (۱۹.۲) به ازای  $\hat{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}$  و  $\hat{\mathbf{B}} = \tilde{\mathbf{B}}$  حاصل می‌شود. بنابراین  $MLE$  پارامتر  $\boldsymbol{\Sigma}$  به صورت

$$\hat{\boldsymbol{\Sigma}} = |\hat{\boldsymbol{\Sigma}}|^{1/N} \hat{\mathbf{B}} = \frac{|\tilde{\boldsymbol{\Sigma}}|^{1/N}}{|\hat{\boldsymbol{\Sigma}}|^{1/N}} \tilde{\boldsymbol{\Sigma}}$$

است. با جایگذاری مقادیر  $|\hat{\boldsymbol{\Sigma}}|^{1/N}$  و  $|\tilde{\boldsymbol{\Sigma}}|^{1/N}$  از رابطه (۱۸.۲) به (۱۶.۲) می‌رسیم. همچنین با جایگذاری مقادیر فوق در رابطه (۱۹.۲)، رابطه (۱۷.۲) نتیجه می‌شود.

## ۲.۳.۲ مدل MR بیضی گون

مدل آمیخته خطی را در نظر بگیرید

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i^T \mathbf{b}_i + \boldsymbol{\epsilon}_i \quad (۲۰.۲)$$

که در آن برداری از  $n_i$  مشاهده،  $\mathbf{X}_i$  ماتریس طرح  $n_i \times p$  از اثرات ثابت و  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$  بردار  $p \times 1$  ضرایب ثابت و نامعلوم، و  $\mathbf{Z}_i$  ماتریس طرح  $q \times n_i$  از

اثرات تصادفی است. همچنین  $b_i$  بردار  $q \times 1$  ضرایب تصادفی و  $\epsilon_i$  بردار  $n \times 1$  خطاهای تصادفی هستند که دارای توزیع بیضی گون به صورت زیر می باشند

$$\begin{aligned} b_i &\sim \mathcal{E}_q(\circ, \mathbf{D}, g) \\ \epsilon_i &\sim \mathcal{E}_n(\circ, \sigma^2 \mathbf{I}, g), \end{aligned}$$

و تابع مشخصه  $\phi_{\epsilon}(t)$  به صورت

$$\phi_{\epsilon}(t) = \psi(\sigma^2 t^T t)$$

است.

اگر  $\epsilon_i$  دارای یک تابع چگالی احتمال باشد آن گاه طبق آنچه گفته شد به صورت

$$f(\epsilon_i) \propto |\sigma^2 \mathbf{I}|^{-\frac{1}{2}} g\left(\frac{1}{\sigma^2} \epsilon_i^T \epsilon_i\right)$$

می باشد که در آن  $g(\cdot)$  یک تابع نامنفی روی  $\mathbb{R}^+$  بوده به طوری که  $f(\cdot)$  یک تابع چگالی است. با توجه به قضیه ۱.۳.۲، هر ترکیب خطی از توزیع بیضی گون، دارای توزیع بیضی گون بوده و در این صورت داریم

$$y_i \sim \mathcal{E}_n(\mathbf{X}_i \beta, \mathbf{V}, g)$$

بنابراین لگاریتم تابع درستنمایی متغیر  $y$  به شرط معلوم بودن تابع  $g(\cdot)$  به صورت زیر است

$$\ell(\beta, \mathbf{V}, \mathbf{y}) = -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \log |\mathbf{V}| + \log g\left(-\frac{1}{2} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta)\right).$$

ماکزیمم کردن تابع فوق بر حسب  $\beta$ ، معادل با می نیمم کردن عبارت وابسته به  $\beta$  است. مشابه قضیه ۲.۳.۲ و به طور دقیق تر طبق آرسی و همکاران (۲۰۱۳) برآوردگر ماکزیمم درستنمایی ( $MLE$ ) پارامتر  $\beta$  به صورت

$$\hat{\beta} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}. \quad (21.2)$$

است. همانطور که مشاهده می شود، تابع  $g$  که نوع توزیع بیضی گون را مشخص می کند در ساختار برآوردگر  $\beta$  هیچ تأثیری ندارد. لذا اولین سوالی که در ذهن خواننده پیش می آید این است که دلیل استفاده از توزیع بیضی گون به عنوان مدل مؤلفه های خطا و اثر تصادفی چیست؟ در پاسخ می توان گفت گاهی اوقات تعیین توزیع دقیق متغیر پاسخ در مثال های واقعی کاری بس دشوار است در حالی که می توان متقارن بودن آن را به راحتی بر اساس داده ها بررسی نمود. در این حالت چنانچه بر اساس ضرایب چولگی پی به متقارن بودن توزیع داده ها ببریم، دیگر لازم نیست در تحلیل داده های طولی نگران توزیع دقیق مؤلفه پاسخ برای یافتن  $MLE$  باشیم و می توان از رابطه (۲۱.۲) به عنوان برآوردگر صحیح بردار ضرایب ثابت  $\beta$  استفاده کرد.

حال لگاریتم تابع درست‌نمایی متغیر  $y$  را بر پایه‌ی مدل آمیخته خطی با اعمال تابع جریمه در رگرسیون ریج، به صورت زیر داریم

$$\ell^{pen}(\beta, \mathbf{V}, \mathbf{y}) = -\frac{N}{\nu} \log(\nu\pi) - \frac{1}{\nu} \sum_{i=1}^n \log |\mathbf{V}| + \log g \left( -\frac{1}{\nu} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) \right) - \lambda \beta^T \beta,$$

بنابراین با انجام محاسباتی ساده مشابه بخش ۲.۲ برآوردگر  $MR$  پارامتر  $\beta$  به صورت

$$\hat{\beta}_{MR} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

حاصل می‌شود.

### ۳.۳.۲ مطالعه شبیه‌سازی

در این قسمت فرض کنید که مؤلفه‌های تصادفی مدل (۲۰.۲) برای شبیه‌سازی مقادیر متغیر پاسخ، بردار ضرایب رگرسیونی به صورت

$$\beta = (0, 0.4, 1, 1/6, 2)$$

،  $b_{ijk} \sim T(0, 0.3, 5)$  و  $\varepsilon_{ijk} \sim T(0, 1, 5)$  فرض شده‌اند به طوری که برای هر واحد آزمایشی  $i$ ام  $n_i = 4$  و هر متغیر پیش‌بین از توزیع نرمال با میانگین ۵ و واریانس ۱ فرض شده‌اند. همچنین برای ایجاد داده‌های هم‌خط از بسته‌ی EnvStat در نرم‌افزار  $R$  با ضریب همبستگی‌های متفاوت موجود در جدول ۴.۲

$$\rho \in \{0/0, 0/2, 0/5, 0/7, 0/9\}$$

استفاده شده است.

لازم به ذکر است که برای محاسبه‌ی خطای استاندارد برآورد بدست آمده با روند شبیه‌سازی مونت کارلویی به تعداد  $B = 100$  تکرار شده و تعداد نمونه آزمایشی  $n = 40$  در نظر گرفته شده است.

جدول ۴.۲: نتایج مطالعه شبیه‌سازی مدل آمیخته خطی ریح با خطای بیضی گون

		مدل آمیخته خطی با خطای بیضی گون		مدل آمیخته خطی ریح با خطای بیضی گون	
$\rho$	$\beta$	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
۰/۰	۰/۰	-۰/۰۱۴۲۵۴۲۳	۰/۱۲۲۶۴۴۶۹	-۰/۰۱۰۴۱۲۴	۰/۱۲۱۷۸۵۱
	۰/۴	۰/۴۰۱۱۱۳۵۵	۰/۱۱۱۹۵۰۳۹	۰/۴۰۳۳۴۷۶	۰/۱۱۱۴۹۲۳
	۱/۰	۰/۹۹۵۱۳۸۹۱	۰/۱۲۷۳۴۴۸۳	۰/۹۹۴۷۲۱۴	۰/۱۲۶۸۱۱۳
	۱/۶	۱/۶۱۷۲۸۱۷۴	۰/۰۹۶۲۹۷۴۵	۱/۶۱۵۱۳۳۰	۰/۰۹۵۶۲۵۲
	۲/۰	۱/۹۹۷۳۸۸۳۲	۰/۱۱۳۳۸۸۱۸	۱/۹۹۳۶۰۳۴	۰/۱۱۳۱۶۶۲
MSE		۱۰/۵۴۷۹۶		۱۰/۴۴۷۴۶	
۰/۳	۰/۰	-۰/۰۱۹۰۹۵۰۳	۰/۱۴۰۴۲۴۲	-۰/۰۱۳۴۹۸۰۳	۰/۱۳۸۹۶۱۷
	۰/۴	۰/۳۹۹۵۶۱۴۳	۰/۱۳۲۲۷۸۶	۰/۴۰۳۰۲۹۱۷	۰/۱۳۱۴۴۵۲
	۱/۰	۰/۹۹۲۵۴۶۲۶	۰/۱۴۷۰۰۱۹	۰/۹۹۱۹۸۶۰۶	۰/۱۴۶۱۳۹۶
	۱/۶	۱/۶۲۱۷۶۵۱۶	۰/۱۱۰۰۰۴۴	۱/۶۱۸۴۹۱۱۱	۰/۱۰۸۸۶۸۷
	۲/۰	۲/۰۰۱۸۸۴۹۵	۰/۱۳۲۷۲۷۹	۱/۹۹۶۳۴۰۸۳	۰/۱۳۲۰۹۹۸
MSE		۱۴/۱۶۷۹۱		۱۳/۹۶۰۲۵	
۰/۶	۰/۰	-۰/۰۲۳۹۹۷۱۷	۰/۱۷۵۱۶۳۶	-۰/۰۱۴۶۹۵۰۱	۰/۱۷۲۳۱۹۱
	۰/۴	۰/۳۹۶۴۳۴۶۹	۰/۱۷۶۹۸۶۱	۰/۴۰۲۲۳۵۷۹	۰/۱۷۵۰۱۴۰
	۱/۰	۰/۹۸۸۱۰۳۵۴	۰/۱۹۲۵۵۴۷	۰/۹۸۷۷۳۴۳۳	۰/۱۹۰۶۳۱۵
	۱/۶	۱/۶۲۹۷۷۰۷۸	۰/۱۴۷۷۸۸۳	۱/۶۲۴۰۷۴۸۳	۰/۱۴۵۴۹۷۲
	۲/۰	۲/۰۰۶۱۱۴۰۷	۰/۱۷۹۳۵۰۷	۱/۹۹۶۶۹۷۴۲	۰/۱۷۷۶۳۲۰
MSE		۲۴/۴۹۶۳۳		۲۳/۹۰۲۸۳	
۰/۹	۰/۰	-۰/۰۴۷۱۲۲۶۱	۰/۳۴۸۴۷۳۸	-۰/۰۱۱۴۲۵۹۵	۰/۳۳۱۹۸۵۴
	۰/۴	۰/۳۹۹۴۳۶۷۱	۰/۳۴۹۵۴۵۶	۰/۴۰۷۴۸۷۳	۰/۳۳۶۴۶۸۰
	۱/۰	۰/۹۸۴۸۸۵۷۴	۰/۳۸۱۱۷۷۲۴	۰/۹۸۲۵۳۹۲۹	۰/۳۹۸۵۳۴۶
	۱/۶	۱/۶۴۳۸۴۶۷۵	۰/۲۸۵۳۲۷۲	۱/۶۲۳۸۵۷۷۲	۰/۲۷۴۱۱۰۴
	۲/۰	۲/۰۱۵۲۳۱۲۶	۰/۳۴۶۰۵۶۷	۱/۹۷۹۵۸۶۰۱	۰/۳۳۴۲۰۷۶
MSE		۹۴/۴۸۹۹		۹۱/۰۵۵۶۲	

نتایج مطالعه شبیه‌سازی در جدول ۴.۲ آمده است. همان‌طور که در این جدول ملاحظه می‌شود برآوردهای حاصل از روش آمیخته ریح با خطاهای بیضی گون دارای انحراف استاندارد کمتری نسبت به برآورد حاصل از روش آمیخته خطی است. مقادیر میانگین توان دوم خطای برآورد (MSE) نیز برتری روش آمیخته ریح را نسبت به روش آمیخته خطی تایید می‌کند. همچنین طبق جدول ۴.۲ واضح است که اگر چه در برخی موارد میزان اریبی  $(|\hat{\beta} - \beta|)$  نسبت به روش آمیخته خطی کمتر است اما هرچه میزان همخطی در داده‌ها قوی‌تر باشد، انحراف استاندارد برآورد در روش آمیخته ریح کاهش بیشتری داشته است.



## ۴.۲ برآوردگرهای انقباضی در مدل آمیخته خطی ريج

در این بخش برآوردگرهای آزمون اولیه<sup>۳۰</sup> و انقباضی نوع استاین<sup>۳۱</sup> را تعريف کرده و به جهت بهبود خطای برآورد در مدل آمیخته خطی ريج، برآوردگرهای آزمون اولیه آمیخته ريج و انقباضی آمیخته و انقباضی آمیخته ريج را برای برآورد ضرایب اثرات ثابت مدل، مطرح می‌کنیم. گاهی در مسائل آماری، شرایط موجود ایجاب می‌کند که در مورد رخداد فرضیه‌ی  $H_0: \beta = \beta_0$  نمی‌توان به‌طور حتمی نظر داد و در مورد اینکه آیا  $\beta = \beta_0$  درست است یا خیر تردید ایجاد می‌شود، یعنی نمی‌توان به‌طور قطعی یکی از دو برآوردگر  $\beta$  یا  $\beta_0$  را انتخاب کرد. در چنین مواردی می‌توان ابتدا فرضیه‌ی  $H_0$  را آزمون کرد و سپس براساس نتیجه‌ی آزمون یکی از دو برآوردگر را انتخاب نمود. بر این اساس بنکرافت (۱۹۴۴) برآوردگر آزمون اولیه  $\beta$  را به صورت زیر ارائه کرد

$$\hat{\beta}^{PT} = \hat{\beta} \mathbf{I}_{\{H_0 \text{ رد}\}} + \beta_0 \mathbf{I}_{\{H_0 \text{ پذیرش}\}}$$

که در آن  $\mathbf{I}_{\{A\}}$  تابع نشان‌گر مجموعه  $A$  است.

در مدل آمیخته خطی ريج فرضیه‌ی زیر را در نظر بگیرید

$$\begin{cases} H_0: \beta = \beta_0 \\ H_1: \beta \neq \beta_0 \end{cases}$$

در این صورت برآوردگر آزمون اولیه آمیخته ريج<sup>۳۲</sup> ( $PTMR$ ) را می‌توان به صورت زیر بیان کرد

$$\hat{\beta}_{MR}^{PT} = \begin{cases} \hat{\beta}_{MR} & \text{فرضیه‌ی } H_0 \text{ رد شود} \\ \beta_0 & \text{فرضیه‌ی } H_0 \text{ پذیرفته شود} \end{cases}$$

برآوردگر  $PTMR$  را می‌توان با استفاده از تابع نشانگر به صورت زیر بازنویسی کرد

$$\begin{aligned} \hat{\beta}_{MR}^{PT} &= \hat{\beta}_{MR} \mathbf{I}_{\{H_0 \text{ رد}\}} + \beta_0 \mathbf{I}_{\{H_0 \text{ پذیرش}\}} \\ &= \hat{\beta}_{MR} - \hat{\beta}_{MR} \mathbf{I}_{\{H_0 \text{ پذیرش}\}} + \beta_0 \mathbf{I}_{\{H_0 \text{ پذیرش}\}} \\ &= \hat{\beta}_{MR} - (\hat{\beta}_{MR} - \beta_0) \mathbf{I}_{\{H_0 \text{ پذیرش}\}} \\ &= \hat{\beta}_{MR} - (\hat{\beta}_{MR} - \beta_0) \mathbf{I}_{(L_n < L_n(\alpha))} \end{aligned}$$

که در آن  $L_n$  آماره والد<sup>۳۳</sup> آزمون فرضیه‌ی  $H_0$ ، با توجه به نرمال بودن توزیع مجانبی  $\hat{\beta}$  به صورت زیر است (آرشی و همکاران ۲۰۱۳)

<sup>30</sup>Preliminary test

<sup>31</sup>Stein-type shrinkage

<sup>32</sup>Preliminary test mixed ridge

<sup>33</sup>Wald statistic

$$\mathbf{L}_n = (\hat{\beta} - \beta_0)^T (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}) (\hat{\beta} - \beta_0) \quad (22.2)$$

همچنین  $\mathbf{L}_n(\alpha)$  صدک بالایی  $100\alpha$  درصد توزیع کی-دو با  $p$  درجه آزادی است رفتار برآوردگر  $PTMR$  به اندازه آزمون یا سطح معناداری  $\alpha$  ( $0 < \alpha < 1$ ) بستگی دارد. با توجه به نتیجه‌ی آزمون،  $\hat{\beta}_{MR}^{PT}$  به یکی از دو مقدار  $\hat{\beta}_{MR}$  یا  $\beta_0$  تبدیل می‌شود که اگر فرضیه‌ی  $H_0$  رد شود  $\hat{\beta}_{MR}$  و در صورتی که  $H_0$  رد نشود  $\beta_0$  نتیجه می‌شود.

خاصیت گسسته‌سازی تابع نشانگر باعث می‌شود که نتوان رفتار دو برآوردگر را توأمأً براساس  $PTMR$  بررسی کرد، لذا ایده استفاده از برآوردگری که بین این دو برآوردگر تغییر کند باعث مطرح شدن برآوردگر انقباضی نوع استاین شد. برای اطلاعات بیشتر جیمز و استاین (۱۹۶۱) را ببینید.

برآوردگر انقباضی نوع استاین آمیخته ریح<sup>۳۴</sup> (SSMR) به صورت زیر تعریف می‌شود

$$\begin{aligned} \hat{\beta}_{MR}^S &= \beta_0 + (1 - c\mathbf{L}_n^{-1}) (\hat{\beta}_{MR} - \beta_0) \\ &= \hat{\beta}_{MR} - c (\hat{\beta}_{MR} - \beta_0) \mathbf{L}_n^{-1} \end{aligned}$$

که در آن  $c = \frac{(p-2)(N-p)}{p(N-p+2)}$  و به ضریب  $(1 - c\mathbf{L}_n^{-1})$  عامل انقباض گویند. مقدار  $c$  برگرفته از صالح (۲۰۰۶) است.

با توجه به ساختار  $\hat{\beta}_{MR}^S$  و  $\hat{\beta}_{MR}^{PT}$  هر دو به ازای مقادیر بزرگ  $\mathbf{L}_n$ ،  $\hat{\beta}_{MR}$  را نتیجه می‌دهند. در حالیکه برای مقادیر کوچک  $\mathbf{L}_n$ ، رفتار  $\hat{\beta}_{MR}^{PT}$  متفاوت از  $\hat{\beta}_{MR}^S$  است. همچنین برخلاف برآوردگر  $PTMR$ ، برآوردگر نوع استاین یک تابع هموار و مستقل از سطح معناداری  $\alpha$  است. برای آگاهی بیشتر در خصوص برآوردگرهای ارائه‌شده، نوروزی راد (۱۳۹۶) را ببینید.

اکنون برای محاسبه‌ی این برآوردگرها الگوریتم ارائه‌شده توسط الیوت و همکاران (۲۰۱۱) را به صورت زیر توسعه داده‌ایم. همچنین برای مقدار اولیه پارامتر  $\lambda$  در الگوریتم محاسبه برآوردگر ریح این بار از روش پیشنهادی هورل و کنارد (۱۹۷۰) به صورت

$$\hat{\lambda}_0 = \frac{\hat{\sigma}^2}{\hat{\alpha}_{\max}^2}$$

استفاده می‌کنیم که در آن  $\hat{\alpha}_{\max} = \max(\hat{\alpha}_1, \dots, \hat{\alpha}_p)$  به طوری که  $\hat{\alpha} = \Gamma^T \hat{\beta}$  و  $\Gamma$  ماتریس متعامد شامل بردارهای ویژه در تجزیه طیفی  $\mathbf{X}^T \mathbf{X}$  است. لازم به یادآوری است که در تجزیه طیفی  $\mathbf{X}^T \mathbf{X}$  داریم

$$\mathbf{X}^T \mathbf{X} = \Gamma \Lambda \Gamma^T$$

به طوری که  $\Lambda = \text{diag}(\lambda_1^*, \dots, \lambda_p^*)$  و  $\lambda_i^*$  مقادیر ویژه  $i$ ام ماتریس  $\mathbf{X}^T \mathbf{X}$  است.

<sup>34</sup>Stein-type shrinkage mixed ridge

**الگوریتم ۴** الگوریتم EM در برآورد آزمون اولیه و نوع استاین پارامتر مدل آمیخته خطی ریبج

**گام ۱)** قرار می‌دهیم  $\hat{\theta}^{(t)} = \theta_0$  و  $\hat{\lambda}^{(t)} = \lambda_0$  سپس مقدار آماره‌های  $\hat{t}_\psi^{(t)}$  و  $\hat{t}_\lambda^{(t)}$  را برای مرحله tام به صورت زیر محاسبه می‌کنیم:

$$\begin{aligned}\hat{t}_\lambda^{(t)} &= E \left( \sum_{i=1}^n \epsilon_i^T \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \epsilon_i^T(\hat{\theta}^{(t)}) \epsilon_i(\hat{\theta}^{(t)}) + \text{trVar} \left( \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \right],\end{aligned}$$

$$\begin{aligned}\hat{t}_\psi^{(t)} &= E \left( \sum_{i=1}^n \mathbf{b}_i \mathbf{b}_i^T \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \mathbf{b}_i(\hat{\theta}^{(t)}) \mathbf{b}_i^T(\hat{\theta}^{(t)}) + \text{Var} \left( \mathbf{b}_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \right].\end{aligned}$$

**گام ۲)** برای به دست آوردن  $\hat{\theta}^{(t+1)}$ ، مقدار  $\hat{\sigma}^2(t+1)$  و  $\hat{\mathbf{D}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم:

$$\hat{\sigma}^2(t+1) = \hat{t}_\lambda^{(t)} / N,$$

$$\hat{\mathbf{D}}^{(t+1)} = \hat{t}_\psi^{(t)} / n.$$

که در آن  $N = \sum_{i=1}^n n_i$  و  $n$  برابر تعداد واحدهای آزمایشی است. حال برای محاسبه  $\hat{\mathbf{V}}^{(t+1)}$  داریم

$$\hat{\mathbf{V}}^{(t+1)} = \mathbf{Z} \hat{\mathbf{D}}^{(t+1)} \mathbf{Z}^T + \hat{\sigma}^2(t+1) \mathbf{I}.$$

**گام ۳)** ماتریس  $\hat{\mathbf{S}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم

$$\begin{aligned}\hat{\mathbf{S}}^{(t+1)} &= \mathbf{X}(\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{\lambda}^{(t)} \mathbf{I})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \\ &+ \mathbf{Z} \hat{\mathbf{D}}^{(t+1)} \mathbf{Z}^T \left( \mathbf{I} - \mathbf{X} \left[ (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{\lambda}^{(t)} \mathbf{I})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \right] \right),\end{aligned}$$

سپس برای محاسبه  $\hat{\lambda}^{(t+1)}$  داریم

$$\hat{\lambda}^{(t+1)} = \arg \min_{\lambda} \left\{ n^{-1} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) / (1 - \text{tr}(\hat{\mathbf{S}}^{(t+1)})/n) \right\}$$

**گام ۴)** در نهایت مقدار  $\hat{\beta}_{MR}^{(t+1)}$  از رابطه زیر بدست می‌آید

$$\hat{\beta}_{MR}^{(t+1)} = (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{\lambda}^{(t+1)})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{y}.$$

**گام ۵)** گام‌های (۱) تا (۴) را تا رسیدن به همگرایی  $(|\hat{\beta}_{MR}^{(t+1)} - \hat{\beta}_{MR}^{(t)}| < \gamma)$  به ازای یک مقدار  $\gamma$  کوچک و از پیش تعیین شده، تکرار می‌کنیم.

**گام ۶)** پس از محاسبه آماره  $L_n$  از رابطه (۲۲.۲)، داریم

$$\hat{\beta}_{MR}^{PT} = \hat{\beta}_{MR} - (\hat{\beta}_{MR} - \beta_0) \quad I(L_n < L_n(\alpha))$$

$$\hat{\beta}_{MR}^S = \hat{\beta}_{MR} - c(\hat{\beta}_{MR} - \beta_0) L_n^{-1}$$

## ۱.۴.۲ مطالعه شبیه‌سازی

برای مقایسه برآوردگرهای پیشنهادی  $\beta_{MR}^{PT}$  و  $\beta_{MR}^S$  با برآوردگر ارائه شده توسط الیوت و همکاران (۲۰۱۱) از مطالعه شبیه‌سازی مونت کارلویی استفاده می‌کنیم. در این شبیه‌سازی برای هر  $n = 40$  واحد آزمایشی  $n_i = 4$  را ثابت در نظر گرفته و همچنین

$$\beta = (0/0, 0/1, 0/2, 0/4, 0/8)$$

و  $b_{ijk} \sim \mathcal{N}(0, 0/6)$  و  $\varepsilon_{ijk} \sim \mathcal{N}(0, 1)$  فرض شده‌اند. برای تولید داده‌های هم‌خط از بسته‌ی EnvStat موجود در نرم‌افزار R استفاده شده است که هر متغیر پیش‌بین دارای توزیع نرمال با میانگین ۵ و واریانس ۱ می‌باشد که مقدار همبستگی بین متغیرهای پیش‌بین طبق جدول‌های ۵.۲ و ۶.۲ به صورت  $\rho \in \{0/0, 0/2, 0/5, 0/7, 0/9\}$  است. همچنین برای محاسبه‌ی خطای استاندارد (sd) و میانگین توان‌های دوم خطا (MSE)، روند شبیه‌سازی به تعداد  $B = 100$  بار تکرار شده است. لازم به ذکر است که در این مطالعه، به جهت بدست آوردن برآوردهای  $\hat{\beta}_{MR}^{PT}$  و  $\hat{\beta}_{MR}^S$  فرضیه  $H_0: \beta = 0$  در نظر گرفته شده است. نتایج این مطالعه در جدول‌های ۵.۲ و ۶.۲ آمده است.

جدول ۵.۲: نتایج مطالعه شبیه‌سازی برآوردگر نوع استاین در مدل آمیخته خطی ریج

$\rho$	$\beta$	مدل آمیخته خطی		مدل آمیخته خطی ریج		مدل آمیخته خطی نوع استاین		مدل آمیخته خطی ریج نوع استاین	
		برآورد	انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
0/0	0/0	-0/01651	0/01897	0/03828	0/01708	-0/01585	0/01833	0/03707	0/01651
	0/1	0/09016	0/01864	0/11907	0/01593	0/08725	0/01802	0/11513	0/01530
	0/2	0/19841	0/01832	0/21519	0/01559	0/19189	0/01777	0/20805	0/01498
	0/4	0/38544	0/02153	0/37015	0/01870	0/37298	0/02110	0/35803	0/01861
	0/8	0/83942	0/01785	0/75063	0/01868	0/81200	0/01640	0/72591	0/01979
MSE		0/294153		0/243236		0/271858		0/243966	
0/2	0/0	-0/00962	0/02271	0/06112	0/02014	-0/00914	0/02192	0/05919	0/01947
	0/1	0/09251	0/02057	0/12655	0/01728	0/08954	0/01989	0/12335	0/01659
	0/2	0/19622	0/02042	0/21971	0/01648	0/18980	0/01980	0/21244	0/01581
	0/4	0/37928	0/02446	0/35959	0/02064	0/36707	0/02396	0/34783	0/02057
	0/8	0/82806	0/01875	0/72513	0/02178	0/81081	0/01734	0/70121	0/02312
MSE		0/370953		0/312257		0/344432		0/307289	
0/5	0/0	-0/00016	0/02968	0/07565	0/02530	0/00004	0/02866	0/07331	0/02445
	0/1	0/09685	0/02508	0/14022	0/02006	0/09381	0/02424	0/12562	0/01926
	0/2	0/19709	0/02523	0/22428	0/01961	0/19065	0/02445	0/21697	0/01883
	0/4	0/35966	0/03151	0/34586	0/02530	0/34827	0/03089	0/33473	0/02511
	0/8	0/84249	0/02363	0/70636	0/02399	0/81521	0/02205	0/68334	0/02541
MSE		0/595086		0/442459		0/55458		0/440721	
0/7	0/0	0/01952	0/03965	0/11723	0/03073	0/01919	0/03831	0/11320	0/02973
	0/1	0/09350	0/03343	0/16157	0/02493	0/09062	0/03233	0/15625	0/02394
	0/2	0/18086	0/03260	0/22593	0/02314	0/17494	0/03164	0/21845	0/02226
	0/4	0/35057	0/03940	0/32776	0/02927	0/33946	0/03852	0/31720	0/02904
	0/8	0/85111	0/03275	0/65915	0/03311	0/82374	0/03098	0/62787	0/03349
MSE		1/026224		0/691952		0/958517		0/686111	
0/9	0/0	0/06336	0/06238	0/21256	0/02809	0/06172	0/05929	0/21040	0/02704
	0/1	0/05429	0/05000	0/21684	0/02302	0/05289	0/04841	0/21351	0/02192
	0/2	0/15672	0/05127	0/20109	0/02179	0/15134	0/04986	0/19881	0/02068
	0/4	0/21727	0/05493	0/29482	0/02799	0/30729	0/05360	0/28893	0/02728
	0/8	0/90250	0/05226	0/52754	0/05286	0/87430	0/05059	0/51431	0/05294
MSE		2/387184		1/09869		2/226858		1/083424	

جدول ۶.۲: نتایج مطالعه شبیه‌سازی برآوردگر آزمون اولیه در مدل آمیخته خطی ریبج

		PTMR, $\alpha = 0.01$		PTMR, $\alpha = 0.01$		PTMR, $\alpha = 0.05$		PTMR, $\alpha = 0.1$	
$\rho$	$\beta$	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
0/0	0/0	0/01246	0/00856	0/03605	0/01578	0/02828	0/01708	0/02828	0/01708
	0/1	0/02731	0/01393	0/010728	0/01466	0/011907	0/01593	0/011907	0/01593
	0/2	0/05582	0/02647	0/018226	0/01764	0/021519	0/01559	0/021519	0/01559
	0/4	0/09815	0/04943	0/032858	0/02632	0/037015	0/01870	0/037015	0/01870
	0/8	0/19518	0/09782	0/066143	0/04317	0/075063	0/01868	0/075063	0/01868
MSE		2/55996		0/55901		0/44334		0/44334	
0/2	0/0	0/01804	0/01010	0/05750	0/01853	0/06112	0/02014	0/06112	0/02014
	0/1	0/03039	0/01426	0/01355	0/01621	0/012655	0/01728	0/012655	0/01728
	0/2	0/06393	0/02630	0/018625	0/01822	0/021971	0/01648	0/021971	0/01648
	0/4	0/10230	0/04901	0/031808	0/02757	0/035959	0/02064	0/035959	0/02064
	0/8	0/20385	0/09668	0/063905	0/04445	0/072513	0/02178	0/072513	0/02178
MSE		2/50707		0/62413		0/31226		0/31226	
0/5	0/0	0/02520	0/01120	0/07434	0/02222	0/07565	0/02530	0/07565	0/02530
	0/1	0/02988	0/01521	0/012458	0/01873	0/014022	0/02006	0/014022	0/02006
	0/2	0/05970	0/02746	0/018715	0/02045	0/022438	0/01961	0/022438	0/01961
	0/4	0/09445	0/04990	0/031025	0/03054	0/034586	0/02530	0/034586	0/02530
	0/8	0/17984	0/09789	0/061851	0/04538	0/070636	0/02399	0/070636	0/02399
MSE		2/61252		0/72805		0/44246		0/44246	
0/7	0/0	0/03759	0/01581	0/01163	0/02780	0/011723	0/03073	0/011723	0/03073
	0/1	0/02963	0/01675	0/014117	0/02351	0/016157	0/02493	0/016157	0/02493
	0/2	0/06801	0/02819	0/018568	0/02332	0/022593	0/02314	0/022593	0/02314
	0/4	0/09917	0/04982	0/029229	0/03410	0/032776	0/02927	0/032776	0/02927
	0/8	0/20379	0/09593	0/058347	0/04987	0/065915	0/03311	0/065915	0/03311
MSE		2/55001		0/95592		0/69195		0/69195	
0/9	0/0	0/05136	0/01859	0/017949	0/03486	0/021356	0/03809	0/021356	0/03809
	0/1	0/06283	0/01780	0/018461	0/02241	0/021684	0/02302	0/021684	0/02302
	0/2	0/04737	0/02795	0/015783	0/02346	0/020109	0/02179	0/020109	0/02179
	0/4	0/07633	0/05073	0/025388	0/03453	0/029482	0/02799	0/029482	0/02799
	0/8	0/14872	0/01014	0/047422	0/06376	0/052754	0/05286	0/052754	0/05286
MSE		2/82244		1/37271		1/10987		1/10987	

همان‌طور که از نتایج جدول‌های ۵.۲ و ۶.۲ برمی‌آید، برآوردگر  $PTMR$  به سطح معنی‌داری  $\alpha$  وابسته است در حالی که این امر در خصوص برآوردگرهای  $SSMR$  صادق نیست. آنچه در تحلیل نتایج جدول‌ها حائز اهمیت است این است که برآوردگر  $SSMR$  در بین تمامی برآوردگرهای مورد مطالعه دارای میانگین توان دوم خطا ( $MSE$ ) و انحراف استاندارد ( $sd$ ) کمتری است. همچنین واضح است که هرچه میزان همخطی در داده‌ها قوی‌تر باشد،  $MSE$  برآورد در برآوردگر  $SSMR$  نسبت به دیگر برآوردگرها به میزان قابل ملاحظه‌ای کمتر است. به عبارت دیگر با افزایش همبستگی بین متغیرهای پیش‌بین، سرعت افزایش میزان  $MSE$  در برآوردگر  $SSMR$  کمتر از دیگر برآوردگرها است. در ادامه روش‌های برآوردیابی ذکر شده را برای تحلیل یک مثال واقعی به کار می‌گیریم.

## ۲.۴.۲ مثال واقعی

لنفوسیت‌ها<sup>۳۵</sup> گروهی از گلبول‌های سفید خون می‌باشند که نقش مهمی را در اجرای پاسخ‌های ایمنی بدن بر ضد عوامل بیگانه به عهده دارند. این سلول‌ها به دو گروه نوع T و B تقسیم می‌شوند که لنفوسیت‌های نوع T شامل دو گروه CD4 و CD8 هستند. سلول‌های CD4 را عموماً با عنوان سلول‌های T کمک کننده<sup>۳۶</sup> می‌شناسند که نقش مهمی را در دفاع علیه بسیاری از عوامل عفونی مانند باکتری‌ها، قارچ‌ها و حتی عوامل غیر عفونی دارند. این سلول‌ها مهمترین گروه سلولی تحت تأثیر ویروس HIV<sup>۳۷</sup> می‌باشند. ویروس HIV وارد این گروه از لنفوسیت‌ها شده و در این سلول‌ها شروع به همانند سازی کرده و در نهایت باعث از بین رفتن و کاهش تعداد آن‌ها می‌شود. هدف اصلی در درمان بیماران مبتلا به ویروس HIV، کند کردن روند کاهش سلول‌های CD4 و ثابت نگه‌داشتن آن‌ها در یک وضعیت، به منظور جلوگیری از ابتلای فرد به بیماری‌ها و عفونت‌های مختلف می‌باشد. از آن جایی که CD4 اولین گروه از سلول‌هایی هستند که تحت تأثیر ویروس HIV قرار می‌گیرند و تعداد آن‌ها سریعاً کاهش می‌یابد، اندازه‌گیری آن‌ها عامل مهمی برای ارزیابی میزان تأثیر درمان‌ها و پیشرفت بیماری به سوی فاز نهایی ایدز<sup>۳۸</sup> می‌باشد. بنابراین لازم است که همزمان با درمان در این افراد، تعداد سلول‌های CD4 به روش‌های خاص خود، اندازه‌گیری شود.

زالسیتابین<sup>۳۹</sup> و دیدانوزین<sup>۴۰</sup>، از جمله داروهای ضد ویروسی است که برای مقابله با ویروس HIV در فرد مبتلا، استفاده می‌شود. این داروها روند پیشروی ویروس HIV را کند می‌کنند و در نتیجه سبب تقویت سیستم ایمنی بدن شده و خطر ایجاد عفونت‌های وابسته به HIV را کاهش می‌دهند. هدف در این مطالعه تأثیر این دو نوع دارو بر کنترل ویروس HIV در بدن افراد مبتلا، مورد بررسی قرار می‌گیرد.

این مجموعه داده در بسته‌ی JMBayes موجود در نرم‌افزار R وجود دارد که شامل ۷ متغیر پیش‌بین برای تعداد  $n = 467$  واحد آزمایشی گردآوری شده که در این مطالعه متغیرهای gender و prevOI را با توجه به آنچه در فصل اول گفته شد به عنوان اثرات تصادفی در نظر گرفته‌ایم.

توضیحات مربوط به متغیرهای این مجموعه داده در جدول‌های ۷.۲ و ۸.۲ آمده است.

<sup>35</sup>Lymphocytus

<sup>36</sup>T-helper

<sup>37</sup>Human immunodeficiency virus

<sup>38</sup>Aids: Acquired immune deficiency syndrome

<sup>39</sup>Zalcitabine

<sup>40</sup>Didanosine

جدول ۷.۲: توصیف متغیرهای مجموعه داده بخش ۲.۴.۲

متغیرها	توضیحات
<i>patient</i>	کد شناسایی مربوط به شخص مورد آزمایش
<i>Time</i>	زمان مرگ یا سانسور
<i>death</i>	۰ نشان دهنده سانسور و ۱ نشان دهنده مرگ
<i>CD4</i>	تعداد سلول‌های نوع CD4 هر فرد
<i>obstime</i>	زمان ثبت تعداد سلول‌های CD4 (با اعمال تبدیل برای پیوسته شدن متغیر)
<i>drug</i>	عاملی با دو سطح <i>ddC</i> و <i>ddI</i> به ترتیب نشان دهنده تزریق داروی زالسیتابین و دیدانوزین به شخص مورد آزمایش
<i>gender</i>	جنسیت شخص مورد آزمایش
<i>prevOI</i>	عاملی با دو سطح <i>AIDS</i> و <i>noAIDS</i> به جهت نمایش نتیجه آزمایش قبلی
<i>AZT</i>	عاملی با دو سطح <i>intolerance</i> و <i>failure</i> برای تفکیک اشخاص از لحاظ مقاومت بدنی

جدول ۸.۲: اطلاعات خلاصه شده از متغیرهای مجموعه داده بخش ۲.۴.۲

	<i>Time</i>	<i>CD4</i>	<i>obstime</i>	
مینیمم	۰/۴۷	۰/۰۰	۰/۰۰	
چارک اول	۱۲/۲۳	۳/۱۶	۰/۰۰	
میانه	۱۴/۰۷	۵/۴۷	۲/۰۰	
میانگین	۱۳/۸۹	۷/۰۲	۴/۲۱	
چارک سوم	۱۷/۰۰	۱۰/۴۴	۶/۰۰	
ماکزیمم	۲۱/۴۰	۲۴/۱۲	۱۸/۰۰	

<i>death</i>	<i>drug</i>	<i>gender</i>	<i>prevOI</i>	<i>AZT</i>
death: ۴۱۲	ddI: ۶۸۸	male: ۱۲۸۸	AIDS: ۸۶۳	failure: ۴۹۱
censoring: ۹۹۳	ddC: ۷۱۷	female: ۱۱۷	noAIDS: ۵۴۲	intolerance: ۹۱۴

اندازه‌گیری هر کدام از متغیرها به دفعات برای هر واحد آزمایشی *z*ام تکرار شده است. لازم به ذکر است که در این مطالعه متغیر *CD4* به عنوان متغیر پاسخ در نظر گرفته شده است.

در جدول ۹.۲ برآوردگرهای آمیخته، آمیخته ریح، آمیخته نوع استاین، آمیخته ریح نوع استاین و آزمون اولیه آمیخته ریح حاصل از تحلیل این مجموعه داده به ترتیب با  $\hat{\beta}_M^S$ ،  $\hat{\beta}_{MR}^S$ ،  $\hat{\beta}_M^{PT}$  و  $\hat{\beta}_{MR}^{PT}$  نشان داده شده‌اند. برای بررسی عملکرد این برآوردگرها در تحلیل مثال واقعی، میانگین خطای پیش‌گویی<sup>۴۱</sup> (MPE) هر برآوردگر نیز محاسبه شده است که هرچه مقدار آن کمتر باشد نشان‌دهنده برتری برآوردگر است. طریقه محاسبه‌ی معیار MPE در ادامه آمده است.

به منظور محاسبه MPE با روش اعتبارسنجی متقابل<sup>۴۲</sup> K-fold، مجموعه داده را به‌طور تصادفی به K زیرمجموعه با اندازه‌های برابر تقسیم می‌کنیم، به ترتیب هر زیرمجموعه را به عنوان داده‌های آزمون  $\{(\mathbf{X}^{test}, \mathbf{y}^{test})\}$  و مابقی را به عنوان داده‌های آزمایش در نظر می‌گیریم و برای مدل‌بندی از آن‌ها استفاده می‌کنیم. برآورد حاصل مدل‌بندی داده‌های آزمایشی را  $\hat{\beta}_{train}$  می‌نامیم سپس مدل برازش شده را به جهت پیش‌گویی پاسخ مجموعه داده آزمون به کار می‌گیریم. در نهایت خطای پیش‌گویی (PE) را برای زیرمجموعه‌ی  $k = 1, \dots, K$  ام به صورت زیر محاسبه می‌کنیم

$$PE^k = \|\mathbf{X}_k^{test} \beta - \hat{\mathbf{y}}_k^{test}\|^2$$

که در آن  $\hat{\mathbf{y}}_k^{test} = \hat{\mathbf{X}}_k^{test} \hat{\beta}_k^{train}$ . این مقدار را برای همه‌ی  $k = 1, \dots, K$  زیرمجموعه‌ها محاسبه می‌کنیم. این فرآیند را  $N = 500$  بار تکرار کرده و مقدار MPE را از رابطه‌ی زیر بدست می‌آوریم

$$MPE = \text{median} \left\{ \frac{1}{K} \sum_{k=1}^K PE_1^k, \dots, \frac{1}{K} \sum_{k=1}^K PE_N^k \right\}$$

که در آن  $PE_i^k$  خطای پیش‌گویی  $k$ امین مجموعه‌ی آزمون در تکرار  $i$ ام است. مقدار MPE مربوط به هر برآوردگر در جدول ۹.۲ گزارش شده است.

همان‌طور که در جدول ۹.۲ ملاحظه می‌شود مقدار MPE برای برآوردگر آمیخته ریح نوع استاین  $\hat{\beta}_{MR}^S$  کمترین است.

جدول ۹.۲: برآورد ضرایب در داده‌های بخش ۲.۴.۲

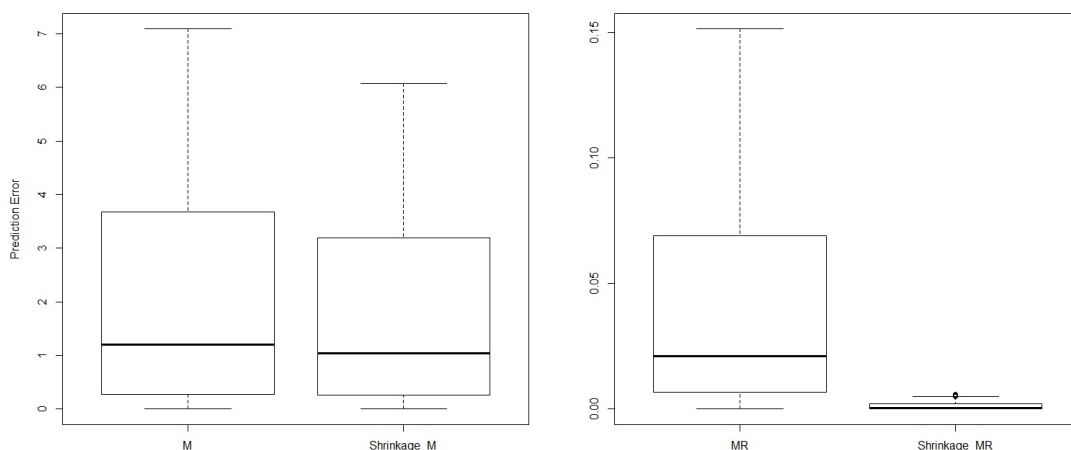
	$\hat{\beta}_M$	$\hat{\beta}_{MR}$	$\hat{\beta}_M^S$	$\hat{\beta}_{MR}^S$	$\hat{\beta}_{MR}^{PT} (\alpha > 0.2)$
Time	۰/۰۶۸۶۹۶۹۸	۰/۰۷۳۵۸۳	۰/۰۲۸۵۶۸۹۷	۰/۰۸۰۳۱۰۵۰	۰/۰۷۳۵۸۳
death	-۲/۰۴۹۱۸۵۰۹	-۱/۰۰۹۳۶۷۶	-۰/۰۸۵۲۱۹۳۳۸	-۰/۰۳۹۰۹۳۱۰۹	-۱/۰۰۹۳۶۷۶
obstime	-۰/۱۴۹۸۹۰۴۰	-۰/۱۴۳۴۹۹۱	-۰/۰۶۲۳۳۴۸۳	-۰/۰۵۵۵۷۷۶۳	-۰/۱۴۳۴۹۹۱
drug	۴/۵۲۵۳۱۱۹۹	۳/۳۳۳۵۶۹۲	۱/۸۸۱۹۳۸۷۹	۱/۲۹۱۱۰۱۳۷	۳/۳۳۳۵۶۹۲
AZT	۰/۶۶۳۸۱۸۹۲	۰/۱۷۲۴۱۱۸	۰/۲۷۶۰۶۱۹۸	۰/۰۶۶۷۷۵۶۰	۰/۱۷۲۴۱۱۸
MPE	۱/۲۰۶۸۷	۰/۰۲۱۱۴	۱/۰۴۶۴۴	۰/۰۰۰۵۲	۰/۰۲۱۱۴

<sup>۴۱</sup>Mean prediction error

<sup>۴۲</sup>Cross validation



شکل ۵.۲ نمودار جعبه‌ای مربوط به مقدار خطای پیش‌گویی برآوردگرها در هر مرحله از تکرار فرآیند محاسبه MPE را نشان می‌دهد که برتری برآوردگر آمیخته ريج نوع استاین را از نظر این معیار تایید می‌کند.



شکل ۵.۲: نمودار جعبه‌ای مربوط به مقدار MPE برآوردها

جدول ۱۰.۲ مربوط به مقدار انحراف استاندارد برآوردها است که همچنان نشان دهنده برتری برآوردگر آمیخته ريج نوع استاین نسبت به دیگر برآوردگرها است.

جدول ۱۰.۲: انحراف استاندارد برآوردهای جدول ۹.۲

	$sd(\beta_M)$	$sd(\beta_{MR})$	$sd(\beta_M^S)$	$sd(\beta_{MR}^S)$	$sd(\beta_{MR}^{PT})$
Time	۰/۰۰۰۶۱۳۰۱۲۷	$۱/۶۶۶۷۴۹ \times 10^{-۴}$	۰/۰۰۰۵۷۹۰۵۹۲	$۹/۲۳۵۰۲۷ \times 10^{-۵}$	$۱/۶۶۶۷۴۹ \times 10^{-۴}$
death	۰/۰۱۵۷۲۱۵۵۹۶	$۱/۱۷۰۷۷۷ \times 10^{-۲}$	۰/۰۱۵۲۳۳۵۸۶۶	$۳/۴۷۸۲۸۰ \times 10^{-۴}$	$۱/۱۷۰۷۷۷ \times 10^{-۲}$
obstime	۰/۰۰۱۱۲۶۲۴۳۱	$۴/۹۰۸۷۹۹ \times 10^{-۴}$	۰/۰۰۱۰۷۳۴۲۲۷	$۸/۷۷۹۷۲۲ \times 10^{-۵}$	$۴/۹۰۸۷۹۹ \times 10^{-۴}$
drug	۰/۰۳۷۴۳۸۵۵۶۷	$۱/۶۳۷۰۱۹ \times 10^{-۲}$	۰/۰۳۴۴۴۶۰۱۲۹	$۴/۵۸۹۶۹۴ \times 10^{-۴}$	$۱/۶۳۷۰۱۹ \times 10^{-۲}$
AZT	۰/۰۰۴۴۴۲۴۲۳۳	$۴/۲۳۴۵۳۴ \times 10^{-۳}$	۰/۰۰۴۱۲۱۵۷۲۲	$۸/۹۷۷۰۰۲ \times 10^{-۴}$	$۴/۲۳۴۵۳۴ \times 10^{-۳}$

## ۵.۲ مدل آمیخته خطی تعمیم‌یافته ريج

همان‌طور که در بخش ۵.۱ گفته شد، گاهی متغیر پاسخ از توزیع نرمال تبعیت نمی‌کند، بنابراین در چنین مواردی استفاده از مدل‌های خطی ممکن نخواهد بود. اما اگر توزیع متغیر پاس متعلق به خانواده توزیع‌های نمایی باشد، در این صورت

$$f_{y_i}(y_i) = \exp \left\{ (y_i \gamma_i - b(\gamma_i)) / \tau^2 - c(y_i, \tau) \right\} \quad (۲۳.۲)$$

لذا لگاریتم تابع درست‌نمایی به صورت زیر است

$$\ell(\beta) = \sum_{i=1}^n [y_i \gamma_i - b(\gamma_i)] / \tau^2 - \sum_{i=1}^n c(y_i, \tau). \quad (24.2)$$

### مدل خطی تعمیم‌یافته ریج

برای بدست آوردن برآوردگر خطی تعمیم‌یافته ریج<sup>۴۳</sup> (GLR) می‌توان لگاریتم تابع درست‌نمایی جریمه شده را به صورت زیر نوشت

$$\ell^{pen}(\beta) = \tau \left( \sum_{i=1}^n [y_i \gamma_i - b(\gamma_i)] / \tau^2 - \sum_{i=1}^n c(y_i, \tau) \right) - \lambda \beta^T \beta. \quad (25.2)$$

$$\begin{aligned} \frac{\partial \ell^{pen}(\beta)}{\partial \beta} &= \frac{\partial \ell^{pen}(\beta)}{\partial \gamma_i} \cdot \frac{\partial \gamma_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial g(\mu_i)} \cdot \frac{\partial g(\mu_i)}{\partial \beta} \\ &= \frac{\tau}{\tau^2} \sum \left( y_i - \frac{\partial b(\gamma_i)}{\partial \gamma_i} \cdot \frac{1}{\nu(\mu_i)} \cdot \frac{\partial \mu_i}{\partial g(\mu_i)} \cdot \mathbf{X}_i^T \right) - 2\lambda \beta \\ &= \frac{\tau}{\tau^2} \sum \left( y_i - \mu_i \cdot \frac{1}{\nu(\mu_i)} \cdot \frac{\partial \mu_i}{\partial g(\mu_i)} \cdot \frac{\partial g(\mu_i)}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial g(\mu_i)} \cdot \mathbf{X}_i^T \right) - 2\lambda \beta \\ &= \tau \sum (w_i \mathbf{X}_i^T (t_i - g(\mu_i))) - 2\lambda \beta \end{aligned} \quad (26.2)$$

که در آن  $t_i = g(\mu_i) + g_\mu(\mu_i)(y_i - \mu_i)$  و  $w_i = [\tau^2 \nu(\mu_i) g_\mu^2(\mu_i)]^{-1}$  شکل ماتریسی رابطه (۲۶.۲) را می‌توان به صورت زیر نوشت

$$\frac{\partial \ell^{pen}(\beta)}{\partial \beta} = \tau \mathbf{X}^T \mathbf{W} (\mathbf{T} - \mathbf{X}\beta) - 2\lambda \beta \quad (27.2)$$

با مساوی صفر قرار دادن رابطه (۲۷.۲) داریم

$$\begin{aligned} 0 &= \frac{\partial \ell^{pen}(\beta)}{\partial \beta} = \lambda \beta - (\mathbf{X}^T \mathbf{W} (\mathbf{T} - \mathbf{X}\beta)) \\ &= \lambda \beta - \mathbf{X}^T \mathbf{W} \mathbf{T} + (\mathbf{X}^T \mathbf{W} \mathbf{X}) \beta \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I}_p) \beta - (\mathbf{X}^T \mathbf{W} \mathbf{T}) \end{aligned}$$

در نتیجه برآوردگر GLR به صورت زیر حاصل می‌شود

$$\hat{\beta}_{GLR} = (\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{T}) \quad (28.2)$$

که در آن  $\mathbf{W} = \{n w_i\}$  ماتریس  $n$ -بعدی با عناصر قطری

$$w_i = [\tau^2 \nu(\mu_i) g_\mu^2(\mu_i)]^{-1}$$

<sup>43</sup>Generalized linear ridge

و بردار  $\mathbf{T} = \mathbf{X}\beta + \Delta(\mathbf{y} - \mu)$  بردار  $n \times 1$  وابسته به پارامتر  $\beta$ ، با دارایی‌های

$$t_i = g(\mu_i) + g_\mu(\mu_i)(y_i - \mu_i)$$

است.

از طرفی با تعریف بردار  $\mathbf{T}$ ، رابطه (۵۲.۱) را می‌توان به صورت زیر بازنویسی کرد

$$\begin{aligned} \hat{\beta}_{GL}^{(m+1)} &= \hat{\beta}_{GL}^{(m)} + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \Delta (\mathbf{y} - \mu) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} [\mathbf{X} \hat{\beta}_{GL}^{(m)} + \Delta (\mathbf{y} - \mu)] \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{T}^{(m)}. \end{aligned} \quad (29.2)$$

در نتیجه با استفاده از رابطه‌های (۲۸.۲) و (۲۹.۲) داریم

$$\begin{aligned} \hat{\beta}_{GLR} &= (\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{T}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{T}) \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{X}) \hat{\beta}_{GL}. \end{aligned} \quad (30.2)$$

برای آگاهی بیشتر در مورد برآورد پارامترهای مدل خطی تعمیم‌یافته در رگرسیون ریج به سرجرستد (۱۹۹۲) مراجعه کنید.

### مدل آمیخته خطی تعمیم‌یافته ریج

حال برای بدست آوردن برآوردگر آمیخته خطی تعمیم‌یافته ریج<sup>۴۴</sup> (GLMR)، لگاریتم تابع درست‌نمایی جریمه شده طبق رابطه‌های (۶۰.۱) و (۲۶.۲) به صورت زیر است

$$\ell^{pen}(\beta) = \lambda \beta^T \beta - \ell(\beta). \quad (31.2)$$

با مشتق‌گیری از رابطه (۳۱.۲) نسبت به  $\beta$  داریم

$$\begin{aligned} \frac{\partial \ell^{pen}(\beta)}{\partial \beta} &= 2\lambda\beta - 2 \int \mathbf{X}^T \mathbf{W} (\mathbf{T} - \mathbf{X}\beta - \mathbf{Z}\mathbf{b}) f_{\mathbf{b}|\mathbf{y}}(\mathbf{b}|\mathbf{y}) d\mathbf{b} \\ &= \lambda\beta + (\mathbf{X}^T E[\mathbf{W}|\mathbf{y}] \mathbf{X})\beta - (\mathbf{X}^T E[\mathbf{W}(\mathbf{T} - \mathbf{Z}\mathbf{b})|\mathbf{y}]) \\ &= (\mathbf{X}^T E[\mathbf{W}|\mathbf{y}] \mathbf{X} + \lambda \mathbf{I}_p) \beta - (\mathbf{X}^T E[\mathbf{W}(\mathbf{T} - \mathbf{Z}\mathbf{b})|\mathbf{y}]), \end{aligned} \quad (32.2)$$

در نتیجه برآوردگر GLMR به صورت زیر حاصل می‌شود

$$\hat{\beta}_{GLMR} = (\mathbf{X}^T E[\mathbf{W}|\mathbf{y}] \mathbf{X} + \lambda \mathbf{I}_p)^{-1} (\mathbf{X}^T E[\mathbf{W}(\mathbf{T} - \mathbf{Z}\mathbf{b})|\mathbf{y}]). \quad (33.2)$$

<sup>44</sup>Generalized linear mixed ridge

از طرفی، بردار  $\mathbf{T}$  به پارامتر  $\beta$  وابسته است لذا برآوردگر ارائه‌شده مشابه حالت  $GLR$  با استفاده از الگوریتم‌های تکرار شونده محاسبه خواهد شد. لذا

$$\hat{\beta}_{GLMR}^{(m+1)} = (\mathbf{X}^T E[\mathbf{W}|\mathbf{y}] \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \left( \mathbf{X}^T E[\mathbf{W}(\mathbf{T}^{(m)} - \mathbf{Zb})|\mathbf{y}] \right). \quad (34.2)$$

برای آگاهی بیشتر در مورد روش‌هایی مانند الگوریتم  $EM$ ، الگوریتم‌های  $MCMC$ <sup>۴۵</sup>، تقریب تصادفی<sup>۴۶</sup> برای محاسبه برآورد ماکزیمم درست‌نمایی در مدل‌های خطی تعمیم‌یافته به مک کولاک و همکاران (۲۰۰۸) مراجعه کنید.

---

<sup>45</sup>Markov chain Monte Carlo

<sup>46</sup>Stochastic approximation



## فصل ۳

# رگرسیون لیو برای داده‌های طولی

### ۱.۳ مقدمه

همان‌طور که در فصل قبل بیان شد، وجود همخطی در داده‌ها سبب بزرگ شدن و فاصله گرفتن مقادیر برآورد شده از پارامتر واقعی می‌شود. این امر در فصل قبل با استفاده از رگرسیون ریج پوشانده شد، اما ایده‌ای دیگر برای مرتفع کردن این مشکل، توسط لیو (۱۹۹۳) مطرح شد که با قرار دادن ضریبی مانند  $d$  ( $0 < d < 1$ )، مقدار برآورد شده را به پارامتر واقعی نزدیک می‌کند. به خاطر بیاورید که برآوردگر ریج برحسب پارامتر آن یک تابع غیرخطی است. حال چنانچه بخواهیم این رابطه خطی باشد می‌توانیم این طور فرض کنیم از آن جایی که  $\hat{\beta}$  از  $\beta$  فاصله دارد و مقادیر آن به دلیل وجود همخطی، بزرگ است،  $d\hat{\beta}$  به ازای یک  $d$  مناسب و کوچک می‌تواند تا حد امکان به  $\beta$  نزدیک شود لذا به جای این که تابع جریمه را به صورت  $\beta^T \beta$  در مدل وارد کنیم، بزرگی فاصله  $d\hat{\beta} - \beta$  را جریمه کرده و تابع جریمه را به صورت

$$(d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta)$$

در نظر می‌گیریم. در این صورت برآوردگر حاصل بر اساس مدل (۱.۲) عبارت است از

$$\hat{\beta}_{Liu} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + (d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta) \right\} \quad (1.3)$$

با مشتق‌گیری از عبارت داخل آکولاد در رابطه (۱.۳) نسبت به  $\beta$  داریم

$$\begin{aligned}
 & \frac{d}{d\beta} \left[ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + (d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta) \right] \\
 &= \frac{d}{d\beta} \left[ \mathbf{y}^T \mathbf{y} - 2\beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X} \beta + d^T \hat{\beta}^T \hat{\beta} - 2d\beta^T \hat{\beta} + \beta^T \beta \right] \\
 &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta - 2d\hat{\beta} + 2\beta \\
 &= -2\mathbf{X}^T \mathbf{y} - 2d\hat{\beta} + 2(\mathbf{X}^T \mathbf{X} + \mathbf{I}_p)\beta,
 \end{aligned}$$

که با مساوی صفر قرار دادن نتیجه می‌شود

$$\begin{aligned}
 \hat{\beta}_{Liu} &= (\mathbf{X}^T \mathbf{X} + \mathbf{I}_p)^{-1} (d\hat{\beta} + \mathbf{X}^T \mathbf{y}) \\
 &= (\mathbf{X}^T \mathbf{X} + \mathbf{I}_p)^{-1} \left( d(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{y} \right) \\
 &= (\mathbf{X}^T \mathbf{X} + \mathbf{I}_p)^{-1} \left( d(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X})(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \right) \\
 &= (\mathbf{X}^T \mathbf{X} + \mathbf{I}_p)^{-1} (d\mathbf{I}_p + \mathbf{X}^T \mathbf{X})(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\
 &= (\mathbf{X}^T \mathbf{X} + \mathbf{I}_p)^{-1} (d\mathbf{I}_p + \mathbf{X}^T \mathbf{X}) \hat{\beta}.
 \end{aligned} \tag{۲.۳}$$

آنچه مورد توجه است، خطی بودن برآوردگر نوع ليو نسبت به پارامتر  $d$  است. از این رو در این فصل، بررسی برآوردگر نوع ليو در مدل آمیخته خطی حائز اهمیت است. برای آگاهی بیشتر در خصوص برآوردگر ليو، برزویی بیدگلی (۱۳۹۳) را ببینید.

## ۲.۳ مدل آمیخته خطی ليو

مدل آمیخته خطی زیر را در نظر بگیرید

$$y_i = \mathbf{X}_i \beta + \mathbf{Z}_i^T \mathbf{b}_i + \varepsilon_i \tag{۳.۳}$$

که در آن برداری از  $n_i$  مشاهده،  $\mathbf{X}_i$  ماتریس طرح  $n_i \times p$  از اثرات ثابت و  $\beta = (\beta_1, \dots, \beta_p)^T$  بردار  $p \times 1$  ضرایب ثابت و نامعلوم، و  $\mathbf{Z}_i$  ماتریس طرح  $q \times n_i$  از اثرات تصادفی،  $\mathbf{b}_i \sim \mathcal{N}_q(0, \mathbf{D})$  بردار  $q \times 1$  ضرایب تصادفی و  $\varepsilon_i \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I})$  بردار  $n \times 1$  خطاهای تصادفی است. حال لگاریتم تابع درستنمایی را در رگرسیون ليو به صورت زیر داریم

$$\ell^{pen}(\beta, \mathbf{X}; \mathbf{y}) = -\frac{N}{\nu} \log(\nu \pi) - \frac{1}{\nu} \sum_{i=1}^n \log |\mathbf{V}| - \frac{1}{\nu} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) - (d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta) \tag{۴.۳}$$

ماکزیمم کردن تابع (۴.۳) معادل با مینیمم کردن عبارت وابسته به پارامتر  $\beta$  است. لذا برآوردگر آمیخته ليو (MLiu) از رابطه زیر محاسبه می‌شود

$$\hat{\beta}_{MLiu} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) + (d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta) \right\} \tag{۵.۳}$$

با مساوی صفر قرار دادن مشتق عبارت داخل آکولاد در رابطه (۵.۳) بر حسب  $\beta$  داریم

$$\begin{aligned} \circ &= \frac{d}{d\beta} \left[ (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) + (d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta) \right] \\ &= \frac{d}{d\beta} \left[ \mathbf{Y}^T \mathbf{V}^{-1} \mathbf{y} - \mathbf{y}^T \mathbf{V}^{-1} \mathbf{X}\beta - \beta^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}\beta \right. \\ &\quad \left. + d^T \hat{\beta}^T \hat{\beta} - d\hat{\beta}^T \beta + \beta^T \beta \right] \\ &= -2\mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} + 2\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}\beta - 2d\hat{\beta}^T + 2\beta \end{aligned}$$

بنابراین مشابه محاسبات رابطه (۱.۳) نتیجه می‌شود

$$\hat{\beta}_{MLiu} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) \hat{\beta} \quad (۶.۳)$$

همچنین امید ریاضی برآوردگر  $MLiu$  عبارت است از

$$\begin{aligned} E(\hat{\beta}_{MLiu}) &= E\left( (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) \hat{\beta} \right) \\ &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) \cdot E(\hat{\beta}) \\ &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) \beta \end{aligned}$$

به همین ترتیب

$$\begin{aligned} Var(\hat{\beta}_{MLiu}) &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \\ &\quad \times (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1}. \end{aligned}$$

واضح است که محاسبه  $\hat{\beta}_{MLiu}$  وابسته به ماتریس کواریانس متغیر پاسخ، یعنی

$$\mathbf{V} = \mathbf{ZDZ}^T + \sigma^2 \mathbf{I}$$

است که مستلزم برآورد پارامتر  $\theta = (\sigma, \mathbf{D})$  می‌باشد. لذا الگوریتم ۵ برای محاسبه  $\hat{\beta}_{MLiu}$  پیشنهاد می‌شود.

همچنین برای برآورد مقدار بهینه  $0 < d < 1$  در رابطه (۶.۳)، به روش اعتبارسنجی متقابل داریم

$$\hat{d} = \arg \min_{d \in (0,1)} \left\{ n^{-1} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) + 2n^{-2} tr(\mathbf{S}) (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) \right\} \quad (۷.۳)$$

که در آن  $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}_{MLiu}$  و

$$\begin{aligned} \mathbf{S} &= \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \\ &\quad + \mathbf{ZDZ}^T \left[ \mathbf{I} - \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \right]. \end{aligned}$$



**الگوریتم ۵ الگوریتم EM در برآورد پارامتر مدل آمیخته خطی لیو**

**گام ۱)** قرار می‌دهیم  $\hat{\theta}^{(t)} = \theta_0$  و  $\hat{d}^{(t)} = d_0$  سپس مقدار آماره‌های  $\hat{t}_\gamma^{(t)}$  و  $\hat{t}_\gamma^{(t)}$  را برای مرحله tام به صورت زیر محاسبه می‌کنیم:

$$\begin{aligned}\hat{t}_\gamma^{(t)} &= E \left( \sum_{i=1}^n \epsilon_i^T \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \epsilon_i^T(\hat{\theta}^{(t)}) \epsilon_i(\hat{\theta}^{(t)}) + \text{tr} \text{Var} \left( \epsilon_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \right],\end{aligned}\quad (۸.۳)$$

$$\begin{aligned}\hat{t}_\gamma^{(t)} &= E \left( \sum_{i=1}^n \mathbf{b}_i \mathbf{b}_i^T \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \\ &= \sum_{i=1}^n \left[ \mathbf{b}_i(\hat{\theta}^{(t)}) \mathbf{b}_i^T(\hat{\theta}^{(t)}) + \text{Var} \left( \mathbf{b}_i \mid \mathbf{y}_i, \hat{\beta}^{(t)}, \hat{\theta}^{(t)} \right) \right].\end{aligned}\quad (۹.۳)$$

**گام ۲)** برای به دست آوردن  $\hat{\theta}^{(t+1)}$ ، مقدار  $\hat{\sigma}^2(t+1)$  و  $\hat{\mathbf{D}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم:

$$\hat{\sigma}^2(t+1) = \hat{t}_\gamma^{(t)} / N,$$

$$\hat{\mathbf{D}}^{(t+1)} = \hat{t}_\gamma^{(t)} / n.$$

که در آن  $N = \sum_{i=1}^n n_i$  و  $n$  برابر تعداد واحدهای آزمایشی است. حال برای محاسبه  $\hat{\mathbf{V}}^{(t+1)}$  داریم

$$\hat{\mathbf{V}}^{(t+1)} = \mathbf{Z} \hat{\mathbf{D}}^{(t+1)} \mathbf{Z}^T + \hat{\sigma}^2(t+1) \mathbf{I}. \quad (۱۰.۳)$$

**گام ۳)** ماتریس  $\hat{\mathbf{S}}^{(t+1)}$  را از رابطه زیر محاسبه می‌کنیم

$$\begin{aligned}\hat{\mathbf{S}}^{(t+1)} &= \mathbf{X}(\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + d \mathbf{I}) (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \\ &+ \mathbf{Z} \mathbf{D} \mathbf{Z}^T \left[ \mathbf{I} - \mathbf{X}(\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{d}^{(t)} \mathbf{I}) (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \right],\end{aligned}$$

سپس برای محاسبه  $\hat{d}^{(t+1)}$  داریم

$$\hat{d}^{(t+1)} = \arg \min_{d \in (0,1)} \left\{ n^{-1} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) / (1 - \text{tr}(\hat{\mathbf{S}}^{(t+1)})/n) \right\} \quad (۱۱.۳)$$

**گام ۴)** در نهایت مقدار  $\hat{\beta}_{MLiu}^{(t+1)}$  از رابطه زیر بدست می‌آید

$$\hat{\beta}_{MLiu}^{(t+1)} = (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \hat{\mathbf{V}}^{-1(t+1)} \mathbf{X} + \hat{d}^{(t+1)} \mathbf{I}) \hat{\beta}^{(t+1)} \quad (۱۲.۳)$$

**گام ۵)** گام‌های (۱) تا (۴) را تا رسیدن به همگرایی  $\left( |\hat{\beta}_{MLiu}^{(t+1)} - \hat{\beta}_{MLiu}^{(t)}| < \gamma \right)$  به ازای یک مقدار  $\gamma$  کوچک و از پیش تعیین شده، تکرار می‌کنیم.

در ادامه چگونگی محاسبه  $\hat{\beta}_{MLiu}$  را در قالب یک مطالعه شبیه‌سازی و مثال واقعی با استفاده از نرم‌افزار R بررسی می‌کنیم.

### ۱.۲.۳ مطالعه شبیه‌سازی

در این بخش به بررسی عملکرد برآوردگر آمیخته لیو و مقایسه رفتار آن با برآوردگر آمیخته خطی در مدل آمیخته خطی خواهیم پرداخت.

در این مطالعه پس از تولید داده‌های هم‌خط مشابه بخش ۱.۲.۲، برای شبیه‌سازی مقادیر متغیر پاسخ، در رابطه‌ی (۳.۳)،  $\beta = (0, 0.4, 1.0, 1.6, 2.0)$ ،  $\varepsilon_{ijk} \sim \mathcal{N}(0, 1)$  و  $b_{ijk} \sim \mathcal{N}(0, 0.6)$  فرض شده‌اند به طوری که برای هر واحد آزمایشی  $i$  ام  $n_i = 4$  و همچنین هر متغیر پیش‌بین از توزیع نرمال با میانگین ۵ و واریانس ۱ در نظر گرفته شده است.

برای ایجاد همبستگی بین متغیرهای پیش‌بین در چند حالت مختلف مقدار  $\rho$  طبق جدول ۱.۳ از صفر تا ۰/۹ فرض شده است. همچنین برای محاسبه‌ی خطای استاندارد برآوردگر روند شبیه‌سازی مونت کارلویی به تعداد  $B = 100$  تکرار شده است. لازم به ذکر است که در این مطالعه تعداد نمونه‌ی آزمایشی  $n = 40$  فرض شده است. نتایج شبیه‌سازی در جدول ۱.۳ آمده است.

جدول ۱.۳: نتایج مطالعه شبیه‌سازی مدل آمیخته خطی لیو

		مدل آمیخته خطی		مدل آمیخته خطی لیو	
$\rho$	$\beta$	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
۰/۰	۰/۰	۰/۰۱۲۱۵۰۹۴	۰/۰۵۵۲۷۳۶	۰/۰۱۳۷۲۷۴۹	۰/۰۵۵۱۰۲۱
	۰/۴	۰/۴۱۱۱۵۹۱۱	۰/۰۵۶۲۰۲۴	۰/۴۱۱۹۶۸۳۹	۰/۰۵۶۰۰۴۴
	۱/۰	۰/۹۷۱۹۰۰۵۱	۰/۰۲۵۱۴۴۴	۰/۹۷۱۸۸۶۳۲	۰/۰۲۴۹۸۹۱
	۱/۶	۱/۶۰۸۸۹۰۷۴	۰/۰۳۹۸۷۶۹	۱/۶۰۷۹۷۹۲۴	۰/۰۳۹۵۹۳۰
	۲/۰	۱/۹۹۶۲۵۱۵۴	۰/۰۴۷۸۸۸۵	۱/۹۹۴۷۲۸۹۶	۰/۰۴۷۶۶۷۹
MSE		۱/۷۱۷۸۹		۱/۷۰۳۱۳۷	
۰/۳	۰/۰	۰/۰۰۹۴۸۹۱۳۹	۰/۰۸۳۴۹۶۷	۰/۰۱۱۶۱۲۶۶	۰/۰۸۳۱۴۷۵
	۰/۴	۰/۴۱۹۳۲۹۳۵۱	۰/۰۸۹۳۷۷۵	۰/۴۲۰۵۱۴۴۰	۰/۰۸۹۰۴۴۰
	۱/۰	۰/۹۶۵۸۵۷۹۳۴	۰/۰۴۷۸۸۸۵	۰/۹۶۵۸۶۰۰۸	۰/۰۴۷۶۴۲۴
	۱/۶	۱/۶۱۲۳۱۷۵۹۸	۰/۰۶۵۵۲۴۲	۱/۶۱۱۰۳۹۱۷	۰/۰۶۵۰۵۶۳
	۲/۰	۱/۹۹۳۷۵۹۶۴۰	۰/۰۴۴۳۲۴۹	۱/۹۹۱۶۶۲۵۴	۰/۰۴۴۰۵۴۷
MSE		۳/۷۶۳۶۵۴		۳/۷۲۷۵۵۷	
۰/۶	۰/۰	۰/۰۱۲۶۰۴۰۵	۰/۰۵۵۴۶۹۳	۰/۰۱۶۱۱۷۶۲	۰/۰۵۴۶۲۱۲
	۰/۴	۰/۴۳۵۸۹۲۳۰	۰/۰۲۹۱۱۵۶	۰/۴۳۷۴۳۶۳۳	۰/۰۲۸۵۱۴۴
	۱/۰	۰/۹۴۵۴۳۱۴۰	۰/۰۴۴۰۱۰۴	۰/۹۴۵۳۳۵۹۰	۰/۰۴۳۵۱۱۲
	۱/۶	۱/۶۱۴۲۲۲۲۶	۰/۰۱۹۲۲۸۰	۱/۶۱۲۵۲۰۹۵	۰/۰۱۸۹۲۵۹
	۲/۰	۱/۹۹۳۴۵۷۴۰	۰/۰۸۵۹۰۹۹	۱/۹۹۰۱۳۶۳۸	۰/۰۸۵۳۱۰۵
MSE		۲/۱۸۲۵۴۳		۲/۱۳۷۰۳۸	
۰/۹	۰/۰	۰/۰۰۲۲۷۶۵۷۶	۰/۰۷۸۶۰۴۵	۰/۰۱۱۰۷۲۱۷	۰/۰۷۸۰۶۵۴
	۰/۴	۰/۴۳۰۸۰۴۲۲۰	۰/۰۵۱۳۴۰۶	۰/۴۳۵۲۶۷۶۷	۰/۰۵۱۰۶۱۳
	۱/۰	۰/۹۱۰۵۳۹۸۹۲	۰/۰۸۱۲۳۸۷	۰/۹۱۱۳۴۲۷۸	۰/۰۸۰۸۷۲۹
	۱/۶	۱/۶۲۲۶۱۵۱۱۳	۰/۰۵۷۸۴۹۰	۱/۶۱۷۴۲۶۲۵	۰/۰۵۷۲۸۰۶
	۲/۰	۲/۰۳۶۸۷۲۰۸۹	۰/۰۸۵۵۰۰۱	۲/۰۲۷۹۰۵۱۷	۰/۰۸۵۰۰۳۱
MSE		۴/۱۸۲۱۹۲		۴/۱۳۰۰۶۶	

همان‌طور که در جدول ۱.۳ ملاحظه می‌شود، برای هر بخش از داده‌های تولید شده دارای میزان همبستگی متفاوت، برآوردهای حاصل از مدل آمیخته خطی لیو به مراتب دارای انحراف استاندارد کمتری نسبت به برآوردهای حاصل از مدل آمیخته خطی است.

## ۲.۲.۳ مثال واقعی

سلول‌های دستگاه ایمنی بدن در یک فرد، گاهی به دلایل گوناگون ممکن است نتواند برخی از مواد و اجزای بدن فرد را از مواد خارجی تشخیص دهد، از این رو نسبت به بافت‌های خود حساس شده و بر ضد آن‌ها آنتی کور<sup>۱</sup> (پادزهر یا موادی که بر علیه سم وارد شده در بدن به وجود می‌آیند) ترشح می‌کند. نتیجه حاصل را بیماری خود ایمنی<sup>۲</sup> می‌نامند که دستگاه ایمنی بدن به اشتباه، حمله به خود را آغاز می‌کند.

بیماری سیروز صفراوی اولیه<sup>۳</sup> (PBC) نوعی اختلال خود ایمنی است که کبد<sup>۴</sup> را تحت تأثیر قرار می‌دهد. به عبارت ساده‌تر در این نوع بیماری بدن احساس می‌کند که درون مجرای صفراوی کبد، جسم خارجی وجود دارد و سعی در از بین بردن پوشش به وجود آمده در این مجرا را دارد. این بیماری به کندی در بدن پیشرفت می‌کند که بسیاری از بیماران حتی ۱۰ الی ۱۵ سال پس از تشخیص نیز عارضه جدی و مهمی ندارند تا این که نهایتاً باعث مرگ سلول‌های کبدی و به حالت پلاستیکی درآمدن کبد می‌گردد که بیمار را نیازمند به انجام عمل پیوند کبد می‌کند یا منجر به مرگ بیمار می‌شود.

بیماری PBC یک بیماری نسبتاً نادر است که حدوداً از بین هر ۴۰۰۰ نفر یک نفر به آن مبتلا است. همچنین این نوع بیماری در زنان بسیار شایع‌تر است چرا که نسبت زنان مبتلا به مردان مبتلا، ۹ به ۱ است.

در این مطالعه تعداد ۳۱۲ بیمار به طور تصادفی از بین افراد مبتلا به خود ایمنی کبد یا سیروز صفراوی اولیه در یکی از معتبرترین مراکز و دانشگاه‌های علوم پزشکی جهان واقع در شهر روچستر<sup>۵</sup> ایالت مینه سوتای<sup>۶</sup> کشور آمریکا به نام کلینیک مایو<sup>۷</sup> مورد بررسی قرار گرفته‌اند. این داده‌ها در بسته JMBayes نرم‌افزار R موجود هستند و شامل ۱۹۴۵ مشاهده با ۲۰ متغیر می‌باشند که اندازه‌گیری هر کدام از این متغیرها برای بیمار  $i$ ام،  $n_i$  بار تکرار شده‌است. در جدول زیر به طور خلاصه به توضیح این متغیرها پرداخته شده و متغیرهایی که برای تحلیل این مثال، در مدل به عنوان اثرات تصادفی در نظر گرفته شده‌اند با علامت ■ مشخص شده‌اند. همچنین برآورد ضرایب حاصل از تحلیل این مجموعه داده در جدول ۳.۳ آورده شده‌است.

<sup>1</sup>Antibody

<sup>2</sup>Autoimmune

<sup>3</sup>Primary biliary cirrhosis

<sup>4</sup>Liver

<sup>5</sup>Rochester

<sup>6</sup>Minnesota

<sup>7</sup>Mayo Clinic

جدول ۲.۳: توصیف متغیرهای مجموعه داده بخش ۱.۲.۳

اثر تصادفی	توصیف	متغیرها
	کد شناسایی مربوط به فرد مورد آزمایش	<i>id</i>
	بازه زمانی ثبت مشخصات تا زمان مرگ یا انجام عمل پیوند	<i>years</i>
■	وضعیت بیمار (زنده، پیوندی، مرگ)	<i>status</i>
■	نوع داروی مصرفی (D-penicil – Placebo)	<i>drug</i>
	سن	<i>age</i>
■	جنسیت	<i>sex</i>
	بازه زمانی ثبت مشخصات تا زمان شروع درمان	<i>year</i>
■	دارای حفره شکم – از عوارض بیماری (بله – خیر)	<i>ascites</i>
■	بزرگ‌شدگی کبد – از عوارض بیماری (بله – خیر)	<i>hepatomegaly</i>
■	عوارض پوستی (بله – خیر)	<i>spiders</i>
■	افزایش میزان آب دفع‌شده از بدن (بله – خیر)	<i>edema</i>
	میزان بیلی‌روبین – نوعی سرم (میلی‌گرم بر دسی‌لیتر)	<i>serBilir</i>
	میزان کلسترول (میلی‌گرم بر دسی‌لیتر)	<i>serChol</i>
	میزان آلبومین – نوعی پروتئین (گرم بر دسی‌لیتر)	<i>albumin</i>
	آلکالین فسفاتاز – آنزیمی که توسط کبد ساخته می‌شود (یونیت بر لیتر)	<i>alkaline</i>
	اسگات – نوعی آنزیم در سلول‌های کبد (یونیت بر میلی‌لیتر)	<i>SGOT</i>
	پلاسمای خون (مکعب میلی‌لیتر بر ۱۰۰۰)	<i>platelets</i>
	آزمایش پروتوموبین – مدت زمان مکانیسم انعقاد خون بر حسب ثانیه	<i>prothrombin</i>
	مرحله بافت‌شناسی	<i>histologic</i>
■	وضعیت بیمار (زنده یا پیوندی، مرگ)	<i>status<sup>۲</sup></i>

جدول ۳.۳: برآورد ضرایب در مجموعه داده بخش ۱.۲.۳

	مدل آمیخته خطی	مدل آمیخته خطی لیو
age	-۰/۰۰۵۶	-۰/۰۰۵۳
year	۰/۷۰۹۵	۰/۷۰۸۵
serBilir	-۰/۱۷۱۱	-۰/۱۷۲۴
serChol	-۰/۰۰۰۴	-۰/۰۰۰۴
albumin	۲/۳۵۰۴	۲/۳۲۳۶
alkaline	۰/۰۰۰۱	۰/۰۰۰۱
SGOT	۰/۰۰۲۳	۰/۰۰۲۳
platelets	۰/۰۰۰۹	۰/۰۰۱۰
prothrombin	-۰/۰۰۸۰	-۰/۰۰۳۷
histologic	-۰/۵۰۸۶	-۰/۵۰۴۶
MPE	۰/۰۱۵۵۹	۰/۰۱۵۳۵

با توجه به آن چه در بخش ۲.۴.۲ در رابطه با معیار MPE گفته شد، در این مثال نیز برتری عملکرد مدل آمیخته خطی لیو نسبت به مدل آمیخته خطی از لحاظ معیار MPE تایید می شود.

### ۳.۳ مدل آمیخته خطی لیو با خطای بیضی گون

در این بخش با توجه به آن چه در بخش ۳.۲ در مورد توزیع های بیضی گون بیان شد، مدل آمیخته خطی لیو را با خطای بیضی گون تحت مفروضات مدل (۲۰.۲) مورد بررسی قرار می دهیم.

لگاریتم تابع درستنمایی متغیر  $y$  را بر پایه ی مدل آمیخته خطی با اعمال تابع جریمه در رگرسیون لیو، به صورت زیر داریم

$$\ell^{pen}(\mathbf{y}) = -\frac{N}{\nu} \log(\nu\pi) - \frac{1}{\nu} \sum_{i=1}^n \log |\mathbf{V}| + \log g \left( -\frac{1}{\nu} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\beta) \right) - (d\hat{\beta} - \beta)^T (d\hat{\beta} - \beta)$$

لذا برآوردگر  $MLiu$  پارامتر  $\beta$  با ماکزیمم کردن تابع فوق به صورت زیر خواهد بود

$$\hat{\beta}_{MLiu} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} + d\mathbf{I}) \hat{\beta}$$

که در آن  $0 < d < 1$  و از رابطه (۷.۳) برآورد می شود. در ادامه عملکرد این برآوردگر را تحت یک مثال شبیه سازی شده مورد بررسی قرار می دهیم.

### ۱.۳.۳ مطالعه شبیه‌سازی

در این بخش برای شبیه‌سازی مقادیر متغیر پاسخ، در رابطه‌ی (۲۰.۲) بردار ضرایب رگرسیونی به صورت  $\beta = (0, 0.4, 1, 1/6, 2)$ ،  $b_{ijk} \sim T(0, 0.3, 5)$  و  $\varepsilon_{ijk} \sim T(0, 1, 5)$  فرض شده‌اند به طوری که برای هر واحد آزمایشی  $i$  ام  $n_i = 4$  و هر متغیر پیش‌بین از توزیع نرمال با میانگین ۵ و واریانس ۱ فرض شده‌اند. همچنین برای ایجاد داده‌های هم‌خط از بسته‌ی EnvStat در نرم‌افزار R با ضریب همبستگی‌های متفاوت موجود در جدول ۴.۳  $\rho \in \{0/0, 0/2, 0/5, 0/7, 0/9\}$  استفاده شده است.

همچنین برای محاسبه‌ی خطای استاندارد برآورد بدست آمده با روند شبیه‌سازی مونت کارلویی به تعداد  $B = 100$  تکرار شده و تعداد نمونه آزمایشی  $n = 40$  در نظر گرفته شده است. در این راستا برای محاسبه برآوردگرها از الگوریتم ۵ استفاده کرده‌ایم. نتایج حاصل در جدول ۴.۳ آمده است.

جدول ۴.۳: نتایج مطالعه شبیه‌سازی مدل آمیخته خطی لیو با خطای بیضی گون

		مدل آمیخته خطی با خطای بیضی گون		مدل آمیخته خطی لیو با خطای بیضی گون	
$\rho$	$\beta$	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
۰/۰	۰/۰	۰/۰۱۳۲۵۷۸	۰/۱۱۲۹۳۱۰	۰/۰۷۸۰۴۲۰	۰/۱۱۲۷۱۸۰
	۰/۴	۰/۴۰۳۴۴۶۴	۰/۱۶۷۸۲۱۵	۰/۴۳۵۱۲۳	۰/۱۶۷۳۶۹۰
	۱/۰	۱/۰۲۷۷۳۴۸	۰/۱۸۶۱۹۰۶	۰/۹۱۷۰۵۰۵	۰/۱۸۵۹۶۱۰
	۱/۶	۱/۶۱۶۹۶۹۲	۰/۱۶۸۶۳۷۴	۱/۵۶۹۶۳۷۷	۰/۱۶۸۲۰۹۰
	۲/۰	۲/۰۲۰۱۴۶۹	۰/۱۷۵۵۲۷۹	۱/۹۹۳۶۰۳۴	۰/۱۷۵۲۷۰۳
MSE		۲۱/۵۷۴۹۳		۲۱/۵۰۵۲۷	
۰/۳	۰/۰	۰/۰۱۵۵۱۲۶	۰/۱۱۵۳۱۶۰	۰/۰۱۵۸۴۴۵	۰/۱۱۵۰۰۲۴
	۰/۴	۰/۴۴۰۱۳۰۶	۰/۱۳۸۵۱۷۰	۰/۴۴۹۷۰۵۸	۰/۱۳۸۱۵۵۳
	۱/۰	۰/۹۸۳۳۸۶۰	۰/۱۹۳۵۸۰۶	۱/۰۲۱۹۴۷۳	۰/۱۹۳۱۰۰۱
	۱/۶	۱/۶۰۱۴۶۴۲	۰/۱۹۲۷۵۵۴	۱/۵۷۹۵۰۳۲	۰/۱۹۲۲۹۶۹
	۲/۰	۲/۰۲۳۵۵۹۳	۰/۱۷۷۹۵۵۵	۲/۰۳۶۶۷۲۱	۰/۱۷۷۷۳۳۸
MSE		۲۲/۲۰۷۶۴		۲۲/۱۱۱۷۷	
۰/۶	۰/۰	۰/۰۱۲۶۳۳۷	۰/۱۱۹۴۳۸۹	۰/۰۵۲۶۲۷۲	۰/۱۱۹۰۳۳۰
	۰/۴	۰/۳۵۰۰۵۰۰	۰/۱۳۸۲۲۱۱	۰/۴۸۱۹۱۳۴	۰/۱۳۷۷۴۷۶
	۱/۰	۱/۰۴۲۷۵۹۸	۰/۱۱۷۹۷۳۰	۰/۹۱۵۴۹۱۸	۰/۱۱۷۳۵۳۳
	۱/۶	۱/۵۷۶۵۳۹۶	۰/۱۴۹۳۹۵۳	۱/۶۸۱۴۶۱۳	۰/۱۴۸۷۱۸۲
	۲/۰	۲/۰۵۳۰۹۹۲	۰/۱۸۰۵۵۵۱	۱/۹۸۷۱۶۳۸	۰/۱۷۹۶۵۴۵
MSE		۱۶/۳۶۱۰۵		۱۶/۲۳۲۶۹	
۰/۹	۰/۰	-۰/۰۴۳۸۰۳۲	۰/۱۴۴۹۸۲۳	۰/۰۸۰۶۶۰۷	۰/۱۳۵۴۷۰۱
	۰/۴	۰/۴۱۶۳۰۵۷۸	۰/۱۵۲۷۷۵۱	۰/۴۹۱۷۰۱۳	۰/۱۴۰۱۷۱۸
	۱/۰	۱/۰۱۱۰۰۵۹	۰/۱۹۳۷۲۱۷	۱/۱۶۵۵۲۱۷	۰/۱۸۵۰۲۴۲
	۱/۶	۱/۵۵۶۶۵۷۱	۰/۱۶۲۲۸۸۹	۱/۷۲۸۸۹۷۲	۰/۱۵۳۸۷۰۹
	۲/۰	۲/۰۲۳۱۸۶۰۶	۰/۱۹۴۳۸۶۸	۲/۰۹۱۲۵۸۶	۰/۱۸۳۸۴۷۵
MSE		۲۳/۳۴۴۲۴		۲۰/۸۲۰۹۲	

همان‌طور که در جدول ۴.۳ ملاحظه می‌شود برآوردهای حاصل از روش آمیخته لیو با خطاهای بیضی گون دارای انحراف استاندارد کمتری نسبت به برآورد حاصل از روش آمیخته خطی است. مقادیر میانگین توان دوم خطای برآورد (MSE) نیز برتری روش آمیخته لیو را نسبت به روش آمیخته خطی تایید می‌کند.



## ۴.۳ برآوردهای انقباضی در مدل آمیخته خطی لیو

در این بخش برآوردهای آزمون اولیه و انقباضی نوع استاین را در مدل آمیخته خطی لیو مورد بررسی قرار می‌دهیم.

فرض کنید در مدل آمیخته خطی لیو، فرضیه زیر مطرح باشد

$$\begin{cases} H_0 : \beta = \beta_0 \\ H_1 : \beta \neq \beta_0 \end{cases}$$

در این صورت طبق آنچه در بخش (۴.۲) بیان شد، برآوردهای آزمون اولیه آمیخته خطی لیو<sup>۸</sup> ( $PTMR$ ) به صورت زیر تعریف می‌شود.

$$\hat{\beta}_{MLiu}^{PT} = \begin{cases} \hat{\beta}_{MLiu} & \text{فرضیه‌ی } H_0 \text{ رد شود} \\ \beta_0 & \text{فرضیه‌ی } H_0 \text{ پذیرفته شود} \end{cases}$$

برآوردهای  $PTMLiu$  را می‌توان با استفاده از تابع نشانگر به صورت زیر بازنویسی کرد

$$\hat{\beta}_{MLiu}^{PT} = \hat{\beta}_{MLiu} - (\hat{\beta}_{MLiu} - \beta_0) \mathbf{I}_{(\mathbf{L}_n < \mathbf{L}_n(\alpha))} \quad (13.3)$$

که در آن آماره والد آزمون فرضیه‌ی  $H_0$  به صورت

$$\mathbf{L}_n = (\hat{\beta} - \beta_0)^T (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}) (\hat{\beta} - \beta_0) \quad (14.3)$$

و  $\mathbf{L}_n(\alpha)$  صدک بالایی  $100\alpha$  درصد توزیع کی-دو با  $p$  درجه آزادی است. به همین ترتیب برآوردهای انقباضی نوع استاین آمیخته لیو<sup>۹</sup> ( $SSMLiu$ ) به صورت زیر تعریف می‌شود

$$\begin{aligned} \hat{\beta}_{MLiu}^S &= \beta_0 + (1 - c\mathbf{L}_n^{-1}) (\hat{\beta}_{MLiu} - \beta_0) \\ &= \hat{\beta}_{MLiu} - c (\hat{\beta}_{MLiu} - \beta_0) \mathbf{L}_n^{-1} \end{aligned}$$

که در آن  $c = \frac{(p-2)(N-p)}{p(N-p+2)}$

برای محاسبه این برآوردها در مدل آمیخته خطی لیو، گام ۶ام را در الگوریتم ۵ به صورت زیر داریم

<sup>۸</sup>Preliminary test mixed Liu

<sup>۹</sup>Stein-type shrinkage mixed Liu

الگوریتم ۶ الگوریتم EM در برآورد آزمون اولیه و نوع استاین پارامتر مدل آمیخته خطی ليو

گام‌های (۱) - (۵)

به الگوریتم ۵ مراجعه شود.

گام ۶) پس از محاسبه آماره  $L_n$  از رابطه (۱۴.۳)، داریم

$$\hat{\beta}_{MLiu}^{PT} = \hat{\beta}_{MLiu} - (\hat{\beta}_{MLiu} - \hat{\beta}_0) I(L_n < L_n(\alpha)) \quad (15.3)$$

$$\hat{\beta}_{MLiu}^S = \hat{\beta}_{MLiu} - c(\hat{\beta}_{MLiu} - \hat{\beta}_0)L_n^{-1} \quad (16.3)$$

### ۱.۴.۳ مطالعه شبیه‌سازی

در این مطالعه برای هر  $n = 40$  واحد آزمایشی  $n_i = 4$  را ثابت در نظر گرفته و همچنین

$$\beta = (0/0, 0/1, 0/2, 0/4, 0/8)$$

برای تولید داده‌های هم‌خط از بسته EnvStat موجود در نرم‌افزار R استفاده شده است که هر متغیر پیش‌بین دارای توزیع نرمال با میانگین ۵ و واریانس ۱ می‌باشد و مقدار همبستگی بین متغیرهای پیش‌بین طبق جدول ۵.۳ به صورت  $\rho \in \{0/0, 0/2, 0/5, 0/7, 0/9\}$  در نظر گرفته شده است. همچنین برای محاسبه‌ی خطای استاندارد و میانگین خطای برآورد (MSE)، روند شبیه‌سازی به تعداد  $B = 100$  بار تکرار شده است. لازم به ذکر است که در این مطالعه، فرضیه  $H_0: \beta = 0$  در نظر گرفته شده است. نتایج این مطالعه در جدول ۵.۳ آمده است.

جدول ۵.۳: نتایج مطالعه شبیه‌سازی برآوردگر نوع استاین در مدل آمیخته خطی لیو

		مدل آمیخته خطی		مدل آمیخته خطی لیو		مدل آمیخته خطی نوع استاین		مدل آمیخته خطی لیو نوع استاین	
$\rho$	$\beta$	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد	برآورد	انحراف استاندارد
۰/۰	۰/۰	-۰/۰۱۴۱۷	۰/۰۱۴۶۵	۰/۰۲۶۳۷	۰/۰۱۴۲۰	-۰/۰۱۵۲۲	۰/۰۱۴۲۲	۰/۰۲۴۱۲	۰/۰۱۴۰۸
	۰/۱	۰/۰۰۹۲۷۰	۰/۰۱۷۶۲	۰/۱۲۰۶۸	۰/۰۱۷۵۹	۰/۰۰۹۱۳۶	۰/۰۱۷۶۰	۰/۱۱۹۷۲	۰/۰۱۷۲۱
	۰/۲	۰/۰۱۹۷۲۴	۰/۰۱۹۴۴	۰/۱۹۸۳۲	۰/۰۱۹۳۲	۰/۱۹۸۱۴	۰/۰۱۹۳۹	۰/۲۰۲۱۶	۰/۰۱۹۱۸
	۰/۴	۰/۰۴۱۱۷۲	۰/۰۱۸۱۶	۰/۳۹۲۹۸	۰/۰۱۸۱۰	۰/۰۴۰۲۰۹	۰/۰۱۸۱۲	۰/۳۹۷۸۲	۰/۰۱۷۹۴
	۰/۸	۰/۰۸۱۲۵۴	۰/۰۲۱۶۹	۰/۸۳۱۶۲	۰/۰۲۱۰۳	۰/۸۲۰۴۶	۰/۰۲۱۴۲	۰/۸۳۳۷۱	۰/۰۲۰۸۸
MSE		۰/۲۲۳۰۷۵		۰/۲۶۶۸۴۴		۰/۲۶۸۷۴۷		۰/۲۶۱۲۶۹	
۰/۲	۰/۰	۰/۰۰۱۰۲	۰/۰۱۸۲۳	۰/۰۳۰۲۹	۰/۰۱۸۰۶	-۰/۰۰۱۱۶	۰/۰۱۸۲۰	۰/۰۲۹۳۵	۰/۰۱۷۷۴
	۰/۱	۰/۰۱۰۲۴۴	۰/۰۱۸۴۰	۰/۱۱۹۰۳	۰/۰۱۸۲۰	۰/۱۱۳۱۶	۰/۰۱۸۳۴	۰/۱۱۷۲۲	۰/۰۱۸۱۱
	۰/۲	۰/۲۰۸۰۶	۰/۰۱۹۸۷	۰/۲۲۴۱۲	۰/۰۱۹۷۱	۰/۲۰۱۸۴	۰/۰۱۹۸۳	۰/۲۲۳۴۷	۰/۰۱۹۶۵
	۰/۴	۰/۰۴۲۰۱۶	۰/۰۲۲۳۷	۰/۳۸۰۴۶	۰/۰۲۲۱۷	۰/۰۴۱۷۲۹	۰/۰۲۲۳۲	۰/۳۷۹۲۴	۰/۰۲۲۰۳
	۰/۸	۰/۰۸۱۶۰۹	۰/۰۱۸۸۲	۰/۸۳۳۷۲	۰/۰۱۸۶۶	۰/۸۲۱۰۲	۰/۰۱۸۷۹	۰/۸۳۳۸۹	۰/۰۱۸۶۵
MSE		۰/۳۰۷۹۸۸		۰/۳۰۵۰۷۵		۰/۳۰۶۸۵		۰/۳۰۱۱۳۳	
۰/۵	۰/۰	۰/۰۱۱۰۲	۰/۰۳۰۱۶	۰/۰۶۸۲۳	۰/۰۳۰۰۰	۰/۰۱۲۱۶	۰/۰۳۰۰۹	۰/۰۶۹۱۳	۰/۰۲۹۷۵
	۰/۱	۰/۰۱۰۵۵۸	۰/۰۲۴۵۸	۰/۱۳۲۵۹	۰/۰۲۴۳۴	۰/۰۹۸۲۲	۰/۰۲۴۵۲	۰/۱۳۳۰۵	۰/۰۲۴۱۶
	۰/۲	۰/۲۰۲۳۱	۰/۰۲۴۱۳	۰/۲۲۳۱۶	۰/۰۲۳۸۴	۰/۱۹۹۴۶	۰/۰۲۴۰۳	۰/۲۱۹۲۴	۰/۰۲۳۶۹
	۰/۴	۰/۰۳۷۷۸۴	۰/۰۳۰۴۵	۰/۳۷۷۰۲	۰/۰۳۰۲۴	۰/۳۷۷۸۴	۰/۰۳۰۴۱	۰/۳۷۰۱۱	۰/۰۳۰۱۹
	۰/۸	۰/۰۸۱۴۲۹	۰/۰۲۲۴۹	۰/۷۹۶۱۸	۰/۰۲۲۱۸	۰/۸۱۵۰۶	۰/۰۲۲۳۸	۰/۷۹۵۹۸	۰/۰۲۲۱۳
MSE		۰/۵۶۵۵۰۳		۰/۵۶۱۸۰۲		۰/۵۶۲۴۱		۰/۵۵۶۱۳۶	
۰/۷	۰/۰	۰/۰۰۶۲۸	۰/۰۳۲۴۸	۰/۰۹۳۲۴	۰/۰۳۲۰۱	۰/۰۰۶۲۸	۰/۰۳۲۳۷	۰/۰۸۲۵۲	۰/۰۳۱۸۴
	۰/۱	۰/۰۱۱۲۸۰	۰/۰۳۱۳۸	۰/۱۶۶۳۹	۰/۰۳۰۹۷	۰/۱۰۳۰۱	۰/۰۳۱۳۱	۰/۱۶۷۳۹	۰/۰۲۹۹۸
	۰/۲	۰/۲۰۶۴۳	۰/۰۳۳۴۹	۰/۲۰۱۱۹	۰/۰۳۳۰۵	۰/۲۰۵۴۰	۰/۰۳۳۴۰	۰/۲۰۱۰۳	۰/۰۳۲۲۰
	۰/۴	۰/۰۴۱۵۳۵	۰/۰۳۸۹۰	۰/۴۰۳۲۱	۰/۰۳۸۵۱	۰/۰۴۱۴۲۵	۰/۰۳۸۸۴	۰/۰۴۱۵۹۸	۰/۰۳۸۳۲
	۰/۸	۰/۰۸۴۲۱۸	۰/۰۳۱۰۶	۰/۷۲۳۰۹	۰/۰۲۹۸۲	۰/۸۴۱۴۷	۰/۰۳۰۵۴	۰/۷۱۴۸۵	۰/۰۲۹۶۹
MSE		۰/۹۰۴۵۲۶		۰/۸۹۰۷۶۲		۰/۸۹۵۵۶۸		۰/۸۶۶۷۵۳	
۰/۹	۰/۰	۰/۰۳۹۲۸	۰/۰۵۸۲۰	۰/۱۹۱۳۶	۰/۰۵۷۶۵	۰/۰۳۷۱۴	۰/۰۵۸۰۲	۰/۱۹۰۲۹	۰/۰۵۷۳۳
	۰/۱	۰/۰۸۶۱۲	۰/۰۴۸۶۳	۰/۱۵۶۴۲	۰/۰۴۷۷۱	۰/۰۸۲۸۴	۰/۰۴۸۵۰	۰/۱۵۷۱۳	۰/۰۴۴۷۰
	۰/۲	۰/۲۰۷۶۴	۰/۰۵۲۴۹	۰/۲۳۴۰۶	۰/۰۵۱۸۱	۰/۲۰۳۸۴	۰/۰۵۲۳۱	۰/۲۳۵۸۳	۰/۰۵۱۴۳
	۰/۴	۰/۳۶۸۲۷	۰/۰۵۳۸۴	۰/۳۳۷۵۶	۰/۰۵۳۱۱	۰/۳۷۴۰۲	۰/۰۵۳۷۲	۰/۳۳۷۲۹	۰/۰۵۲۶۲
	۰/۸	۰/۸۸۷۶۱	۰/۰۴۹۶۵	۰/۶۹۳۱۶	۰/۰۴۸۷۶	۰/۸۷۲۶۲	۰/۰۴۹۵۱	۰/۶۹۳۰۹	۰/۰۴۸۵۶
MSE		۲/۲۲۹۸۶۶		۲/۲۱۳۴۳۶		۲/۲۱۴۴۵۵		۲/۱۴۵۲۰۸	

آنچه در تحلیل نتایج جدول ۵.۳ حائز اهمیت است این است که برآوردگر  $SSMLiu$  در بین تمامی برآوردگرهای مورد مطالعه دارای میانگین توان دوم خطا (MSE) و انحراف استاندارد کمتری است.

## ۵.۳ نتیجه‌گیری و آینده تحقیق

### ۱.۵.۳ خلاصه و نتیجه‌گیری

مهمترین ویژگی داده‌های طولی، اندازه‌گیری‌های مکرر از واحدهای تحت آزمایش در طول زمان است لذا هدف اول در مطالعات طولی، بررسی تغییرات پاسخ در طول زمان و همچنین عوامل مؤثر بر این تغییرات است. از آنجا که اندازه‌گیری بر روی واحدهای تحت مطالعه تکرار می‌شود، امکان بررسی تغییرات درون هر واحد نیز فراهم می‌شود. در جمع‌آوری داده‌های طولی، اندازه‌گیری مکرر متغیری خاص برای یک واحد آزمایشی در موقعیت‌ها و زمان‌های مختلف، همبستگی بین مشاهدات مربوط به آن واحد را نشان می‌دهد. همان‌طور اشاره شد

یکی از روش‌های مدل‌بندی این داده‌ها استفاده از مدل‌های آمیخته خطی است. این مدل‌ها همبستگی بین مشاهدات را با معرفی یک یا چند متغیر تصادفی به عنوان اثر تصادفی در مدل در نظر می‌گیرند. گاهی اوقات مشکل همخطی در داده‌های جمع‌آوری شده به وجود می‌آید که از روش‌های مقابله با آن استفاده از برآوردگرهای انقباضی است. از این‌رو برای تحلیل داده‌های طولی که دارای مشکل همخطی می‌باشند می‌توان از ترکیب مدل‌های آمیخته خطی با رگرسیون ریح از این مهم عبور کرد. همانطور که در نتایج مربوط به مشاهدات شبیه‌سازی ملاحظه می‌شود، برآوردگرهای آمیخته خطی ریح و آمیخته خطی لیو همواره دارای انحراف استاندارد و  $MSE$  کمتری نسبت به برآوردگرهای آمیخته خطی می‌باشند، مخصوصاً زمانی که همخطی در داده‌ها قوی‌تر باشد. همچنین در قسمت تحلیل داده‌های واقعی نیز این برتری با استفاده از معیار  $MPE$  تأیید می‌شود.

گاهی اوقات تعیین توزیع دقیق متغیر پاسخ در مثال‌های واقعی کار دشواری است در حالی که می‌توان بر اساس داده‌ها به متقارن بودن یا نبودن آن‌ها پی برد لذا از این حیث استفاده از مدل‌های ارائه شده با خطاهای دارای توزیع بیضی‌گون مشکل عدم اطلاع در مورد توزیع دقیق مؤلفه پاسخ برای یافتن برآوردگر  $ML$  را مرتفع می‌کند نتایج حاصل از مطالعات شبیه‌سازی نیز نشان دهنده کمتر بودن انحراف استاندارد و  $MSE$  برآوردگرها در مدل‌های آمیخته خطی با خطای بیضی‌گون نسبت به برآوردگرهای حاصل از مدل‌های آمیخته خطی است.

گاهی با وجود یک اطلاع پیشین در مورد پارامتر مدل، علاقه‌مند به انجام آزمون فرضیه مورد نظر هستیم لذا با هدف بهبود برآوردگرهای حاصل از مدل‌های آمیخته خطی، برآوردگرهای انقباضی نوع استاین در مدل آمیخته خطی ریح و آمیخته لیو مورد بررسی قرار داده شد. برتری برآوردگرهای انقباضی نوع استاین در مدل آمیخته خطی ریح را می‌توان در نتایج حاصل از مطالعات شبیه‌سازی ملاحظه کرد.

### ۲.۵.۳ پیشنهادات برای پژوهش‌های آینده

طبق نتایج بدست آمده پیرامون مباحث و موضوعات مطرح شده در این مجموعه می‌توان به موارد زیر به عنوان پیشنهاد برای پژوهش‌های آینده اشاره کرد:

- استفاده از معیارهای دیگری برای مقایسه برآوردگرهای ارائه شده با برآوردگر حاصل از مدل آمیخته خطی.
- گنجاندن تمامی نتایج حاصل از این مجموعه در قالب یک مطالعه شبیه‌سازی با هدف مقایسه برآوردگرهای ارائه شده با یکدیگر.
- بررسی رفتار برآوردگر آمیخته تعمیم‌یافته ریح تحت یک مطالعه شبیه‌سازی.
- استفاده از روش‌های دیگری همچون رگرسیون لاسو در مدل آمیخته خطی و مقایسه آن با مدل آمیخته خطی ریح.



# مراجع

[۱] آرست، م. (۱۳۹۵)، پایان‌نامه کارشناسی ارشد: ”مقایسه رفتار برخی برآوردگرهای انقباضی بریج در مدل رگرسیون چندگانه“، دانشکده علوم ریاضی، دانشگاه صنعتی شاهرود.

[۲] آرشی، م. (۱۳۸۷)، رساله دکتری: ”برآوردگرهای بهبودیافته در مدل رگرسیون خطی چندگانه با خطاهای دارای توزیع بیضی‌گون“، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد.

[۳] برزویی بیدگلی، م. (۱۳۹۳)، پایان‌نامه کارشناسی ارشد: ”بررسی برآوردگر جک نایف ریج“، دانشکده علوم ریاضی، دانشگاه صنعتی شاهرود.

[۴] تعاونی، م. (۱۳۹۲)، پایان‌نامه کارشناسی ارشد: ”استفاده از مدل‌های مختلف رگرسیونی در تحلیل داده‌های طولی“، دانشکده علوم، دانشگاه گلستان.

[۵] حسن‌زاده بشتیان، م. (۱۳۸۸)، پایان‌نامه کارشناسی ارشد: ”برآورد انقباضی در مدل رگرسیون ریج با خطاهای دارای توزیع بیضی‌گون“، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد.

[۶] روزبه، م. (۱۳۹۰)، رساله دکتری: ”برآورد در مدل‌های خطی جزئی“، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد.

[۷] نوروزی‌راد، م. (۱۳۹۶)، رساله دکتری: ”برآوردگرهای بهبودیافته در برخی مدل‌های رگرسیونی خطی جریمه‌شده“، دانشکده علوم ریاضی، دانشگاه صنعتی شاهرود.

[8] Arashi, M., Saleh, A.K.Md.E., and Tabatabaey, S.M.M. (2013), “Regression model with elliptically contoured errors”, *Statistics*, 47, pp 1266-1284.

[9] Bancroft, T. A. (1944), “On biases in estimation due to the use of preliminary tests of significance”, *Annal. Math. Statist*, 15, pp 195-204.

[10] Chu, K. C. (1973), “Estimation and decision for linear systems with elliptically random process”, *IEEE. Tran. Aut. Cont.*, 18, pp 499-505.

- 
- [11] Cravan, P. and Wahba, G. (1979), "Smoothing noisy data with spline function: Estimating the correct degree of smoothing by the method of generalized cross validation", *Numerische Mathematik*, 31, pp 377-403.
- [12] Diggle, P. J., Heagerty, P., Liang, K. Y. and Zeger, S. L. (2002), "*Analysis of Longitudinal Data*", Oxford University Press, New York..
- [13] Eliot, M., Ferguson, J., Reilly, M. P and Foulkes, A. S. (2011), "Ridge regression for longitudinal biomarker data", *Biostatistics*, 7, pp 1-11.
- [14] Fang, K. T., Kotz, S. and Ng, K. W. (1990), "*Symmetric Multivariate and Related Distributions*", Chapman and Hall, London, New York.
- [15] Gentle, J., Hardle, W. and Mori, Y. (2004): "*Handbook of Computational Statistics*", Springer-Verlag Berlin Heidelberg.
- [16] Gibbons, D. G. (1981), "A simulation study of some ridge estimators", *J. Amer. Statist. Assoc*, 76, pp 1131-1150.
- [17] Gupta, A. K. and Varga, T. (1993), "*Elliptically Contored Models in Statistics*", Kluwer Academic Press, London.
- [18] Hoerl, A. E. and Kennard, R. W. (1970), "Ridge regression biased estimation for nonorthogonal problems", *Technometrics*, 12, pp 69-89.
- [19] James, W. and Stein, C. (1961), "Estimation with quadratic loss", *Proc. of Fouth Berkeley symp. Math. Statist. Prob.*, 1, pp 361-379.
- [20] Laird, N. and Ware, J. (1982), "Random-effects models for longitudinal data", *Biometrics*, 38, pp 963-974.
- [21] Liu, K. J. (1993), "A new class of biased estimate in linear regression", *Comm. Statist. Theo. Meth.*, 22, pp 393-402.
- [22] McCulloch, C. E. and Searle, S. (2001), "*Generalized, Linear, and Mixed Model*", John Wiley, Canada.
- [23] McCulloch, C. E., Searle, R. S. and Neuhaus, M. J. (2008), "*Generalized, Linear, and Mixed Models*", John Wiley, Inc., Hoboken, New Jersey, 2nd ed.
- [24] McDonald, G. C., Galarnaeu D. I. and Keimel K. (1975), "A Monte Carlo evaluation of some ridge-type estimators" *Amer. Statist. Assoc*, 20, pp 407- 416.

- [25] Montgomery, D. C., Peck, E. A., and Vining, G. G. (2012), "*Introduction to Linear Regression Analysis*", 5th Ed., John Wiley, New Jersey.
- [26] Muirhead, R. j. (1982), "*Aspect of Multivariate Statistical Theory*", John Wiley, New York.
- [27] Roozbeh, M., Arashi, M., and Kibria, B. M. G. (2015), "Generalized ridge regression estimator in semiparametric regression models", *J. Irania. Statist. Soc. (JIRSS)*, 14 (1), 25-62.
- [28] Saleh, A. K. Md. Ehsanes, (2006), "*Theory of Preliminary Test and Stein-Type Estimation with Application*", John Wiley, New Jersey.
- [29] Segerstedt, B. (1992), "On ordinary ridge regression in generalized linear models", *Comm. Statist. Theo. Meth.*, 21(8), pp 2227-2246.





# پیوست آ

## گزیده‌های از برنامه‌های کامپیوتری

### کدهای مربوط به نرم افزار R

رگرسیون خطی و آمیخته خطی در داده‌های *Orthodont*

```
> library(nlme)
> head(Orthodont, 8)

> Linear_Model <- lm(distance ~ age * Sex, data=Orthodont)
> summary(Linear_Model)

> Linear_Mixed_Model <- lme(distance ~ age * Sex, random = ~1|Subject
, data=Orthodont)
> summary(Linear_Mixed_Model)
```

تولید داده‌های همخط به روش پیشنهادی مک دونالد و همکاران (۱۹۷۵) و گیونز (۱۹۸۱)

P

n

```
rho=0.9

generator=function(p,n,rho){
  z=matrix(c(rnorm(p*n,mean=5,sd=1)),ncol=p)
  x=matrix(c(rep(0,n*p)),ncol=p)
  for(i in 1:n){
    for(j in 1:p){
      x[i,j]=(1-rho^2)^(1/2)*z[i,j]+rho*z[i,p]
    }
  }
  return(x)
}
data=generator(p,n,rho)
x=data
```

تولید داده‌های همخط با استفاده از بسته *EnvStats*

```
#install.packages("EnvStats")
library(EnvStats)
cor.mat=matrix(rho,5,5)
diag(cor.mat)=1
X <- as.matrix(simulateMvMatrix(ni*n,
  distributions = c("norm","norm","norm","norm","norm"),
  param.list = list( list(mean=5, sd=1),
                    list(mean=5, sd=1),
                    list(mean=5, sd=1),
                    list(mean=5, sd=1),
                    list(mean=5, sd=1)),
  cor.mat = cor.mat))
```

ساختار ماتریس Z

```
###Z
Zi=array(NA,c(q,ni,n))
for(i in 1:n){
  Zi[, ,i]<-matrix(rnorm(ni*q),q,ni)
```

```

}

#library(magic)
Z<-t(Zi[, ,1])
for(i in 2:n){
Z=adiag(Z,t(Zi[, ,i]))
}

```

ساختار بردار اثرات تصادفی  $b$

```

#library(tensr)
D=diag(rep(0.6,q))
bi=array(NA,c(q,1,n))
for(i in 1:n){
bi[, ,i]<-tensr::rmvnorm(1,rep(0,q),D)
}
b=as.vector(bi)

```

ساختار بردار خطای تصادفی  $\epsilon$

```

###E
sigma2=1
Ei=array(NA,c(ni,1,n))
for(i in 1:n){
Ei[, ,i]<-tensr::rmvnorm(1,rep(0,ni),sigma2*diag(ni))
}
E=as.vector(Ei)

```

محاسبه بردار پاسخ  $y$  با تعریف بردار  $\beta$

```

Y=X%*%B+Z%*%b+E
#Y500=Xi[, ,500]%*%B+t(Zi[, ,500])%*%bi[, ,500]+Ei[, ,500]

```

ساختار ماتریس واریانس  $V$

```

###V=var(Y)
U<-D

```

```
for(i in 2:n){
U=adiag(U,D)
}
D=U
V=Z%*%D%*%t(Z)+sigma2*diag(N)
```

بردار ضرایب  $\beta$  در مدل آمیخته خطی

```
###B_Mixed
B_Mixed=solve(t(X)%*%solve(V)%*%X)%*%t(X)%*%solve(V)%*%Y
```

بردار ضرایب  $\beta$  در مدل آمیخته خطی ریج

```
lambda0=seq(0.1,200,0.1)
arg0=numeric(length(lambda0))
for(i0 in 1:length(lambda0)){
BMR0=solve(t(X)%*%solve(V)%*%X+lambda0[i0]*diag(p))%*%t(X)%*%solve(V)%*%Y
Yhat0=X%*%BMR0+Z%*%b

S0=X%*%solve(t(X)%*%solve(V)%*%X+diag(lambda0[i0],p))%*%t(X)
%*%solve(V)+Z%*%D%*%t(Z)%*%(diag(N)X%*%(solve(t(X)
%*%solve(V)%*%X+diag(lambda0[i0],p))%*%t(X)%*%solve(V)))

trS0=sum(diag(S0))
arg0[i0]=(1/n)*t(Y-Yhat0)%*(Y-Yhat0)+2*(n^(-2))*trS0*t(Y-Yhat0)
%*(Y-Yhat0)
}
found0=cbind(arg0,lambda0)
lambda.min0=found0[found0[,1]==min(arg0),2]
#lambda.min0
BMR=solve(t(X)%*%solve(V)%*%X+lambda.min0*diag(p))%*%t(X)%*%solve(V)%*%Y
```

بردار ضرایب  $\beta$  در مدل آمیخته خطی لیو

```
d=seq(0.01,0.99,0.01)
arg0=numeric(length(d))
```

---

```

for(i0 in 1:length(d)){
  B_MLiu=solve(t(X)%*%solve(V)%*%X+diag(p))%*%(t(X)%*%solve(V)%*%(X)
+diag(d[i0],p))%*%B_Mixed

  Yhat0=X%*%B_MLiu+Z%*%b

  S0=X%*%solve(t(X)%*%solve(V)%*%X+diag(p))%*%(t(X)%*%solve(V)%*%X
+diag(d[i0],p))%*%solve(t(X)%*%solve(V)%*%X)%*%t(X)%*%solve(V)
+Z%*%D%*%t(Z)%*(diag(N)X%*%solve(t(X)%*%solve(V)%*%X
+diag(p))%*%(t(X)%*%solve(V)%*%X+diag(d[i0],p))%*%solve(t(X)
%*%solve(V)%*%X)%*%t(X)%*%solve(V))

  trS0=sum(diag(S0))
  arg0[i0]=(1/n)*t(Y-Yhat0)%*(Y-Yhat0)+2*(n^(-2))*trS0*t(Y-Yhat0)
%*(Y-Yhat0)
}
found0=cbind(arg0,d)
d.min=found0[found0[,1]==min(arg0),2]
#d.min=0.3
print(d.min)
B_MLiu=solve(t(X)%*%solve(V)%*%X+diag(p))%*%(t(X)%*%solve(V)%*%(X)
+(d.min*diag(p)))%*%B_Mixed

```

برآورد پارامترهای ماتریس کواریانس  $V$  در مدل آمیخته خطی

```

##go on start for convergence algorithm
iter=10000
B_EM_Mixed=matrix(NA,p,iter)
for(k in 1:iter){
  Yhat_Mixed=X%*%B_Mixed+Z%*%b
  e_Mixed=Y-Yhat_Mixed

  ###EM-algorithm Estimate D & sigma2
  ####sigma2 in EM_Algorithm for Mixed
  a=0

```

```

for(i in 1:n){
o=e_Mixed[(((i-1)*ni)+1):(i*ni),]
bb=t(o)%*%o
a=a+bb
}
t1_Mixed=sum(a)/n
sigma2_Mixed=t1_Mixed/N

####D in EM_Algorithm
aa=matrix((0),q,q)
for(i in 1:n){
O=b[(((i-1)*q)+1):(i*q)]
bb=O%*%t(O)
aa=aa+bb
}
t2=(1/n)*aa
D=t2/n

####Now update Sigma2 & D for E & b
###Update b for Mixed
bi=array(NA,c(q,1,n))
for(i in 1:n){
bi[,i]<-tenstr:::rmvnorm(1,rep(0,q),D)
}
b=as.vector(bi)

###Update E for Mixed
Ei=array(NA,c(ni,1,n))
for(i in 1:n){
Ei[,i]<-tenstr:::rmvnorm(1,rep(0,ni),sigma2_Mixed*diag(ni))
}
E=as.vector(Ei)

#####update Y

```

```

Y=X**B_Mixed+Z**b+E
### product D then Update V
U<-D
for(i in 2:n){
U=adiag(U,D)
}
D=U
V_Mixed=Z**D**t(Z)+sigma2_Mixed*diag(N)

### update B_Mixed
B_Mixed=solve(t(X)**solve(V_Mixed)**X)**t(X)**solve(V_Mixed)**Y
B_EM_Mixed[,k]=B_Mixed
}
B_EM_Mixed

```

برآورد پارامترهای ماتریس کواریانس  $V$  در مدل  $MR$  و محاسبه انحراف استاندارد

```

###SD.....
iter_Sd=500
B_MR_Sd_matrix=matrix((NA),p,iter_Sd)
B_MR_mean_matrix=matrix((NA),p,iter_Sd)
for(s in 1: iter_Sd){
print(s)
D=diag(rep(0.6,q))
bi=array(NA,c(q,1,n))
for(i in 1:n){
bi[, ,i]<-tensr::rmvnorm(1,rep(0,q),D)
}
b=as.vector(bi)

###E
sigma2=1
Ei=array(NA,c(ni,1,n))
for(i in 1:n){
Ei[, ,i]<-tensr::rmvnorm(1,rep(0,ni),sigma2*diag(ni))
}

```



```

}
E=as.vector(Ei)

###Mixed.Model & Y
Y=X%*%B+Z%*%b+E
dim(b)
#Y500=Xi[, , 500]%*%B+t(Zi[, , 500])%*%bi[, , 500]+Ei[, , 500]

###V=var(Y)
U<-D
for(i in 2:n){
U=adiag(U,D)
}
D=U
V=Z%*%D%*%t(Z)+sigma2*diag(N)

###B_Mixed0
lambda0=seq(0,200,1)
arg0=numeric(length(lambda0))
for(i0 in 1:length(lambda0)){
BMR0=solve(t(X)%*%solve(V)%*%X+lambda0[i0]*diag(p))%*%t(X)%*%solve(V)%*%Y
Yhat0=X%*%BMR0+Z%*%b

S0=X%*%solve(t(X)%*%solve(V)%*%X+diag(lambda0[i0],p))%*%t(X)
%*%solve(V)+Z%*%D%*%t(Z)%*%(diag(N)X%*%(solve(t(X)
%*%solve(V)%*%X+diag(lambda0[i0],p))%*%t(X)%*%solve(V)))

trS0=sum(diag(S0))
arg0[i0]=(1/n)*t(Y-Yhat0)%*%(Y-Yhat0)+2*(n^(-2))*trS0*t(Y-Yhat0)
%*%(Y-Yhat0)
}
found0=cbind(arg0,lambda0)
lambda.min0=found0[found0[,1]==min(arg0),2]
#lambda.min0

```

```

BMR=solve(t(X)%*%solve(V)%*%X+lambda.min0*diag(p))%*%t(X)%*%solve(V)%*%Y

##go on start for convergence algorithm
iter=10000
what.lam=matrix(NA,iter_Sd,iter)
B_EM_MR=matrix(NA,p,iter)
for(k in 1:iter){
Yhat_MR=X%*%BMR+Z%*%b
e_MR=Y-Yhat_MR

###EM-algorithm Estimate D & sigma2
####sigma2 in EM_Algorithm for Mixed
a=0
for(i in 1:n){
o=e_MR[(((i-1)*ni)+1):(i*ni),]
bb=t(o)%*%o
a=a+bb
}
t1_MR=sum(a)/n
sigma2_MR=t1_MR/N

####D in EM_Algorithm
aa=matrix((0),q,q)
for(i in 1:n){
O=b[(((i-1)*q)+1):(i*q)]
bb=O%*%t(O)
aa=aa+bb
}
t2=(1/n)*aa
D=t2/n

####Now update Sigma2 & D for E & b
###Update b for Mixed

```

```

bi=array(NA,c(q,1,n))
for(i in 1:n){
bi[, ,i]<-tenstr:::rmvnorm(1,rep(0,q),D)
}
b=as.vector(bi)

### product D then Update V
U<-D
for(i in 2:n){
U=adiag(U,D)
}
D=U
V_MR=Z%*%D%*%t(Z)+sigma2_MR*diag(N)

### update B_Mixed
lambda1=seq(0,200,1)
arg1=numeric(length(lambda1))
for(i1 in 1:length(lambda1)){

BMR1=solve(t(X)%*%solve(V_MR)%*%X+lambda1[i1]*diag(p))
%*%t(X)%*%solve(V_MR)%*%Y

Yhat1=X%*%BMR1+Z%*%b

S1=X%*%solve(t(X)%*%solve(V_MR)%*%X+diag(lambda1[i1],p))
%*%t(X)%*%solve(V_MR)+Z%*%D%*%t(Z)%*%(diag(N)X%*(solve(t(X)
%*%solve(V_MR)%*%X+diag(lambda1[i1],p))%*%t(X)%*%solve(V_MR)))

trS1=sum(diag(S1))
arg1[i1]=(1/n)*t(Y-Yhat1)%*(Y-Yhat1)+2*(n^(-2))*trS1*t(Y-Yhat1)
%*(Y-Yhat1)
}
found1=cbind(arg1,lambda1)

```

---

```

lambda.min1=found1[found1[,1]==min(arg1),2]
#lambda.min1

BMR=solve(t(X)%*%solve(V_MR)%*%X+lambda.min1*diag(p))%*%t(X)
%*%solve(V_MR)%*%Y

B_EM_MR[,k]=BMR
}

B_MR_mean_matrix[,s]=B_EM_MR[,iter]
B_MR_Sd_matrix[,s]=((B_EM_MR[,iter]-B)^2)
}
mean_B_MR=t(t(apply(B_MR_mean_matrix,1,mean)))
Sd_B_MR=t(t(sqrt(apply(B_MR_Sd_matrix,1,mean))))

final_list=list(cbind(rho,B,mean_B_MR,Sd_B_MR))
names(final_list)=c("MIXED_RIDGE")
final_list
#write.table(final_list,"E:/finaltableMixed&MRrho0.0itersd500n50
.txt",sep="\t")

```

بردار خطاهای دارای توزیع بیضی گون (مثال:  $t$  - student چند متغیره)

```

library(mvtnorm)
D=diag(rep(0.3,q))
bi=array(NA,c(q,1,n))
for(i in 1:n){
bi[, ,i]<-rmvt(1, D, df = 5)
}
b=as.vector(bi)

###E
sigma2=1
Ei=array(NA,c(ni,1,n))
for(i in 1:n){

```

```
Ei[, ,i]<-rmvt(1, sigma2*diag(ni), df = 5)
}
E=as.vector(Ei)
```

محاسبه  $\lambda$  به روش هورل و کنارد (۱۹۷۰) – تجزیه طیفی ماتریس  $(X^T X)$

```
t1_Mixed=a
sigma2_Mixed=as.numeric(t1_Mixed/N)
XTX=t(X)%*%X
GAMMA=eigen(XTX)$vectors
LAMBDA=diag(eigen(XTX)$values)
GAMMA%*%LAMBDA%*%t(GAMMA)
ALPHA=t(GAMMA)%*%B_Mixed
ALPHA2MAX=(max(ALPHA))^2
lambda0=sigma2_Mixed/ALPHA2MAX
```

برآوردگرهای *SSMR* و *PTMR*

```
#####Hypotheses
#H0: c (B_Mixed)=d
#in this problem; c=1 ; d=0
#Wald Statistic
W=drop(t(B_Mixed)%*%(t(X)%*%solve(V_Mixed)%*%X)%*(B_Mixed))/N

#####Stein
Ln_Mixed=W
Ln_MR=W
C=((N-p)*(p-2))/(p*(N-p+2))
B_Shrinkage_Mixed=B_Mixed-B_Mixed*(C/as.numeric(Ln_Mixed))
B_Shrinkage_MR=BMR-BMR*(C/as.numeric(Ln_MR))

#####PT
alpha=c(0.001,0.01,0.05,0.1,0.5)
Iln=c(NA)
for(cont in 1:length(alpha)){
Iln[cont]=0
```

```
if (Ln_MR<qchisq(alpha[cont],df=p,lower.tail=FALSE)) {Iln[cont]=1}
}
B_PTMR_alpha0.001=BMR-BMR*Iln[1]
B_PTMR_alpha0.01=BMR-BMR*Iln[2]
B_PTMR_alpha0.05=BMR-BMR*Iln[3]
B_PTMR_alpha0.1=BMR-BMR*Iln[4]
B_PTMR_alpha0.5=BMR-BMR*Iln[5]
```



# واژه‌نامه فارسی به انگلیسی

Preliminary test mixed ridge	آزمون اولیه آمیخته ریج
Wald statistic	آماره والد
Fixed effect	اثر ثابت
Random effect	اثر تصادفی
Interaction	اثر متقابل
Bias	اریبی
Independence	استقلال
Cross-Validation	اعتبار سنجی متقابل
Algorithm	الگوریتم
Standard deviation	انحراف استاندارد
Stein-type shrinkage	انقباضی نوع استاین
Shrinkage estimator	برآوردگر انقباضی
Liu estimator	برآوردگر لیو
Eigen vector	بردار ویژه
Multivariate Bessel	بسل چندمتغیره
Ridge parameter	پارامتر ریج
Scale parameter	پارامتر مقیاس
Multivariate Pearson	پیرسن چندمتغیره
Density generator function	تابع مولد چگالی
Characteristic generator function	تابع مولد مشخصه
Indicator function	تابع نشان‌گر
Exchangeable	تبادل پذیر
Elliptically contoured distribution	توزیع بیضی‌گون
Penalty	جریمه
Longitudinal data	داده‌های طولی
Determinant	دترمینان



Full Rank	رتبه کامل
Ridge regression	رگرسیون ریج
Regularity Conditions	شرایط نظم
Variance inflation factor	عامل تورم واریانس
Condition number	عدد شرطی
Least squares	کمترین توان‌های دوم
Iterative weighted least squares method	کمترین توان‌های دوم موزون تکراری
Multivariate Cauchy	کوشی چندمتغیره
Multivariate Laplace	لاپلاس چندمتغیره
Fisher information matrix	ماتریس اطلاع فیشر
Design matrix	ماتریس طرح
Diagonal matrix	ماتریس قطری
Maximum likelihood	ماکزیمم درست‌نمایی
Orthogonal	متعامد
Response variable	متغیر پاسخ
Independent Variable	متغیر مستقل
Transition model	مدل انتقال
Marginal model	مدل حاشیه‌ای
Linear mixed model	مدل خطی آمیخته
Generalized linear model	مدل خطی تعمیم یافته
Mean squares error	میانگین توان‌های دوم خطا
Mean prediction error	میانگین خطای پیش‌گویی
Eigen value	مقدار ویژه
Multivariate exponential power	نمایی توانی چندمتغیره
Correlation	همبستگی
Collinearity	همخطی

# واژه‌نامه انگلیسی به فارسی

Algorithm	الگوریتم
Bias	اریبی
Characteristic generator function	تابع مولد مشخصه
Collinearity	همخطی
Condition number	عدد شرطی
Cross-Validation	اعتبار سنجی متقابل
Correlation	همبستگی
Density generator function	تابع مولد چگالی
Design matrix	ماتریس طرح
Determinant	دترمینان
Diagonal matrix	ماتریس قطری
Eigen value	مقدار ویژه
Eigen vector	برداری ویژه
Elliptically contoured distribution	توزیع بیضی‌گون
Exchangeable	تبادل‌پذیر
Fisher information matrix	ماتریس اطلاع‌فیشر
Fixed effect	اثر ثابت
Full Rank	رتبه کامل
Generalized linear model	مدل خطی تعمیم‌یافته
Independence	استقلال
Independent Variable	متغیر مستقل
Indicator function	تابع نشان‌گر
Interaction	اثر متقابل
Iterative weighted least squares method	کمترین توان‌های دوم موزون تکراری
Least squares	کمترین توان‌های دوم
Linear mixed model	مدل آمیخته خطی

Liu estimator . . . . .	برآوردگر لیو
Longitudinal data . . . . .	داده‌های طولی
Marginal model . . . . .	مدل حاشیه‌ای
Maximum likelihood . . . . .	ماکزیمم درست‌نمایی
Mean prediction error . . . . .	میانگین خطای پیش‌گویی
Mean squares error . . . . .	میانگین توان‌های دوم خطا
Multivariate Bessel . . . . .	بسل چندمتغیره
Multivariate Cauchy . . . . .	کوشی چندمتغیره
Multivariate Exponential power . . . . .	نمایی توانی چندمتغیره
Multivariate Laplace . . . . .	لاپلاس چندمتغیره
Multivariate Pearson . . . . .	پیرسن چندمتغیره
Orthogonal . . . . .	متعامد
Penalty . . . . .	جریمه
Preliminary test mixed ridge . . . . .	آزمون اولیه آمیخته ریج
Random effect . . . . .	اثر تصادفی
Regularity Conditions . . . . .	شرایط نظم
Response variable . . . . .	متغیر پاسخ
Ridge parameter . . . . .	پارامتر ریج
Ridge regression . . . . .	رگرسیون ریج
Scale parameter . . . . .	پارامتر مقیاس
Shrinkage estimator . . . . .	برآوردگر انقباضی
Standard deviation . . . . .	انحراف استاندارد
Stein-type shrinkage . . . . .	انقباضی نوع استاین
Transition model . . . . .	مدل انتقال
Variance inflation factor . . . . .	عامل تورم واریانس
Wald statistic . . . . .	آماره والد

## **Abstract**

Longitudinal data are those observations in which measurements of each experiment units repeats over the time. Linear mixed model is one of the tools used to analyze such data. However, pre sense of multicollinearity causes to use shrinkage methods for analysis. In this dissertation, the combination of linear mixed elliptical models with shrinkage methods such as ridge and Liu is considered for analysis purposes. The performance of the proposed estimators is evaluated using simulation and real data analysis.

**Keywords:** EM algorithm, Elliptical distribution, Linear mixed models, Liu estimator, Ridge regression, Shrinkage estimator.



**Shahrood University of Technology**

**Faculty of Mathematical Sciences**

**MSc Thesis in: Mathematical Statistics**

**Likelihood-based Inference in Ridge  
Regression for Longitudinal Data**

**By: Mohammad Rahmani**

**Supervisor**

**Dr. Mohammad Arashi**

**September 2018**