

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ
مَنْ عَمِلْ سَئِئْرًا فَلْيَرْجُ
عَذَابَ اللَّهِ الْكَبِيرِ



دانشکده علوم ریاضی

پایان نامه کارشناسی ارشد آمار

تحلیل رگرسیونی داده‌های نرخ و نسبت با استفاده از توزیع سیمپلکس

نگارنده: سیما ملکی

استاد راهنما

دکتر داود شاهسونی

بهمن ۱۳۹۷

اگر شایسته تقدیم باشد، به

خانواده و دوستانم

که همواره در امر تحصیل مشوق من بوده اند.

سپاس گزاری ...

منت خدای راعز و جل که طاعتش موجب قربت است و به شکر اندرش مزید نعمت. هر نفسی که فرومی رود مدحیات است و چون برمی آید مفرح ذات. پس در هر نفسی دو نعمت موجود است و بر هر نعمتی شگری واجب.

سپاس خود را به محضر استاد راهنا، آقای دکتر شایسونی برای رهنمودهایشان تقدیم می دارم. از استاد فرزانه، آقای دکتر باغشینی و دانشجوی گرامی ایشان آقای میعادولی پور کمال قدردانی را دارم. از استاد کرامتقدر، آقای دکتر ربیعی که همواره در این راه دلگرمی بنده بودند و مرا به صبر و سکینایی دعوت نمودند؛ کمال تشکر را دارم.

کمال تشکر را دارم از خانم اسدی، خانم لاله زاده و آقای حسین زاده به پاس کمک های ارزنده شان که هم چون یک دوست و راهنمای بنده بودند. در پایان از همه دوستان که همواره باعث دلگرمی بودند، نهایت تشکر را دارم.

تعهد نامه

این جانب سیما ملکی دانشجوی کارشناسی ارشد رشته آمار دانشکده علوم ریاضی دانشگاه صنعتی شاهرود، نویسنده پایان نامه با عنوان **تحلیل رگرسیونی داده‌های نرخ و نسبت با استفاده از توزیع سیمپلکس**، تحت راهنمایی **داود شاهسونی** متعهد می‌شوم:

- تحقیقات در این پایان نامه توسط این جانب 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 به دست 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.

سیما ملکی

بهمن ۱۳۹۷

مالکیت نتایج و حق نشر

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

چکیده

در مسائل رگرسیون، هنگامی که متغیر پاسخ از نوع نرخ یا نسبت در دامنه تغییرات (۱, ۰) است، به طور معمول از مدل‌های رگرسیونی لجستیک یا پروبیت استفاده می‌شود. در حالتی که توزیع این داده‌ها چوله باشد، مدل‌های مذکور مناسب نبوده و نمی‌توانند به خوبی به داده‌ها برازش شوند، زیرا توزیع‌های زیر بنایی آن‌ها یعنی لجستیک و نرمال متقارن بوده و با ماهیت چوله بودن در تعارض است. در این پایان‌نامه، مدل رگرسیونی سیمپلکس، به‌عنوان یک راه حل جایگزین، معرفی شده است. ابتدا توزیع سیمپلکس، مبتنی بر توزیع گوسی معکوس آمیخته را مطرح نموده و نحوه شبیه‌سازی داده‌های سیمپلکس را بر اساس شبیه‌سازی گوسی معکوس آمیخته بیان می‌کنیم. از آن‌جا که مدل رگرسیونی سیمپلکس در قالب «مدل‌های پراکندگی» قابل طرح است، لذا مدل‌های پراکندگی را معرفی نموده و همانند روش مدل‌های خطی تعمیم‌یافته، پارامترهای مرتبط را برآورد می‌کنیم. چون توزیع سیمپلکس شامل دو پارامتر است، مسئله برآورد را در دو حالت ثابت و متغیر بودن پارامتر واریانس، بررسی کرده و از مثال‌های شبیه‌سازی متنوع برای درک مطلب بهره می‌جوییم. از رگرسیون سیمپلکس برای مدل‌سازی داده‌های طولی نیز استفاده کرده و با استفاده از دو مثال واقعی، کارایی و عملکرد رگرسیون سیمپلکس را با این نوع داده‌ها، بررسی می‌کنیم.

کلمات کلیدی: مدل‌های پراکندگی، داده‌های نرخ و نسبت، تولید متغیر تصادفی، مدل‌های خطی تعمیم‌یافته، توزیع سیمپلکس، مدل رگرسیونی سیمپلکس، داده‌های طولی، معادلات برآوردیابی تعمیم‌یافته.

لیست مقالات مستخرج از پایان نامه

۱. ملکی، س.، شاهسونی، د. (۱۳۹۷)، ”مدل بندی داده های نرخ و نسبت با استفاده از توزیع سیمپلکس”، اولین سمینار دانشجویی آمار ایران، دانشگاه تربیت مدرس.

فهرست مطالب

ق	فهرست تصاویر
ش	فهرست جداول
۱	۱ تعاریف و مفاهیم اولیه
۱	۱.۱ مقدمه
۴	۲.۱ مدل‌های پراکندگی
۷	۱.۲.۱ خواص مدل‌های پراکندگی
۹	۲.۲.۱ برآورد درستی‌نمایی ماکسیمم
۱۲	۳.۱ معیارهای اندازه‌گیری دقت مدل
۱۵	۲ توزیع سیمپلکس
۱۵	۱.۲ توزیع گوسی معکوس
۱۷	۱.۱.۲ تابع مولد گشتاور توزیع گوسی معکوس
۱۸	۲.۱.۲ وجه تسمیه معکوس
۱۹	۳.۱.۲ برخی خواص توزیع گوسی معکوس
۲۰	۴.۱.۲ تولید مقادیر تصادفی از توزیع گوسی معکوس
۲۲	۲.۲ توزیع گوسی معکوس آمیخته
۲۴	۱.۲.۲ تولید مقادیر تصادفی از توزیع گوسی معکوس آمیخته
۲۴	۳.۲ توزیع سیمپلکس
۲۷	۱.۳.۲ رابطه بین توزیع گوسی معکوس آمیخته و توزیع سیمپلکس
۲۸	۲.۳.۲ تولید مقادیر تصادفی از توزیع سیمپلکس
۳۱	۳.۳.۲ برخی خواص توزیع سیمپلکس
۳۳	۳ مدل رگرسیونی سیمپلکس
۳۳	۱.۳ مدل خطی تعمیم‌یافته سیمپلکس
۳۴	۱.۱.۳ برازش مدل

۳۷	آماره انحراف	۲.۱.۳
۳۹	برآورد پارامتر پراکندگی	۳.۱.۳
۴۰	تحلیل باقی مانده‌ها برای مدل رگرسیونی سیمپلکس	۴.۱.۳
۴۰	مطالعه شبیه‌سازی	۵.۱.۳
۴۴	مدل رگرسیونی سیمپلکس با پارامتر پراکندگی متغیر	۲.۳
۴۴	معرفی مدل	۱.۲.۳
۴۴	برازش مدل	۲.۲.۳
۴۶	مطالعه شبیه‌سازی	۳.۲.۳
۵۱		مدل حاشیه‌ای برای تحلیل داده‌های نسبت طولی	۴
۵۲	معرفی داده‌های طولی	۱.۴
۵۲	ساختار داده‌های طولی	۱.۱.۴
۵۳	اهداف و ضرورت تحلیل داده‌های طولی	۲.۱.۴
۵۴	روش معادلات برآوردیابی تعمیم‌یافته	۲.۴
۵۵	خانواده توزیع‌های نمایی	۱.۲.۴
۵۵	تابع شبه‌درست‌نمایی	۲.۲.۴
۵۷	مدل حاشیه‌ای سیمپلکس برای تحلیل داده‌های نسبت طولی	۳.۴
۵۷	مدل حاشیه‌ای	۱.۳.۴
۵۹	روش GEEs برای برآورد پارامترها	۲.۳.۴
		روش معادلات برآوردیابی تعمیم‌یافته مرتبه اول (GEE1) و مرتبه دوم (GEE2)	۳.۳.۴
۶۳		
۶۴	انواع معادلات برآوردیابی تعمیم‌یافته	۴.۳.۴
۶۵	تحلیل باقی مانده‌ها	۵.۳.۴
۶۷		کاربرد مدل رگرسیونی سیمپلکس و برازش آن با R	۵
۶۷	مقدمه	۱.۵
۷۰	کاربرد مدل در دو مثال واقعی	۲.۵
۷۰	داده‌های سلول بنیادی خون محیطی	۱.۲.۵
۷۴	داده‌های جراحی چشم	۲.۲.۵
۸۱	نتیجه‌گیری	۳.۵
۸۳			آ
۸۳	تولید عدد تصادفی از توزیع MIG	۱.آ
۸۴	ماتریس اطلاع گودامب	۲.آ
۸۵	دستورات نرم‌افزار R	۳.آ

فهرست تصاویر

۱۷ نمودار توابع چگالی گوسی معکوس	۱.۲
۲۶ نمودار تابع چگالی سیمپلکس	۲.۲
۲۷ نمودار توابع چگالی سیمپلکس (منحنی توپر) و نرمال (منحنی مقطع)	۳.۲
۲۹ نمودار چگالی (منحنی توپر) به همراه بافت‌نگار اعداد تصادفی از توزیع سیمپلکس	۴.۲
 نمودار چگالی (منحنی توپر) به همراه هسته چگالی (منحنی مقطع) اعداد تصادفی	۵.۲
۳۰ از توزیع سیمپلکس	
۳۱ نمودار چندک-چندک	۶.۲
 نمودار چندک-چندک برآوردهای پارامترهای مدل (۱۱.۲)، با پیشگوی خطی	۱.۳
۴۳ (۱۵.۳)	
 نمودار چندک-چندک برآوردهای ضرایب رگرسیونی میانگین و پارامتر پراکندگی	۲.۳
۴۹ در مدل (۲۴.۳) برای $n = 400$ و تعداد تکرار ۵۰۰	
۵۴ نمودار قدرت خواندن کودکان	۱.۴
 بررسی نیکویی برازش مدل‌های رگرسیون سیمپلکس با باقی‌مانده‌های پیرسون،	۱.۵
۷۴ در مثال PBSC	
۸۰ بررسی خودهمبستگی در چهار تاخیر مختلف، در مثال جراحی چشم	۲.۵
۸۱ بررسی تابع پیوند و نیکویی برازش مدل در مثال جراحی چشم	۳.۵

فهرست جداول

۶	تابع انحراف و تابع واریانس برخی مدل‌های پراکندگی	۱.۱
۱۰	برخی توابع پیوند و مشتق اول آن‌ها	۲.۱
۳۹	آماره انحراف مقیاس‌بندی‌شده برخی از توزیع‌های متعلق به خانواده نمایی	۱.۳
۴۱	خلاصه نتایج شبیه‌سازی برای برآورد ضرایب مدل (۱۱.۲) با پیشگوی خطی (۱۴.۳)	۲.۳
۴۲	خلاصه نتایج شبیه‌سازی برای برآورد ضرایب مدل (۱۱.۲) با پیشگوی خطی (۱۵.۳)	۳.۳
	خلاصه نتایج شبیه‌سازی برای برآورد ضرایب مدل (۱۱.۲) با پیشگوی خطی (۲۴.۳)	۴.۳
۴۸		
۵۲	ساختار داده‌ها در مطالعات طولی	۱.۴
۵۷	ضوابط تابع واریانس و شبه‌درست‌نمایی برای چند توزیع از خانواده توزیع‌های نمایی	۲.۴
۷۲	باقی‌مانده استاندارد پیرسون مدل (۱.۵)	۱.۵
۷۲	برآورد پارامترهای مدل (۱.۵)	۲.۵
۷۳	باقی‌مانده استاندارد پیرسون مدل (۲.۵)	۳.۵
۷۳	برآورد پارامترهای مدل (۲.۵)	۴.۵
۷۳	مقایسه AIC برای دو مدل	۵.۵
۷۳	مقایسه BIC برای دو مدل	۶.۵
۷۶	باقی‌مانده استاندارد پیرسون مدل (۳.۵)	۷.۵
۷۶	برآورد پارامترهای مدل (۳.۵)	۸.۵
۷۷	باقی‌مانده استاندارد پیرسون مدل (۴.۵)	۹.۵
۷۷	برآورد پارامترهای مدل (۴.۵)	۱۰.۵
۷۸	باقی‌مانده استاندارد پیرسون مدل (۵.۵)	۱۱.۵
۷۸	برآورد پارامترها با ساختار همبستگی Exc در مدل (۵.۵)	۱۲.۵
۷۹	باقی‌مانده استاندارد پیرسون مدل (۵.۵) با ساختار همبستگی AR(1)	۱۳.۵
۷۹	برآورد پارامترها با ساختار همبستگی AR(1) در مدل (۵.۵)	۱۴.۵

فصل ۱

تعاریف و مفاهیم اولیه

۱.۱ مقدمه

مدل‌سازی آماری، روابط بین متغیر پاسخ و متغیرهای تبیینی را بیان می‌کند که صورت کلی این مدل‌ها عبارت است از

مولفه خطای تصادفی + مولفه سیستماتیک = متغیر پاسخ

مولفه سیستماتیک حاوی تغییرات متغیر پاسخ برحسب متغیرهای تبیینی است و مولفه خطای تصادفی حاوی سهمی از تغییرات پاسخ است که توسط مولفه سیستماتیک بیان نشده است. فرمول‌بندی کلی یک مدل آماری برای n مقدار مشاهده‌شده متغیر پاسخ Y ، به صورت

$$y_i = \eta_i + \epsilon_i, \quad i = 1, \dots, n$$

بیان می‌شود که η_i همان مولفه سیستماتیک و ϵ_i مولفه خطاست. مولفه η_i می‌تواند نسبت به پارامترهای خود، صورت‌های خطی نظیر

$$\eta_i = \beta_0 + \beta_1 x_i$$

$$\eta_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2$$

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

داشته باشد یا این که نظیر

$$\eta_i = \beta_0 + \exp\{\beta_1 x_i\}$$

$$\eta_i = x_i / (\beta_0 + \beta_1 x_i)$$

نسبت به پارامترها غیرخطی باشد، که در آن x_{i1}, \dots, x_{ip} مقادیر متغیرهای تبیینی X_1, \dots, X_p هستند.

فرض کنید Y_i به صورت خطی به مقادیر این p متغیر تبیینی X_1, \dots, X_p وابسته باشد. آن گاه با فرض این که

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} = \mathbf{x}_i^T \boldsymbol{\beta}$$

مدل مفروض را یک مدل خطی کلی^۱، $\eta_i = \sum_{j=0}^p \beta_j x_{ij}$ را «پیشگوی خطی»^۲ نامند. این مدل طیف وسیعی از مباحث آماری را به خود اختصاص داده به طوری که در وضعیت متداول هم خطی متغیرهای تبیینی، مباحثی از قبیل رگرسیون جریمه شده، رگرسیون مولفه‌های اصلی، رگرسیون کمترین توان‌های دوم جزئی خودنمایی می‌کنند.

در مدل‌های خطی، پذیره‌هایی حاکم هستند که موجب تسهیل در استنباط آماری می‌شوند. اغلب فرض می‌شود که ϵ_i ها متغیرهای تصادفی ناهمبسته با میانگین صفر و واریانس مشترک σ^2 هستند. یعنی

$$\mu_i = E(Y_i) = \eta_i, \quad V(Y_i) = \sigma^2.$$

علاوه بر این، با اعمال پذیره توزیعی نرمال برای ϵ ها، $\epsilon_i \sim N(0, \sigma^2)$ ، داریم

$$Y_i \sim N(\eta_i, \sigma^2).$$

موارد فوق را به صورت زیر خلاصه می‌کنیم

۱. Y_i ها متغیرهای تصادفی مستقل با میانگین μ_i هستند، $E(Y_i) = \mu_i$ ، و همچنین $Y_i \sim N(\mu_i, \sigma^2)$.

۲. تاثیر p متغیر تبیینی بر متغیر پاسخ، با مولفه سیستماتیک خطی مدل‌سازی می‌شود، یعنی $\eta_i = \sum_{j=0}^p \beta_j x_{ij}$.

۳. ارتباط بین ۱ و ۲ عبارت است از $\mu_i = \eta_i$.

صور سه‌گانه بالا را می‌توان به صورت زیر تعمیم داد

۱. توزیع Y_i لزوماً نرمال نیست.

¹ General Linear Model

² Linear Predictor

۲. μ_i می‌تواند تابعی از η_i باشد، یعنی $\mu_i = g(\eta_i)$.

۳. مولفه سیستماتیک می‌تواند غیرخطی باشد.

این مفاهیم مبانی «مدل‌های خطی تعمیم‌یافته»^۳ (GLM) (مک کالا و نلدر، ۱۹۸۹) را شکل می‌دهند. در واقع مدل‌های خطی تعمیم‌یافته، تعمیمی از مدل‌های خطی کلی هستند. نظریه مدل‌های خطی تعمیم‌یافته، حاکی از آن است که تحلیل رگرسیون نیازمند تشخیص نوع متغیر پاسخ است. علیرغم معروف بودن و پرکاربرد بودن توزیع نرمال، جرجنسن (۱۹۹۷) به این نکته اشاره داشت که به‌جز برای داده‌های با پراکندگی کم، توزیع نرمال یک استثنا است تا این که یک قانون باشد. فیشر (۱۹۵۳)، آماردان برجسته، نیز به ما یادآوری کرد که توصیف داده‌ها باید با توجه به ذات و ماهیت طبیعی آن‌ها صورت گیرد. نلدر و ودربرن (۱۹۷۲) اولین کسانی بودند که توانستند با معرفی رده GLM نشان دهند که طیف وسیعی از داده‌های غیر نرمال را می‌توان با یک چارچوب واحد تحلیل کرد. در واقع GLM برای خانواده توزیع‌های نمایی بسط داده شده، یعنی

۱. متغیر تصادفی Y_i دارای توزیعی در خانواده توزیع‌های نمایی است.

۲. متغیرهای تصادفی از طریق مولفه سیستماتیک وارد مدل می‌شوند، $\eta_i = \sum_{j=0}^p \beta_j x_{ij}$.

۳. میانگین Y_i یعنی μ_i ، از طریق تابعی به نام «تابع پیوند»^۴ g به η_i مرتبط می‌شود، یعنی $\eta_i = g(\mu_i)$.

هنگامی که متغیر پاسخ از جنس نرخ و نسبت است، مدل‌های معدودی برای برازش چنین داده‌هایی معرفی شده و سودمند هستند. در توزیع بتا که اغلب در آمار بیزی مورد استفاده قرار می‌گیرد، می‌توان با تبدیل مناسبی از پارامترها، آن را برای مدل‌سازی چنین متغیرهای پاسخی به کار برد. فراری و سریباری (۲۰۰۴) مدل رگرسیونی بتا را برای پاسخ‌های نرخ و نسبت به کار بردند و روش‌های برآورد و استنباط آماری مرتبط، در نرم‌افزار R توسط بسته betareg عملیاتی شد (سریباری و همکاران، ۲۰۱۰).

با این وجود، وقتی پارامتر پراکندگی به مدل‌سازی اضافه می‌شود، مدل GLM بر مبنای توزیع سیمپلکس^۵ (بارندروف، نیلسن و جرجنسن، ۱۹۹۱) برای تحلیل داده‌های پیوسته از جنس نسبت، نیرومندتر شناخته شد (ژنگ و ژو، ۲۰۱۴). توزیع سیمپلکس شامل رده وسیعی از توزیع‌هایی است که دامنه آن‌ها در بازه (۰، ۱) است. در واقع توزیع سیمپلکس یک مدل پراکندگی است و خواص تحلیلی مشترکی با مدل‌های پراکندگی نمایی دارد.

³ Generalized Linear Models

⁴ Link Function

⁵ Simplex Distribution

در بحث تحلیل داده‌های طولی نیز، توزیع سیمپلکس نقش موثری را ایفا کرده است. مدل‌های حاشیه‌ای سیمپلکس^۶ (سانگ و تن، ۲۰۰۰) و مدل‌های اثرات آمیخته سیمپلکس (ژو و همکاران، ۲۰۰۸) نیز از جمله این موارد هستند.

در این پایان‌نامه، با معرفی توزیع سیمپلکس در فصل ۲ آغاز کرده و رابطه آن با توزیع گوسی معکوس آمیخته^۷ و نحوه شبیه‌سازی اعداد تصادفی سیمپلکس را شرح می‌دهیم. سپس مدل خطی تعمیم‌یافته سیمپلکس را در فصل ۳ مورد بررسی قرار می‌دهیم. در فصل ۴ به معرفی داده‌های طولی^۸ و روش معادلات برآوریابی تعمیم‌یافته^۹ (GEEs) می‌پردازیم. سرانجام در فصل ۵ به برازش و تحلیل داده‌ها در مدل‌های رگرسیونی سیمپلکس با استفاده از بسته simplexreg موجود در نرم‌افزار R، می‌پردازیم و سپس کاربرد مدل را با دو مثال واقعی نمایش می‌دهیم. برای ورود به فصل ۳، نیازمند معرفی اجمالی مدل‌های پراکندگی هستیم که در بخش ۲.۱ به آن می‌پردازیم. لازم به ذکر است که فصل ۲، کاملاً مستقل از مفهوم رگرسیون و مدل رگرسیونی سیمپلکس است و در صورتی که خواننده صرفاً به بحث رگرسیونی آن علاقه‌مند باشد می‌تواند فقط با دانستن چگالی سیمپلکس در بخش ۳.۲، بحث رگرسیون سیمپلکس را از فصل ۳ دنبال کند.

۲.۱ مدل‌های پراکندگی

ایده اصلی حاکم بر GLM، توسط جرجسن (۱۹۹۷) به رده بزرگ‌تری از مدل‌ها که «مدل‌های پراکندگی»^{۱۰} (DM) نامیده می‌شوند، تعمیم داده شد. در این مدل‌ها عقیده کلی در مورد «مکان» و «مقیاس» به «موقعیت»^{۱۱} و «پراکندگی» تعمیم می‌یابد. در این صورت طیف وسیعی از داده‌های غیر نرمال، پوشش داده می‌شود و گونه‌هایی نظیر داده‌های مثبت، مثبت و صفر، داده‌های شمارشی، داده‌های دودویی و داده‌های جهتی^{۱۲} می‌توانند توسط مدل‌های پراکندگی و با حفظ ماهیت ذاتی آن‌ها مورد تحلیل و مدل‌سازی قرار گیرند؛ در واقع رده مهمی از مدل‌های پراکندگی، مدل‌های خطی تعمیم‌یافته است که در آن نقش مجموع توان‌های دوم، به‌عنوان مبنا، به مفهوم دیگری به نام «انحراف»^{۱۳} که ابزاری مهم برای برازش و انتخاب مدل است، تعمیم داده می‌شود.

برای تشریح مدل‌های پراکندگی، توزیع نرمال را که نقش اساسی در رگرسیون خطی کلاسیک

⁶ Simplex Marginal Models

⁷ Mixture Inverse Gaussian Distribution

⁸ Longitudinal Data

⁹ Generalized Estimating Equations

¹⁰ Dispersion Models

¹¹ Position

¹² Directional

¹³ Deviance

ایفا می‌کند، در نظر بگیرید. در چگالی $N(\mu, \sigma^2)$ یعنی

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(y-\mu)^2}{2\sigma^2}\right\}$$

جمله $(y-\mu)^2$ یک فاصله اقلیدسی است که تفاوت^{۱۴} بین مشاهده y و میانگین μ را اندازه‌گیری می‌کند. این معیار اندازه‌گیری تفاوت، در تحلیل رگرسیون، برای تشکیل آماره F و دستیابی به نیکویی برازش استفاده می‌شود.

با الگوپذیری از چگالی نرمال، جرجنسن (۱۹۹۲) مدل‌های پراکندگی را ارائه داد که در آن فاصله اقلیدسی $(y-\mu)^2$ به صورت کلی‌تری از مفهوم تفاوت و با نماد $d(y; \mu)$ در نظر گرفته می‌شود.

تعریف ۱.۲.۱. مدل پراکندگی $DM(\mu, \sigma^2)$ با پارامتر مکان μ و پارامتر پراکندگی σ^2 ، خانواده‌ای از توزیع‌هاست که دارای تابع چگالی به صورت زیر است:

$$f(y; \mu, \sigma^2) = a(y; \sigma^2) \exp\left\{\frac{-1}{2\sigma^2} d(y; \mu)\right\} \quad y \in \mathcal{C} \quad (1.1)$$

که در آن $\mu \in \Omega$ ، $\sigma^2 > 0$ و ثابت نرمال‌ساز $a \geq 0$ و مستقل از μ است و عموماً $\Omega \subseteq \mathcal{C} \subseteq \mathbb{R}$. از آن‌جا که ثابت نرمال‌ساز شامل μ نیست، لذا برآورد μ (یا β در GLM‌ها) به‌طور مجزا از برآورد σ^2 انجام می‌شود که سهولت برآورد پارامترها را در بر خواهد داشت. این خاصیت بسیار مفید که به‌عنوان متعامدسازی درست‌نمایی^{۱۵} شناخته می‌شود، در توزیع نرمال برقرار بوده و در مدل‌های پراکندگی نیز پا برجاست.

تابع دو متغیره $d(\cdot, \cdot)$ را تابع انحراف^{۱۶} نامیده که $(y, \mu) \in \mathcal{C} \times \Omega$ و دارای دو شرط زیر است:

۱. وقتی که مشاهده y و μ مورد انتظار باهم برابر باشند، $y = \mu$ ، آن‌گاه

$$d(y; y) = 0, \quad \forall y \in \Omega.$$

۲. وقتی که مشاهده y و μ مورد انتظار باهم برابر نباشند، آن‌گاه

$$d(y; \mu) > 0, \quad \forall y \neq \mu.$$

تابع انحراف، منظم^{۱۷} نامیده می‌شود اگر مشتق‌های مرتبه دوم تابع $d(y; \mu)$ نسبت به (y, μ) با هم برابر باشند

$$\frac{\partial^2 d}{\partial y^2}(y; y) = \frac{\partial^2 d}{\partial \mu^2}(y; \mu) \Big|_{\mu=y} > 0, \quad \forall y \in \Omega.$$

¹⁴ Discrepancy

¹⁵ Likelihood Orthogonality

¹⁶ Deviance

¹⁷ Regular

برای تابع انحراف منظم، تابع واریانس^{۱۸} $V : \Omega \rightarrow (0, \infty)$ به شرح زیر تعریف می‌شود

$$V(\mu) = \frac{2}{\frac{\partial^2 d}{\partial \mu^2}(y; \mu)|_{y=\mu}}, \quad \mu \in \Omega. \quad (2.1)$$

برخی از مدل‌های پراکندگی معروف به همراه تابع انحراف d و تابع واریانس V متناظرشان در جدول ۱.۱ معرفی شده‌اند.

جدول ۱.۱: تابع انحراف و تابع واریانس برخی مدل‌های پراکندگی

توزیع	d	\mathcal{C}	Ω	$V(\mu)$
نرمال	$(y - \mu)^2$	$(-\infty, \infty)$	$(-\infty, \infty)$	۱
پواسون	$2(y \log \frac{y}{\mu} - y + \mu)$	$\{0, 1, \dots\}$	$(0, \infty)$	μ
گوسی معکوس	$\frac{(y - \mu)^2}{y\mu^2}$	$(0, \infty)$	$(0, \infty)$	μ^3
سیمپلکس	$\frac{(y - \mu)^2}{y(1-y)\mu^2(1-\mu)^2}$	$(0, 1)$	$(0, 1)$	$\mu^3(1-\mu)^3$

در دو مثال زیر، چگونگی قرار گرفتن توزیع‌های نرمال و پواسون را در قالب مدل‌های پراکندگی بررسی می‌کنیم.

مثال ۱.۲.۱. توزیع نرمال $N(\mu, \sigma^2)$ را در نظر بگیرید که در آن $d(y; \mu) = (y - \mu)^2$ و $y \in \mathcal{C} = \mathbb{R}$ و هم‌چنین $\mu \in \Omega = \mathbb{R}$ واضح است که تابع d نامنفی است و در صورتی که $y = \mu$ ، برابر با صفر است. در نتیجه مشتق‌های مرتبه اول و دوم d به صورت زیر است:

$$\frac{\partial d}{\partial \mu} = -2(y - \mu), \quad \frac{\partial^2 d}{\partial \mu^2} = 2$$

سپس تابع واریانس را با توجه به رابطه (۲.۱) می‌توان به دست آورد

$$V(\mu) = \frac{2}{2} = 1$$

مثال ۲.۲.۱. توزیع پواسون با پارامتر μ را در نظر بگیرید

$$f(y; \mu) = \frac{\mu^y}{y!} e^{-\mu}, \quad y \in \{0, 1, \dots\}; \mu \in \Omega = (0, \infty),$$

که معادل است با

$$f(y; \mu) = \frac{1}{y!} \exp \{y \log \mu - \mu\}.$$

در این جا $\{y \log \mu - \mu\}$ یک تابع انحراف نیست؛ زیرا با توجه به شرایطی که d باید دارا باشد، وقتی $y = \mu$ داریم $d \neq 0$. حال رابطه فوق را در عبارت $\exp\{y \log y - y\}$ ضرب و تقسیم می‌کنیم؛ در نتیجه داریم

$$p(y; \mu) = \left\{ \frac{1}{y!} \exp(y \log y - y) \right\} \exp \left\{ -\frac{1}{\mu} \Psi(y \log y - y - y \log \mu + \mu) \right\}.$$

حال با توجه به رابطه بالا و تعریف DM در (۱.۱) داریم

$$d(y; \mu) = \Psi(y \log \frac{y}{\mu} - y + \mu)$$

$$a(y, \sigma^2) = \frac{1}{y!} \exp \{y \log y - y\}$$

$$\sigma^2 = 1.$$

برای نشان دادن این که تابع d یک تابع انحراف منظم است، کافی است نشان دهیم که این تابع محدب دارای مینیمم یکتای صفر است. برای هر μ داریم

$$\frac{\partial d}{\partial y} = \Psi(\log y - \log \mu), \quad \frac{\partial^2}{\partial y^2} = \frac{\Psi}{y}.$$

واضح است که مشتق اول منفی است اگر $y < \mu$ و مثبت است اگر $y > \mu$ ، و به این معنی است که تابع انحراف d یک تابع محدب با مینیمم مقدار یکتای صفر در $y = \mu$ است؛ در نتیجه با توجه به تعریف تابع واریانس در رابطه (۲.۱) داریم

$$V(\mu) = \frac{\Psi}{\frac{\partial^2 d}{\partial \mu^2}} = \frac{\Psi}{\frac{\Psi y}{\mu^2}} = \frac{\Psi}{\mu} = \mu.$$

۱.۲.۱ خواص مدل‌های پراکندگی

خاصیت ۱. اگر d تابع انحراف منظم باشد آن‌گاه

$$\frac{\partial^2 d}{\partial y^2}(y; y) = \frac{\partial^2 d}{\partial \mu^2}(y; y) = -\frac{\partial^2 d}{\partial \mu \partial y}(y; y), \quad \forall y \in \Omega. \quad (3.1)$$

برهان. با استفاده از تعریف انحراف داریم

$$d(y; y) = d(\mu; \mu) = 0, \quad d(y; \mu) \geq 0, \quad \forall y, \mu \in \Omega$$

یعنی $d(y; \cdot)$ دارای یک مینیمم یکتا در y است و به‌طور مشابه $d(\cdot; \mu)$ دارای یک مینیمم یکتا در μ است، بنابراین

$$\frac{\partial d}{\partial \mu}(y; y) = \frac{\partial d}{\partial y}(y; y) = 0. \quad (4.1)$$

نتیجه (۳.۱) به سادگی با استفاده از مشتق دو معادله (۴.۱) نسبت به y حاصل می‌شود.

□

خاصیت ۲. بسط تیلور انحراف منظم در نزدیکی (μ_0, μ_0) به صورت زیر است

$$d(\mu_0 + x\delta; \mu_0 + m\delta) = \frac{\delta^2}{V(\mu_0)}(x - m)^2 + O(\delta^2)$$

که در آن $V(\cdot)$ همان تابع واریانس است.

برهان. با توجه به (۴.۱) داریم

$$\begin{aligned} d(\mu_0 + x\delta; \mu_0 + m\delta) &= d(\mu_0, \mu_0) + \frac{\partial d}{\partial \mu}(x\delta) + \frac{\partial d}{\partial y}(m\delta) \\ &+ \frac{1}{2} \frac{\partial^2 d}{\partial \mu^2}(\mu_0, \mu_0)(\delta^2 x^2) + \frac{1}{2} \frac{\partial^2 d}{\partial \mu \partial y}(\mu_0, \mu_0)(\delta m)(\delta x) \\ &+ \frac{1}{2} \frac{\partial^2 d}{\partial y^2}(\mu_0, \mu_0)(\delta^2 m^2) + O(\delta^2) \end{aligned}$$

از آنجا که سه جمله اول سمت راست عبارت بالا صفر هستند، بنابراین داریم:

$$\begin{aligned} d(\mu_0 + x\delta; \mu_0 + m\delta) &= \frac{\delta^2}{V(\mu_0)}x^2 + \frac{\delta^2}{V(\mu_0)}2xm + \frac{\delta^2}{V(\mu_0)}m^2 + O(\delta^2) \\ &= \frac{\delta^2}{V(\mu_0)}(x - m)^2 + O(\delta^2). \end{aligned}$$

□

از آنجایی که ثابت نرمال‌ساز $a(\cdot)$ صورت بسته‌ای ندارد و به دست آوردن پارامتر پراکندگی σ^2 کار دشواری است؛ تقریبی از ثابت نرمال‌ساز $a(\cdot)$ ، که تقریب نقطه‌زینی برای چگالی با پراکندگی کوچک است، را با استفاده از گزاره زیر به دست می‌آوریم.

خاصیت ۳. (تقریب نقطه‌زینی^{۱۹}) هرگاه $\sigma^2 \rightarrow 0$ ، چگالی یک مدل پراکندگی منظم، DM، می‌تواند به صورت زیر تقریب شود

$$f(y; \mu, \sigma^2) \simeq \{2\pi\sigma^2 V(y)\}^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}d(y; \mu)\right\}$$

که رابطه فوق معادل است با این که بگوییم وقتی $\sigma^2 \rightarrow 0$ ، ثابت نرمال‌ساز به صورت زیر تقریب می‌شود:

$$a(y; \sigma^2) \simeq \{2\pi\sigma^2 V(y)\}^{-\frac{1}{2}} \quad (5.1)$$

¹⁹ Saddlepoint Approximation

که در آن $V(\cdot)$ همان تابع واریانس است. برای اثبات به بارندروف، نیلسن و کاکس (۱۹۸۹) رجوع کنید.

از خواص ۲ و ۳ می‌توان خاصیت نرمال مجانبی را که در زیر آمده است، نتیجه گرفت.

خاصیت ۴. (نرمال مجانبی) فرض کنید $Y \sim DM(\mu_0 + \sigma\mu, \sigma^2)$ یک مدل پراکندگی با تقریب نقطه زینی همگرای یکنواخت است، یعنی همگرایی در (۵.۱) یکنواخت است؛ داریم

$$\frac{Y - \mu_0}{\sigma} \xrightarrow{d} N(\mu, V(\mu_0)), \quad \sigma^2 \rightarrow 0. \quad (6.1)$$

به عبارتی دیگر، برای پراکندگی‌های کوچک داریم

$$DM(\mu_0 + \sigma\mu, \sigma^2) \simeq N(\mu_0 + \sigma\mu, \sigma^2 V(\mu_0)).$$

۲.۲.۱ برآورد درست‌نمایی ماکسیمم

در این بخش به برآورد درست‌نمایی ماکسیمم در GLM بر اساس مدل‌های پراکندگی می‌پردازیم. یک مجموعه داده مقطعی (y_i, x_i) ، $i = 1, \dots, m$ را در نظر بگیرید، که در آن y_i ها مصداقی از متغیرهای تصادفی مستقل و هم توزیع $Y_i \sim DM(\mu_i, \sigma^2)$ هستند و $g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta}$ و $Y_i \sim DM(\mu_i, \sigma^2)$ قرار می‌دهیم $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma^2)$ درست‌نمایی برای بردار پارامتر $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma^2)$ به صورت زیر است

$$L(\boldsymbol{\theta}; \mathbf{y}) = \prod_{i=1}^m a(y_i; \sigma^2) \exp \left\{ -\frac{1}{2\sigma^2} d(y_i; \mu_i) \right\}, \boldsymbol{\beta} \in \mathbb{R}^{p+1}, \sigma^2 > 0.$$

لگاریتم درست‌نمایی نیز به صورت زیر بیان می‌شود

$$\begin{aligned} \ell(\boldsymbol{\theta}; \mathbf{y}) &= \sum_{i=1}^m \log a(y_i; \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m d(y_i; \mu_i) \\ &= \sum_{i=1}^m \log a(y_i; \sigma^2) - \frac{1}{2\sigma^2} D(\mathbf{y}; \boldsymbol{\mu}) \end{aligned} \quad (7.1)$$

که در آن $\mu_i = \mu_i(\boldsymbol{\beta})$ یک تابع غیرخطی از $\boldsymbol{\beta}$ و $D(\mathbf{y}; \boldsymbol{\mu}) = \sum_{i=1}^m d(y_i; \mu_i)$ فقط وابسته به $\boldsymbol{\beta}$ است. تابع D مشابه مجموع توان‌های دوم باقی‌مانده در مدل رگرسیون خطی است. تابع امتیاز برای ضریب رگرسیونی $\boldsymbol{\beta}$ به صورت زیر است:

$$U(\mathbf{y}; \boldsymbol{\beta}) = \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\beta}} = -\frac{1}{2\sigma^2} \sum_{i=1}^m \frac{\partial d(y_i; \mu_i)}{\partial \mu_i} \frac{\partial \mu_i}{\partial \boldsymbol{\beta}}.$$

i امین پیشگوی خطی را با $\eta_i = \mathbf{x}_i^T \boldsymbol{\beta}$ و همچنین تابع امتیاز انحراف^{۲۰} را با u نمادگذاری می‌کنیم. داریم

$$u(y_i; \mu_i) = -\frac{1}{2} \frac{\partial d(y_i; \mu_i)}{\partial \mu_i}, \quad i = 1, \dots, m. \quad (8.1)$$

$$\frac{\partial \mu_i}{\partial \beta} = \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta} = \{g'(\mu_i)\}^{-1} \mathbf{x}_i$$

که در آن مشتق مرتبه اول تابع پیوند، نسبت به μ است. در جدول ۲.۱ برخی از توابع پیوند و مشتق آن‌ها آمده است.

جدول ۲.۱: برخی توابع پیوند و مشتق اول آن‌ها

مدل	تابع پیوند	مشتق تابع پیوند	Ω
سیمپلکس، دوجمله‌ای	$\log\left(\frac{\mu}{1-\mu}\right)$	$\frac{1}{\mu(1-\mu)}$	$\mu \in (0, 1)$
پواسون، گوسی معکوس، دوجمله‌ای منفی	$\log(\mu)$	$\frac{1}{\mu}$	$\mu \in (0, \infty)$
گاما	$\frac{1}{\mu}$	$-\frac{1}{\mu^2}$	$\mu \in (0, \infty)$

حال تابع امتیاز را به صورت زیر می‌توان نوشت

$$U(\mathbf{y}; \beta) = \frac{1}{\sigma^2} \sum_{i=1}^m \mathbf{x}_i \frac{1}{g'(\mu_i)} u(y_i; \mu_i). \quad (9.1)$$

برآورد درست‌نمایی ماکسیمم پارامتر β از حل معادله زیر حاصل می‌شود

$$\sum_{i=1}^m \mathbf{x}_i \frac{1}{g'(\mu_i)} u(y_i; \mu_i) = 0. \quad (10.1)$$

باید دقت داشت که معادله (۱۰.۱) شامل پارامتر پراکندگی σ^2 نیست؛ تحت شرایط نظم، برآوردگر درست‌نمایی ماکسیمم $\hat{\beta}_m$ یک برآوردگر سازگار برای β است، یعنی

$$\hat{\beta}_m \xrightarrow{P} \beta, m \rightarrow \infty$$

و دارای نرمال مجانبی با میانگین صفر و ماتریس کواریانس $I^{-1}(\theta)$ است که $I(\theta)$ ماتریس اطلاع فیشر است که به صورت زیر بیان می‌شود

$$\begin{aligned} I(\theta) &= -E\{U'(\mathbf{Y}; \beta)\} \\ &= \frac{1}{\sigma^2} \sum_{i=1}^m \mathbf{x}_i \frac{1}{\{g'(\mu_i)\}^2} E\{-u'(Y_i; \mu_i)\} \mathbf{x}_i^T \\ &= \sum_{i=1}^m \mathbf{x}_i w_i \mathbf{x}_i^T \\ &= \mathbf{X}^T \mathbf{W} \mathbf{X}. \end{aligned}$$

که در آن X یک ماتریس با بعد $m \times (p+1)$ است که i امین سطر آن x_i می‌باشد. هم‌چنین W یک ماتریس قطری است که i امین عنصر قطری آن w_i است که به صورت زیر بیان می‌شود

$$w_i = \frac{1}{\sigma^2} \frac{1}{\{g'(\mu_i)\}^2} E\{-u'(Y_i, \mu_i)\}, \quad i = 1, \dots, m.$$

زمانی که پارامتر پراکندگی در مدل وجود داشته باشد، یعنی اگر $a(y, \sigma^2)$ به گونه‌ای باشد که مشتق‌گیری آن مقدور باشد، برای برآورد درست‌نمایی ماکسیمم آن مشابه ثابت نرمال‌ساز، کافی است مانند توزیع نرمال از آن مشتق گرفت؛ اما زمانی که ثابت نرمال‌ساز دارای صورت بسته نباشد، مشتق آن نسبت به σ^2 دارای حل عددی پیچیده‌ای است که برای رفع این مشکل و برآورد پارامتر σ^2 دو روش پیشنهاد می‌شود. اولین روش، استفاده از خاصیت نرمال‌جانبی کوچک پراکندگی است که در آن

$$\log a(y; \sigma^2) \simeq -\frac{1}{\sigma^2} \log \sigma^2.$$

با اعمال این تقریب در رابطه (۷.۱)، تقریب لگاریتم درست‌نمایی نسبت به σ^2 به صورت زیر است

$$-\frac{m}{2\sigma^2} + \frac{1}{2\sigma^4} D(y; \mu) = 0$$

که منجر به برآورد σ^2 به صورت زیر می‌شود

$$\widehat{\sigma^2} = \frac{1}{m} D(y; \hat{\mu}) = \frac{1}{m} \sum_{i=1}^m d(y_i, \hat{\mu}_i).$$

این برآوردگر را برآوردگر جرجنسن می‌نامند که میانگینی از برآورد انحراف است. برآوردگر جرجنسن در حالت کلی نارایب نیست، حتی اگر روی درجه‌های آزادی آن تصحیح انجام دهیم، یعنی $m - (p+1)$ را جایگزین m کنیم. این رابطه فقط وقتی پیشنهاد می‌شود که پارامتر σ^2 کوچک (< 5) باشد.

روش دوم برای به دست آوردن یک برآوردگر نارایب از σ^2 استفاده از ویژگی گشتاوری است که در زیر به آن می‌پردازیم.

گزاره ۱.۲.۱. اگر $Y \sim DM(\mu, \sigma^2)$ و $d(y; \mu)$ یک انحراف منظم باشد، آن‌گاه

$$E\{u(Y; \mu)\} = 0,$$

$$V\{u(Y; \mu)\} = \sigma^2 E\{-u'(Y; \mu)\}$$

که در آن u' مشتق اول نسبت به μ ، در رابطه (۸.۱) است.

برهان. می‌دانیم که $\int f(y; \mu, \sigma^2) dy = 1$ ، حال نسبت به μ از آن مشتق گرفته و برابر صفر قرار می‌دهیم

$$-\frac{1}{2\sigma^2} \int d'(y; \mu) f(y; \mu, \sigma^2) dy = 0,$$

یا این که می توان گفت $E\{d'(Y; \mu)\} = 0$. بار دیگر نسبت به μ مشتق گرفته و نتیجه می شود

$$-\frac{1}{\sigma^2} \int \{d'(y; \mu)\}^2 f(y; \mu, \sigma^2) dy + \int d''(y; \mu) f(y; \mu, \sigma^2) dy = 0$$

یا به طور معادل، گوییم

$$E\{d''(Y; \mu)\} = \frac{1}{\sigma^2} E\{d'(Y; \mu)\}^2 = \frac{1}{\sigma^2} V\{d'(Y; \mu)\}.$$

با توجه به رابطه (۸.۱)، می توان رابطه بالا را بازنویسی کرد؛ که عبارت است از

$$V\{u(Y; \mu)\} = \sigma^2 E\{-u'(Y; \mu)\}. \quad (11.1)$$

□

بر اساس این نتیجه می توان پارامتر پراکندگی σ^2 را با روش گشتاورها برآورد کرد. این برآوردگر عبارت است از

$$\widehat{\sigma^2} = \frac{\sum_{i=1}^m (u_i - \bar{u})^2}{\sum_{i=1}^m (-u'_i)},$$

که در آن $\bar{u} = \frac{1}{m} \sum_{i=1}^m u_i$ و $u'_i = u'(y_i, \hat{\mu}_i)$ ، $u_i = u(y_i, \hat{\mu}_i)$

۳.۱ معیارهای اندازه گیری دقت مدل

در آمار معیارهایی برای ارزیابی توانایی مدل ها در برآورد متغیر پاسخ وجود دارند که با استفاده از آن ها می توان نتایج مدل های مختلف را با یکدیگر مقایسه کرد. این معیارها با اندازه گیری خطای برآورد یا دقت برآورد یک مدل، نتایج را ارزیابی و مقایسه می کنند. معیارهای متفاوتی برای این هدف وجود دارند که در این بخش، به چند معیار اشاره خواهیم کرد.

میانگین توان دوم خطا

یکی از راه های اندازه گیری خطای برآورد مدل، معیار میانگین توان دوم خطا^{۲۱} (MSE) است. برای مقایسه برآوردگرهای اریب با برآوردگرهای ناریب و یا مقایسه دو برآوردگر یک پارامتر، با اریبی های متفاوت، MSE معیار مفیدی است و مقدار آن به صورت زیر محاسبه می گردد:

$$MSE(\hat{\theta}) = E\left[(\theta - \hat{\theta})^2\right]$$

که در آن θ مقدار واقعی پارامتر و $\hat{\theta}$ مقدار برآورد شده پارامتر است. در صورتی که $MSE(\theta)$ به سمت صفر میل کند، بیان گر سازگار^{۲۲} بودن $\hat{\theta}$ می باشد.

²¹ Mean Squared Error

²² Consistency

اریبی

آماره $\delta(X)$ که X یک متغیر تصادفی با توزیعی از خانواده توزیع‌های $\theta \in \Theta$ است، برای $\gamma(\theta)$ اریب^{۲۳} گفته می‌شود اگر برای یک $\theta \in \Theta$ ،

$$E_{\theta}[\delta(X)] \neq \gamma(\theta)$$

آن‌گاه اریبی برآوردگر $\delta(X)$ در برآورد $\gamma(\theta)$ به صورت زیر است:

$$b(\theta) = E_{\theta}[\delta(X)] - \gamma(\theta)$$

معیار اطلاع آکاییک

معیار اطلاع آکاییک^{۲۴} (AIC)، معیاری برای انتخاب بهترین مدل آماری از بین چند نامزد است. این معیار، نیکویی برازش مدل را همراه با پیچیدگی، به‌طور توأم، در نظر می‌گیرد که توسط آکاییک در سال ۱۹۷۴ پیشنهاد شد. با داده‌های مشاهده‌شده، چند مدل رقیب بر حسب مقدار AIC رتبه‌بندی می‌شوند و مدل با کمترین AIC بهترین است. در حالت کلی، AIC به صورت

$$AIC = -2 \log(L) + 2p$$

تعریف می‌شود که در آن p تعداد پارامترهای مدل و L مقدار ماکسیمم تابع درست‌نمایی برای مدل برآوردشده است. به‌طور مجانبی مینیمم کردن AIC معادل مینیمم کردن آماره CV است. این نتیجه برای هر مدلی، نه فقط مدل خطی، برقرار است. این ویژگی باعث شده است زمانی که هدف پیشگویی است، استفاده از AIC توصیه شود.

معیار اطلاع بیزی

معیار اطلاع بیزی^{۲۵} (BIC) به صورت زیر

$$BIC = -2 \log(L) + p \log(n)$$

تعریف می‌شود که در آن n حجم نمونه است. به دلیل جریمه سنگین‌تر این معیار نسبت به AIC، مدلی که توسط BIC انتخاب می‌شود، یا همان مدل منتخب بر اساس AIC است یا مدلی است با تعداد پارامتر کمتر.

²³ Bias

²⁴ Akaike Information Criterion

²⁵ Bayesian Information Criterion

فصل ۲

توزیع سیمپلکس

مقدمه

با توجه به هدف این پایان نامه که بررسی مدل رگرسیونی با متغیر پاسخ از نوع سیمپلکس است، این فصل به شناخت توزیع سیمپلکس اختصاص دارد. از آن جا که توزیع سیمپلکس، تبدیلی از توزیع گوسی معکوس آمیخته است، ابتدا توزیع گوسی معکوس^۱ و توزیع گوسی معکوس آمیخته را مطرح می کنیم و سپس توزیع سیمپلکس، خواص آن و تولید اعداد تصادفی سیمپلکس را تشریح می کنیم.

۱.۲ توزیع گوسی معکوس

توزیع گوسی معکوس یکی از پرکاربردترین توزیع ها در تحلیل قابلیت اعتماد و هم چنین به عنوان مدلی برای طول عمر و زمان خرابی است. هم چنین این توزیع در تحلیل داده های مثبت چوله به راست نیز مورد استفاده قرار می گیرد. توزیع گوسی معکوس توسط شرودینگر در سال ۱۹۱۵ به دلیل خواص مهم فیزیکی آن تحت عنوان توزیع «اولین لحظه گذر در حرکت براونی» معرفی شد. فولکز و چیکارا (۱۹۸۸) به شباهت های قابل توجهی بین توزیع گوسی معکوس و خانواده گوسی اشاره کرده اند. جرجنسن (۱۹۹۱) تابع توزیع گوسی معکوس تعمیم یافته را

^۱Inverse Gaussian Distribution

معرفی و خواص آن را بیان کرد. همچنین مطالعه‌ای گسترده درباره پیدایش و خواص نظری توزیع گوسی معکوس توسط سشادری (۱۹۹۹) صورت گرفت. مادولکار و ناتارجان (۲۰۰۲) خصوصیات مختلف این توزیع از جمله تقارن، کشیدگی و چولگی را بررسی کردند و نشان دادند که توزیع گوسی معکوس متعلق به خانواده توزیع‌های نمایی است و در صورتی که داده‌ها نامنفی و دارای چولگی مثبت باشند، می‌تواند به‌عنوان جایگزینی برای توزیع گوسی استفاده شود. مادولکار و وانگ (۲۰۰۷) درباره توزیع وارون متغیر گوسی معکوس مطالعاتی انجام دادند. بالاکیریشن و همکاران (۲۰۰۹) به معرفی توزیع گوسی معکوس آمیخته و تبدیل‌ها و گشتاورهای آن پرداختند و یک روش جدید برای محاسبه گشتاورهای توزیع معرفی کردند و از گشتاورهای توزیع برای انتخاب مدل مناسب استفاده کردند. در تحلیل داده‌ها، در بسیاری از موارد داده‌ها متقارن نبوده و دارای چولگی هستند. یکی از کاربردهای مهم این توزیع در مسئله داده‌های چوله است.

تعریف ۱.۱.۲. اگر متغیر تصادفی پیوسته و مثبت X دارای تابع چگالی احتمال زیر باشد

$$f(x; \xi, \epsilon^2) = \left(2\pi\epsilon^2 x^3\right)^{-\frac{1}{3}} \exp\left\{-\frac{(x-\xi)^2}{2\xi^2\epsilon^2 x}\right\}, \quad x > 0, \epsilon^2 > 0, \xi > 0 \quad (1.2)$$

آن‌گاه X دارای توزیع گوسی معکوس با نماد $X \sim IG(\xi, \epsilon^2)$ است که پارامتر ξ میانگین توزیع و ϵ^2 پارامتر مقیاس است.

نمودار تابع چگالی گوسی معکوس در شکل ۱.۲ به ازای مقادیر مختلفی از ξ و ϵ^2 نشان داده شده است. امید ریاضی و واریانس توزیع گوسی معکوس عبارتند از

$$E(X) = \xi$$

$$V(X) = \xi^3 \epsilon^2$$

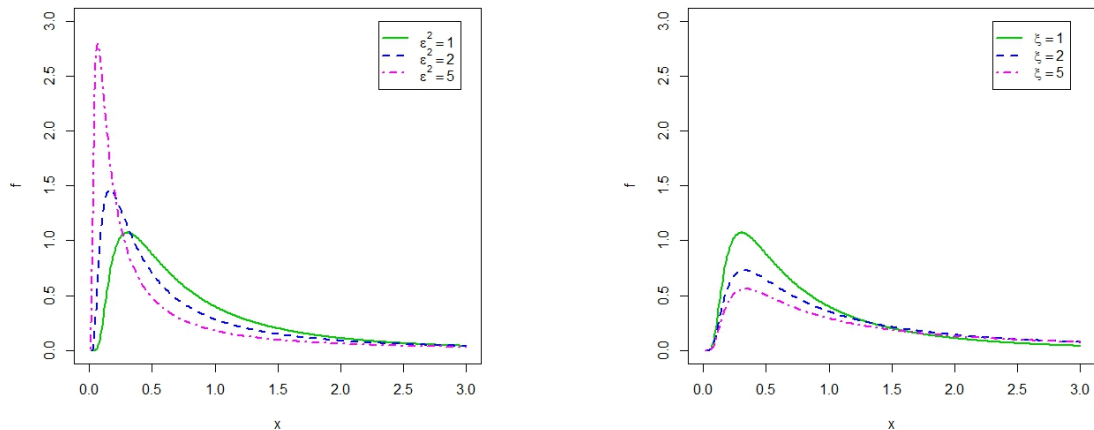
توزیع گوسی معکوس یک خانواده نمایی دو پارامتری با آماره‌های طبیعی $\frac{1}{X}$ و پارامترهای طبیعی $\frac{-1}{\xi^2 \epsilon^2}, \frac{-1}{\xi^2}$ است. با قرار دادن $\xi = 1$ و $\epsilon = 1$ در چگالی (۱.۲) صورت استاندارد توزیع به صورت

$$f(x; 1, 1) = \left(2\pi x^3\right)^{-\frac{1}{3}} \exp\left\{-\frac{(x-1)^2}{2x}\right\}$$

حاصل می‌شود که در واقع شکل حدی توزیع نمونه در یک آزمون نسبت احتمال دنباله‌ای می‌باشد که توسط والد به‌دست آمده است. از این‌رو، گاهی اوقات آنرا توزیع والد استاندارد گویند.

برای $\epsilon^2 = 1$ ، تابع چگالی و تابع توزیع تجمعی گوسی معکوس به ترتیب به صورت زیر هستند

$$f_{IG}(x; \xi) = \left(2\pi x^3\right)^{-\frac{1}{3}} \exp\left\{-\frac{(x-\xi)^2}{2\xi^2 x}\right\} \quad (2.2)$$



(آ) نمودار تابع چگالی گوسی معکوس برای $\epsilon^2 = 1$ و مقادیر مختلف ξ .
 (ب) نمودار تابع چگالی گوسی معکوس برای $\xi = 1$ و مقادیر مختلف ϵ^2 .

شکل ۱.۲: نمودار توابع چگالی گوسی معکوس

$$F_{IG}(x; \xi) = \Phi\left(\frac{1}{\sqrt{x}} \frac{x - \xi}{\xi}\right) + \exp\left(\frac{\xi}{x}\right) \Phi\left(-\frac{1}{\sqrt{x}} \frac{x + \xi}{\xi}\right) \quad x > 0, \xi > 0$$

که در آن $\Phi(\cdot)$ تابع توزیع تجمعی نرمال استاندارد است.

۱.۱.۲ تابع مولد گشتاور توزیع گوسی معکوس

اگر متغیر تصادفی $X \sim IG(\xi, \epsilon^2)$ باشد، در این صورت

$$M_X(t) = \int_0^{\infty} e^{tx} f(x; \xi, \epsilon^2) dx = e^{\frac{\xi}{\epsilon^2}} \int_0^{\infty} \left\{ \frac{1}{\sqrt{2\pi\epsilon^2 x^3}} \right\}^{\frac{1}{2}} \exp\left\{ \frac{-(x^2 + \xi^2) - 2\xi^2 x^2 t}{2\xi^2 \epsilon^2 x} \right\} dx$$

که می‌توان نوشت

$$\begin{aligned} \frac{-(x^2 + \xi^2) - 2\xi^2 x^2 t}{2\xi^2 \epsilon^2 x} &= \frac{-\left(\sqrt{1 - 2\xi^2 \epsilon^2 t} x - \xi\right)^2}{2\xi^2 \epsilon^2 x} - \frac{1}{\xi \epsilon^2} \sqrt{1 - 2\xi^2 \epsilon^2 t} \\ &= \frac{-(x - \xi')^2}{2\xi'^2 \epsilon^2 x} - \frac{1}{\xi \epsilon^2} \sqrt{1 - 2\xi^2 \epsilon^2 t} \end{aligned}$$

که در آن $\xi' = \frac{\xi}{(1 - 2\xi^2 \epsilon^2 t)^{\frac{1}{2}}}$. با جایگذاری رابطه بالا در انتگرال نمایش $M_X(t)$ داریم

$$M_X(t) = \exp\left\{ \frac{1}{\xi \epsilon^2} \right\} \exp\left\{ \frac{-1}{\xi \epsilon^2} \sqrt{1 - 2\xi^2 \epsilon^2 t} \right\} \int_0^{\infty} \sqrt{\frac{1}{2\pi\epsilon^2 x^3}} \exp\left\{ \frac{-(x - \xi')^2}{2\xi'^2 \epsilon^2 x} \right\} dx$$

تابع زیر انتگرال بالا تابع چگالی احتمال یک توزیع گوسی معکوس با پارامترهای ξ و ϵ^2 است. بنابراین مقدار انتگرال، برابر با یک است و در نتیجه تابع مولد گشتاور گوسی معکوس به صورت زیر است:

$$M_X(t) = \exp \left\{ \frac{1}{\xi \epsilon^2} \left[1 - \left(1 - 2 \xi^2 \epsilon^2 t \right)^{\frac{1}{2}} \right] \right\}, \quad t < \frac{1}{2 \xi^2 \epsilon^2}$$

۲.۱.۲ وجه تسمیه معکوس

اگر $M_X(t)$ نشان دهنده تابع مولد گشتاور متغیر تصادفی X باشد، تابع مولد انباشتگی^۲ متغیر تصادفی X به صورت زیر تعریف می شود

$$\psi(t) = \log M_X(t) = \log E[e^{tX}]$$

حال اگر $L_X(t)$ نشان دهنده لگاریتم تبدیل لاپلاس متغیر تصادفی نامنفی X باشد، بنابراین داریم

$$L_X(t) = \log E[e^{-tX}] = \psi(-t)$$

در سال ۱۹۴۵ تویدی ارتباطی وارون بین تابع مولد انباشتگی جفت‌هایی از متغیر تصادفی دو جمله‌ای و دو جمله‌ای منفی، پواسون و گاما، نرمال و گوسی معکوس ارائه نمود. در این جا ابتدا حالت کلی متغیرهای تصادفی معکوس معرفی می شود، سپس این ارتباط وارون در مورد جفت متغیرهای تصادفی که دارای این خاصیت هستند مورد بررسی قرار گیرند.

تعریف ۲.۱.۲. متغیرهای تصادفی X و Y معکوس یکدیگرند اگر و تنها اگر $L_X(t)$ و $L_Y(t)$ به ازای هر مقدار t در دامنه مشترک در شرایط زیر صدق کند

$$i) L_X(t) = \alpha L(t) \quad ii) L_Y(t) = \beta L^{-1}(t) \quad (۳.۲)$$

که در آن $L_X(t)$ نشان دهنده لگاریتم تبدیل لاپلاس متغیر تصادفی X است، α و β مقادیر ثابت و $L^{-1}(t)$ وارون تابع $L(t)$ می باشد، یعنی $L(L^{-1}(t)) = t$. توزیع‌های متناظر با متغیرهای تصادفی معکوس را توزیع‌های معکوس می نامند.

فرض کنید متغیر تصادفی X دارای توزیع نرمال با میانگین μ و واریانس σ^2 با تابع چگالی احتمال زیر باشد

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x-\xi)^2}{2\sigma^2} \right\}, \quad x \in \mathbb{R}, \xi > 0, \sigma > 0$$

آن گاه

$$E[e^{tX}] = \exp \left\{ \xi t + \frac{\sigma^2 t^2}{2} \right\} \implies L_X(t) = \xi \left(\frac{\sigma^2 t^2}{2\xi} - t \right), \quad t \in \mathbb{R}$$

² Cumulant Generating Function

با فرض $\alpha = \xi$ و $L(t) = \frac{\sigma^2 t^2}{2\xi} - t$ ، به راحتی می توان دید که

$$L^{-1}(t) = \frac{\xi}{\sigma^2} \left[1 - \left(1 + \frac{2\sigma^2}{\xi} t \right)^{\frac{1}{2}} \right] \quad t < \frac{\xi}{2\sigma^2}$$

با قرار دادن $\frac{1}{\epsilon^2} = \frac{\xi^3}{\sigma^2}$ داریم

$$L^{-1}(t) = \frac{1}{\xi^2 \epsilon^2} \left[1 - \left(1 + 2\xi^2 \epsilon^2 t \right)^{\frac{1}{2}} \right] \quad t < \frac{1}{2\xi^2 \epsilon^2}$$

حال اگر متغیر تصادفی Y دارای تابع چگالی احتمال (۱.۲) باشد، آن گاه

$$L_Y(t) = \frac{1}{\xi \epsilon^2} \left[1 - \left(1 + 2\xi^2 \epsilon^2 t \right)^{\frac{1}{2}} \right] = \xi L^{-1}(t) \quad t < \frac{1}{2\xi^2 \epsilon^2}$$

با فرض $\beta = \xi$ به ازای تمام مقادیر $t < \frac{1}{2\xi^2 \epsilon^2}$ روابط (۳.۲) در مورد متغیرهای تصادفی X و Y حفظ می شود؛ یعنی متغیر تصادفی Y معکوس متغیر تصادفی نرمال است.

۳.۱.۲ برخی خواص توزیع گوسی معکوس

توزیع گوسی معکوس در مقایسه با بسیاری از توزیع های کلاسیک از خواصی برخوردار می باشد که بعضی از آنها مشابه توزیع نرمال است. در زیر به بررسی برخی از مهم ترین ویژگی های این توزیع می پردازیم.

ویژگی ۱- اگر $X \stackrel{iid}{\sim} IG(\xi, \epsilon^2)$ آن گاه $kX \sim IG(k\xi, k\epsilon^2)$ ، $k > 0$

ویژگی ۲- اگر $X_i \stackrel{iid}{\sim} IG(\xi, \epsilon^2)$ ، $i = 1, \dots, n$ آن گاه $\sum X_i \sim IG(n\xi, n\epsilon^2)$

ویژگی ۳- اگر $X_i \stackrel{iid}{\sim} IG(\xi, \epsilon^2)$ ، $i = 1, \dots, n$ آن گاه $\bar{X} \sim IG(\xi, n\epsilon^2)$

ویژگی ۴-

قضیه ۱.۱.۲. (شاستر، ۱۹۶۸) اگر متغیر تصادفی X دارای تابع چگالی به صورت رابطه (۱.۲) باشد، آن گاه

$$\frac{(X - \xi)^2}{X\xi^2\epsilon^2} \sim \chi_{(1)}^2$$

ویژگی ۵- (وارون توزیع گوسی معکوس^۳) اگر $X \sim IG(\xi, 1)$ که تابع چگالی آن در رابطه (۲.۲) داده شد، آن گاه تابع چگالی احتمال متغیر تصادفی X^{-1} (وارون توزیع گوسی معکوس) که آن را با نماد $RIG(\xi, 1)$ نشان می دهیم، به صورت زیر است

$$f_{RIG}(x; \xi) = (2\pi x)^{-\frac{1}{2}} \exp \left\{ \frac{-(1 - \xi x)^2}{2\xi^2 x} \right\} \quad (۴.۲)$$

³ Reciprocal Inverse Gaussian Distribution

و تابع توزیع تجمعی آن نیز به صورت

$$\begin{aligned} F_{RIG}(x; \xi) &= 1 - F_{IG}\left(\frac{1}{x}, \xi\right) \\ &= 1 - \Phi\left(\frac{1}{\sqrt{x}} \frac{1 - \xi x}{\xi}\right) - \exp\left(\frac{2}{\xi}\right) \Phi\left(-\frac{1}{\sqrt{x}} \frac{1 + \xi x}{\xi}\right) \end{aligned} \quad (5.2)$$

است.

برهان. برای اثبات رابطه (۴.۲)، با توجه به رابطه (۲.۲) و بنا به اصول روش تغییر متغیر می‌دانیم

$$f_{RIG}(x; \xi) = \frac{1}{x^2} f_{IG}\left(\frac{1}{x}; \xi\right) \quad (6.2)$$

که $\frac{1}{x^2}$ ژاکوبین تبدیل است و لذا

$$\begin{aligned} f_{RIG}(x; \xi) &= \frac{1}{x^2} \left[(2\pi x^{-2})^{-\frac{1}{2}} \exp\left\{ \frac{-(x^{-1} - \xi)^2}{2\xi^2 x^{-1}} \right\} \right] \\ &= \left[(2\pi x)^{-\frac{1}{2}} \exp\left\{ \frac{-(x^{-1} - \xi)^2 x^2}{2\xi^2 x^{-1} x^2} \right\} \right] \\ &= \left[(2\pi x)^{-\frac{1}{2}} \exp\left\{ \frac{-x^2(x^{-2} + \xi^2 - 2\xi x^{-1})}{2\xi^2 x} \right\} \right] \\ &= \left[(2\pi x)^{-\frac{1}{2}} \exp\left\{ \frac{-(1 + x^2\xi^2 - 2\xi x)}{2\xi^2 x} \right\} \right] \\ &= \left[(2\pi x)^{-\frac{1}{2}} \exp\left\{ \frac{-(1 - \xi x)^2}{2\xi^2 x} \right\} \right] \end{aligned}$$

حال برای اثبات رابطه (۵.۲) داریم

$$P(X \leq t) = F_X(t) = P\left(\frac{1}{X} \geq \frac{1}{t}\right) = P\left(Y \geq \frac{1}{t}\right) = 1 - F_Y\left(\frac{1}{t}\right).$$

بنابراین

$$F_Y(t) = 1 - F_X\left(\frac{1}{t}\right).$$

□

لازم به ذکر است که از این ویژگی در توصیف توزیع گوسی معکوس آمیخته در بخش ۲.۲ استفاده خواهیم کرد.

۴.۱.۲ تولید مقادیر تصادفی از توزیع گوسی معکوس

فرض کنید متغیر تصادفی پیوسته X دارای تابع توزیع F باشد، به سادگی می‌توان نشان داد که $Y = F(X)$ دارای توزیع یکنواخت روی بازه $(0, 1)$ است. حال اگر تابع وارون F به آسانی

قابل محاسبه باشد می توان با به دست آوردن معکوس F ، اعداد تصادفی که از توزیع F پیروی می کنند را تولید کرد. کافی است اعدادی تصادفی در بازه $(0, 1)$ اختیار شوند و آن ها را در ضابطه تابع معکوس قرار دهند و مقادیر متغیر متناظر محاسبه شود. در مورد توزیع گوسی معکوس به دلیل شکل پیچیده تابع توزیع، محاسبه معکوس F به راحتی امکان پذیر نیست. میکائیل و همکاران در سال ۱۹۷۶ روشی برای تولید مشاهدات تصادفی از یک توزیع گوسی معکوس ارائه کردند. فرض کنید $X \sim IG(\xi, \epsilon^2)$. با توجه به ویژگی ۴ توزیع گوسی معکوس، قضیه ۱.۱.۲، که در آن $Z = \frac{(X - \xi)^2}{X \xi^2 \epsilon^2} \sim \chi_{(1)}^2$ است، اگر X را بر حسب Z به دست آوریم با یک معادله درجه دو بر حسب X مواجه هستیم. فرض کنید $X^{(1)}$ و $X^{(2)}$ ریشه های این معادله باشند؛ داریم

$$X^{(1)} = \xi + \frac{\xi^2 \epsilon^2 Z}{2} - \frac{\xi \epsilon^2}{2} \sqrt{\frac{4 \xi Z}{\epsilon^2} + \xi^2 Z^2}, \quad X^{(2)} = \frac{\xi^2}{X^{(1)}}. \quad (7.2)$$

میکائیل و همکارانش احتمال شرطی انتخاب ریشه ها با فرض تولید یک متغیر تصادفی Z از توزیع کای دو با یک درجه آزادی را محاسبه نمودند و نشان دادند که ریشه کوچک تر با احتمال $\frac{\xi}{\xi + X^{(1)}}$ و ریشه دیگر، با احتمال $\frac{X^{(1)}}{\xi + X^{(1)}}$ انتخاب می شود. پس اگر Z را از توزیع کای دو با یک درجه آزادی تولید کنیم، آن گاه با روابط (۷.۲) می توان به مقادیر متغیر گوسی معکوس دست پیدا کرد؛ به این صورت که یک آزمایش برنولی با احتمال موفقیت $\frac{\xi}{\xi + X^{(1)}}$ انجام می دهیم؛ برای اثبات به میکائیل و همکاران، (۱۹۷۶) رجوع کنید. لازم به ذکر است که برای تولید یک مشاهده از توزیع کای دو با یک درجه آزادی، کافی است یک مشاهده از توزیع نرمال استاندارد را مجذور نماییم و برای تولید یک مشاهده نرمال استاندارد از تبدیلات زیر استفاده می کنیم

$$Z_1 = \sqrt{-2 \ln U_1} \cos 2\pi U_2, \quad Z_2 = \sqrt{-2 \ln U_1} \sin 2\pi U_2$$

که در آن U_1 و U_2 دو مشاهده مستقل از توزیع یکنواخت در بازه $(0, 1)$ ، Z_1 و Z_2 دو متغیر نرمال استاندارد مستقل هستند.

بنابراین می توان الگوریتم تولید اعداد تصادفی گوسی معکوس را به صورت زیر ارائه داد.

الگوریتم ۱.۱.۲. الگوریتم تولید اعداد تصادفی از گوسی معکوس

گام ۱- اعداد تصادفی Z از توزیع کای دو با یک درجه آزادی تولید می کنیم.
گام ۲- برای هر عدد تصادفی در گام ۱، ریشه کوچک تر، $X^{(1)}$ ، را براساس رابطه (۷.۲) محاسبه می کنیم.

گام ۳- یک آزمایش برنولی با احتمال موفقیت $p = \frac{\xi}{\xi + X^{(1)}}$ اجرا می کنیم.

گام ۴- اگر نتیجه آزمایش موفقیت باشد ریشه $X^{(1)}$ را به عنوان مشاهده ای از توزیع گوسی معکوس انتخاب می کنیم و در غیر این صورت ریشه بزرگ تر $X^{(2)}$ انتخاب می شود.

رابطه بین ریشه‌های معادله $X^{(1)}X^{(2)} = \xi^2$ است.

۲.۲ توزیع گوسی معکوس آمیخته

بالاکیریشن و همکاران (۲۰۰۹) توزیع گوسی معکوس آمیخته و تبدیل‌ها و گشتاورهای آن را معرفی کرده و یک روش جدید برای محاسبه گشتاورهای این توزیع معرفی کردند. در تحلیل پارامتری داده‌ها، مدل‌های چوله اهمیت بسیاری دارند. توزیع‌های آمیخته نیز بسیار انعطاف‌پذیرند و دارای خصوصیات آماری خوبی نیز می‌باشند. توزیع گوسی معکوس آمیخته یکی از مدل‌های آماری است که کاربرد زیادی دارد. با انجام تبدیلاتی روی توزیع گوسی معکوس آمیخته توزیع‌های پارامتری جدیدی به‌دست می‌آیند که در موقعیت‌های مختلف مفیدند.

تعریف ۱.۲.۲. فرض کنید $f_{X_1}(x)$ و $f_{X_2}(x)$ به ترتیب توابع چگالی احتمال متغیرهای تصادفی X_1 و X_2 باشند. متغیر تصادفی X را آمیخته گسسته گویند، اگر دارای چگالی زیر باشد:

$$f_X(x) = (1-p)f_{X_1}(x) + pf_{X_2}(x), \quad x > 0, 0 \leq p \leq 1.$$

تعریف ۲.۲.۲. متغیر تصادفی X دارای توزیع گوسی معکوس آمیخته با پارامترهای مکان ξ ، مقیاس ϵ^2 و پارامتر آمیختگی p است، هرگاه تابع چگالی احتمال آن به صورت زیر باشد

$$f_X(x) = \phi\left(\frac{1}{\sqrt{\xi\epsilon^2}}\left[\sqrt{\frac{x}{\xi}} - \sqrt{\frac{\xi}{x}}\right]\right) \frac{1}{\sqrt{\epsilon^2 x^3}} \left[1 - p + \frac{px}{\xi}\right], \quad x > 0 \quad (۸.۲)$$

که در آن $\phi(\cdot)$ تابع چگالی احتمال نرمال استاندارد است. این توزیع را با نماد $X \sim MIG(x; \xi, \epsilon^2, p)$ نشان می‌دهند.

توزیع گوسی معکوس آمیخته، یک خانواده انعطاف‌پذیر از توزیع‌های آمیخته و چوله به راست است.

قضیه ۱.۲.۲. اگر متغیر تصادفی X دارای توزیع گوسی معکوس آمیخته باشد، $X \sim MIG(\xi, \epsilon^2, p)$ ، آن‌گاه

$$CX \sim MIG(C\xi, C\epsilon^2, p) \quad (\text{الف})$$

(ب) $R = X^{-1}$ (وارون توزیع گوسی معکوس آمیخته^۴) دارای تابع چگالی زیر است، که آن را به اختصار با نماد $\frac{1}{X} \sim MIG\left(\frac{1}{\xi}, \frac{1}{\xi^2\epsilon^2}, 1-p\right)$ نشان می‌دهیم

$$f_R(r) = \phi\left(\frac{1}{\sqrt{\xi\epsilon^2}}\left[\sqrt{\frac{1}{r\xi}} - \sqrt{r\xi}\right]\right) \frac{1}{r^2} \left[(1-p)\sqrt{\frac{r^3}{\epsilon^2}} + p\sqrt{\frac{r}{\xi^2\epsilon^2}}\right], \quad r > 0.$$

⁴ Reciprocal Mixture Inverse Gaussian Distribution

برهان. الف: با تغییر متغیر نتیجه مستقیماً حاصل می‌شود.

ب: با تغییر متغیر $x = \frac{1}{r}$ و ژاکوبین تبدیل $|\frac{1}{r^2}|$ در رابطه (۸.۲) داریم

$$\begin{aligned} f_R(r) &= \phi \left(\frac{1}{\sqrt{\xi \epsilon^2}} \left[\sqrt{\frac{1}{\xi r}} - \sqrt{\frac{\xi}{r}} \right] \right) \frac{1}{\sqrt{\epsilon^2 (\frac{1}{r})^3}} \left[1 - p + \frac{p}{\xi} \right] \times \frac{1}{r^2} \\ &= \phi \left(\frac{1}{\sqrt{\xi \epsilon^2}} \left[\sqrt{\frac{1}{r \xi}} - \sqrt{r \xi} \right] \right) \frac{1}{r^2} \left[(1-p) \sqrt{\frac{r^3}{\epsilon^2}} + p \sqrt{\frac{r}{\xi \epsilon^2}} \right] \end{aligned}$$

□

و اثبات کامل می‌شود.

گزاره ۱.۲.۲. اگر $X_1 \sim IG(\xi, \epsilon^2)$ و $X_2^{-1} \sim IG(\xi^{-1}, \xi^2 \epsilon^2)$ ، آن‌گاه متغیر تصادفی X دارای توزیع گوسی معکوس آمیخته با پارامترهای ξ ، ϵ^2 و p که X_1 با احتمال $1-p$ و X_2 با احتمال p است و تابع چگالی احتمال آن به صورت

$$f(x; \xi, \epsilon^2, p) = (2\pi \epsilon^2 x^3)^{-\frac{1}{2}} \left(1 - p + \frac{px}{\xi} \right) \exp \left\{ -\frac{1}{2\epsilon^2} \frac{(x-\xi)^2}{\xi^2 x} \right\}, \quad x > 0 \quad (9.2)$$

است.

برهان. بنا به رابطه (۱.۲) می‌دانیم

$$f_{X_1}(x) = (2\pi \epsilon^2 x^3)^{-\frac{1}{2}} \exp \left\{ -\frac{(x-\xi)^2}{2\xi^2 \epsilon^2 x} \right\}, \quad x > 0, \epsilon^2 > 0, \xi > 0$$

قرار دهید $X_2^{-1} = U$ و $D(x, \xi) = \frac{(x-\xi)^2}{\xi^2 x}$ و $D(x^{-1}, \xi^{-1}) = \frac{D(x, \xi)}{\epsilon^2}$. حال چگالی مربوط به آن را با توجه به این که $X_2^{-1} \sim IG(\xi^{-1}, \xi^2 \epsilon^2)$ و بنا به رابطه بالا به صورت زیر به دست می‌آوریم

$$f_U(u) = (2\pi \xi^2 \epsilon^2 u^3)^{-\frac{1}{2}} \exp \left\{ \frac{-D(u, \xi^{-1})}{2\xi^2 \epsilon^2} \right\}$$

حال نیازمند به دست آوردن توزیع خود X_2 هستیم که از رابطه (۶.۲) می‌توان آن را به صورت زیر به دست آورد

$$\begin{aligned} f_{X_2}(u) &= \frac{1}{u^2} f_U\left(\frac{1}{u}\right) \\ &= \frac{1}{u^2} (2\pi \xi^2 \epsilon^2 u^{-3})^{-\frac{1}{2}} \exp \left\{ \frac{-D(u^{-1}, \xi^{-1})}{2\xi^2 \epsilon^2} \right\} \\ &\Rightarrow f_{X_2}(x) = (2\pi \xi^2 \epsilon^2 x)^{-\frac{1}{2}} \exp \left\{ \frac{-D(x, \xi)}{2\epsilon^2} \right\} \end{aligned}$$

حال با توجه به تعریف (۱.۲.۲) داریم:

$$f_X(x) = (1-p) (2\pi \epsilon^2 x^3)^{-\frac{1}{2}} \exp \left\{ \frac{-D(x, \xi)}{2\xi^2 \epsilon^2 x} \right\} + p (2\pi \xi^2 \epsilon^2 x)^{-\frac{1}{2}} \exp \left\{ \frac{-D(x, \xi)}{2\epsilon^2} \right\}$$

□

بنابراین با فاکتور گرفتن از جمله مشترک و ساده‌سازی، رابطه (۹.۲) به دست می‌آید.

۱.۲.۲ تولید مقادیر تصادفی از توزیع گوسی معکوس آمیخته

متغیرهای تصادفی مستقل از هم X_1 و X_2 را در نظر بگیرید که توزیع آن‌ها $X_1 \sim IG(\xi, \epsilon^2)$ و $X_2 \sim IG(\xi^{-1}, \epsilon^2)$ است. اکنون متغیر تصادفی X را به صورت زیر تعریف می‌کنیم

$$X = \begin{cases} X_1 & \text{با احتمال } 1-p \\ X_2 & \text{با احتمال } p \end{cases}$$

که آمیخته‌ای از X_1 و X_2 است. جانسون و کاتز (۱۹۷۰) نشان دادند که X_2 را می‌توان به صورت ترکیب خطی $X_1 + X_3$ نوشت (ضمیمه ۱.۱ را ببینید)، که $X_3 \sim \xi^2 \epsilon^2 \chi_{(1)}^2$. پس می‌توان نوشت $X = X_1 + V$ ، که در آن

$$V = \begin{cases} 0 & \text{با احتمال } 1-p \\ X_3 & \text{با احتمال } p \end{cases}$$

و هم‌چنین X_1 و V مستقلند.

طبق مطالب فوق یک روش شبیه‌سازی ساده برای تولید مقادیر تصادفی X از توزیع گوسی معکوس آمیخته به شرح زیر قابل بیان است.

الگوریتم ۱.۲.۲. الگوریتم تولید اعداد تصادفی از گوسی معکوس آمیخته

گام ۱- تولید مقادیر متغیر تصادفی از توزیع $X_1 \sim IG(\xi, \epsilon^2)$ (زیر بخش ۴.۱.۲)

گام ۲- تولید مقادیر متغیر تصادفی از تابع توزیع $W \sim \chi_{(1)}^2$ و قرار می‌دهیم $X_3 = \xi^2 \epsilon^2 W$

گام ۳- تولید اعداد تصادفی از توزیع یکنواخت، $U \sim U(0, 1)$

گام ۴- اگر $U < p$ ، $X = X_1 + X_3$ ، در غیر این صورت $X = X_1$.

در این صورت X متغیر تصادفی تولیدشده از توزیع گوسی معکوس آمیخته با پارامترهای مذکور می‌باشد.

۳.۲ توزیع سیمپلکس

توزیع سیمپلکس برگرفته از توزیع گوسی معکوس تعمیم‌یافته است (بارندروف، نیلسن و جرجنسن، ۱۹۹۱) که تابع چگالی احتمال آن به صورت زیر تعریف می‌شود

$$f(y; \alpha_1, \alpha_2, \mu, \sigma^2) = c(\alpha_1, \alpha_2, \mu, \sigma^2) y^{\alpha_1-1} (1-y)^{\alpha_2-1} \exp\left\{-\frac{1}{\sigma^2} d_{\alpha_1, \alpha_2}(y; \mu)\right\}, y \in (0, 1) \quad (10.2)$$

که در آن μ, α_2, α_1 و σ^2 پارامترهای توزیع و $c(\alpha_1, \alpha_2, \mu, \sigma^2)$ ثابت نرمال‌ساز و $d_{\alpha_1, \alpha_2}(y; \mu)$ انحراف است که به صورت

$$d_{\alpha_1, \alpha_2}(y; \mu) = \mu^{\alpha_1-1} (1-\mu)^{\alpha_2-1} \frac{(y-\mu)^2}{y(1-y)}$$

تعریف می‌شود و تابع واریانس آن نیز به صورت زیر است:

$$V_{\alpha_1, \alpha_2}(\mu) = \mu^{2(1-\alpha_1)}(1-\mu)^{2(1-\alpha_2)}.$$

دامنه پارامترهای توزیع عبارتند از $\mu \in (0, 1)$ ، $(\alpha_1, \alpha_2)^T \in \mathbb{R}^2$ و $\sigma^2 > 0$. چگالی (۱۰.۲) را هم‌چنین توزیع سیمپلکس کلی^۵ می‌نامند و آن را با $Y \sim S(\alpha_1, \alpha_2, \mu, \sigma^2)$ نمایش می‌دهند. اگر $\alpha_1, \alpha_2 > 0$ و $\sigma^2 \rightarrow \infty$ آن‌گاه Y دارای توزیع بتا با پارامترهای α_1 و α_2 و $Y \sim \text{Beta}(\alpha_1, \alpha_2)$ است.

حالت خاصی از توزیع (۱۰.۲) زمانی است که $(\alpha_1, \alpha_2) = (-\frac{1}{2}, -\frac{1}{2})$ که آن را سیمپلکس استاندارد گویند و از این پس با توزیع سیمپلکس استاندارد کار خود را پیش برده و برای سهولت آن را «توزیع سیمپلکس» می‌نامیم.

تعریف ۱.۳.۲. (متغیر تصادفی سیمپلکس) متغیر تصادفی Y را دارای توزیع سیمپلکس با پارامتر مکان $\mu \in (0, 1)$ و پارامتر پراکندگی $\sigma^2 > 0$ گویند و با نماد $Y \sim S^-(\mu, \sigma^2)$ نشان می‌دهند اگر تابع چگالی آن به صورت زیر باشد

$$f(y; \mu, \sigma^2) = [2\pi\sigma^2 \{y(1-y)\}^3]^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}d(y; \mu)\right\}, \quad y \in (0, 1) \quad (11.2)$$

که در آن

$$d(y; \mu) = \frac{(y-\mu)^2}{y(1-y)\mu^2(1-\mu)^2}.$$

جرجنسن (۱۹۹۷)، واریانس Y را به صورت زیر

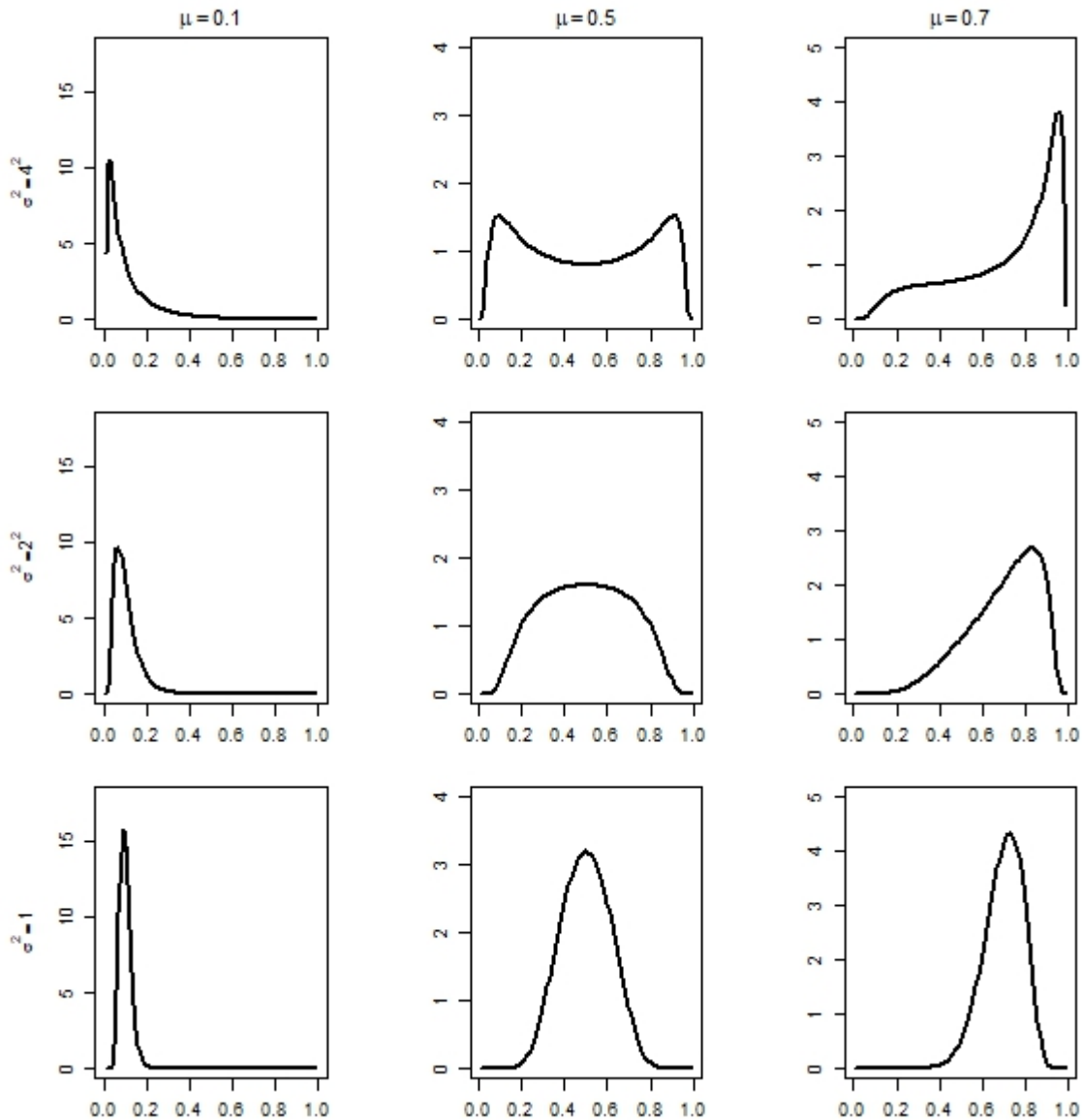
$$\tau^2 = \mu(1-\mu) - \frac{1}{\sqrt{2}\sigma^2} \exp\left\{-\frac{1}{2\sigma^2\mu^2(1-\mu)^2}\right\} \Gamma\left\{\frac{1}{2}, \frac{1}{2\sigma^2\mu^2(1-\mu)^2}\right\}$$

محاسبه کرد که در آن $\Gamma(a, b)$ ، تابع گامای ناکامل، به صورت $\Gamma(a, b) = \int_b^\infty t^{a-1} e^{-t} dt$ است.

شکل ۲.۲ نمودار تابع چگالی سیمپلکس را با پارامترهای متفاوت مکان و پراکندگی نشان می‌دهد. همان‌گونه که ملاحظه می‌شود، توزیع سیمپلکس شامل یک کلاس بزرگ از توزیع‌های در محدوده $(0, 1)$ است به طوری که اگر $\sigma \leq \frac{1}{\sqrt{3}}$ ، توزیع تک‌مُدی و در غیر این صورت چند مُدی خواهد بود. در این شکل، ردیف اول و دوم و سوم به ترتیب متناظر با σ^2 برابر با ۱، ۲، ۴ است، هم‌چنین ستون اول و دوم و سوم از راست به ترتیب متناظر با μ برابر با ۰/۱، ۰/۵، ۰/۷ است.

در مقایسه رابطه (۱۱.۲) و (۱.۱) به آسانی می‌توان دریافت که توزیع سیمپلکس متعلق به خانواده مدل‌های پراکندگی است.

⁵General Simplex Distribution



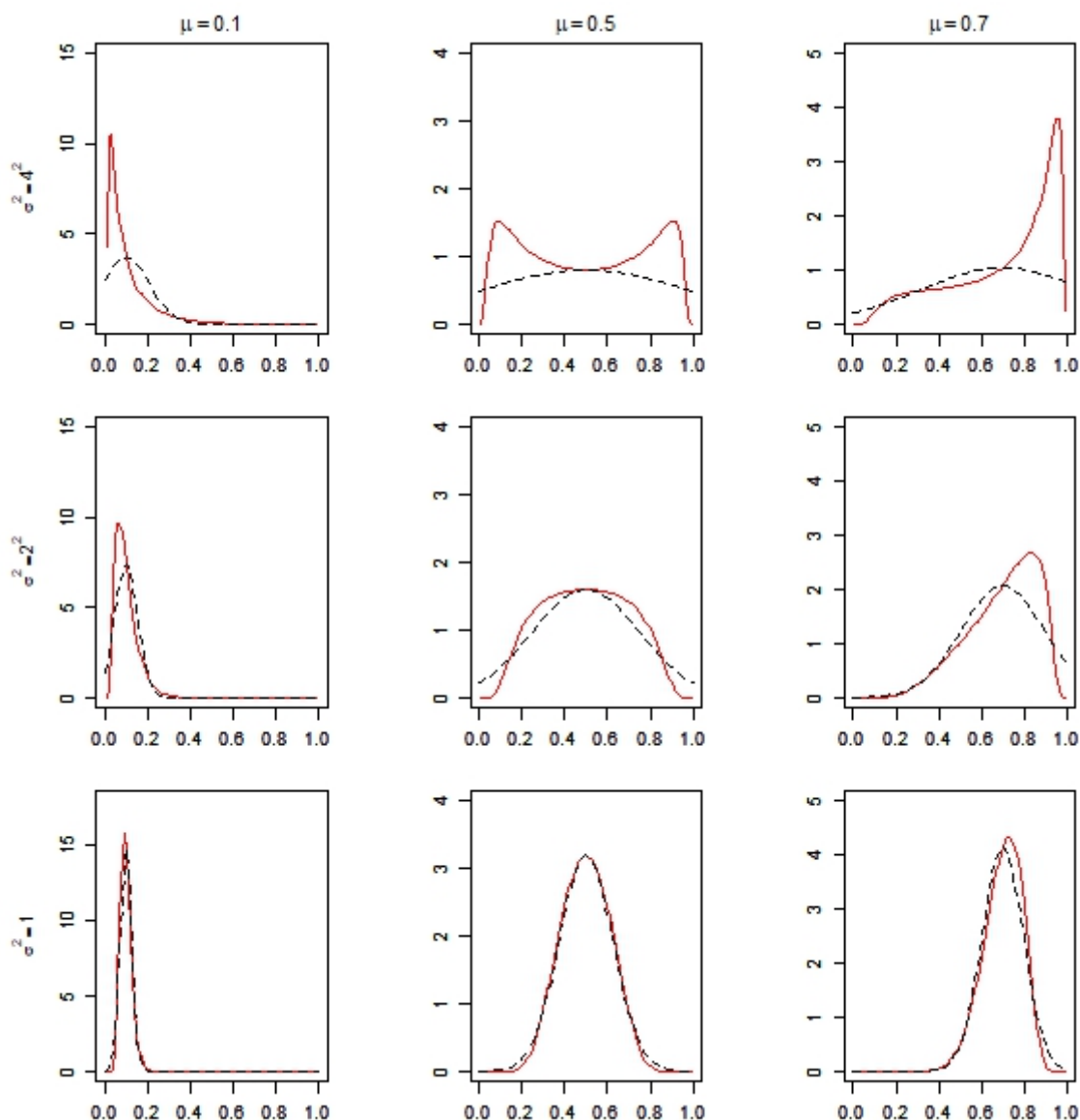
شکل ۲.۲: نمودار تابع چگالی سیمپلکس

با توجه به خاصیت ۴ مدل های پراکندگی، رابطه (۶.۱)، و رابطه (۷.۱) (بخش ۲.۲.۱ را ببینید) می توان گفت که اگر $\sigma^2 \rightarrow 0$ ، آن گاه

$$\frac{Y - \mu}{\sigma \sqrt{V(\mu)}} \xrightarrow{d} N(0, 1) \quad (۱۲.۲)$$

که در آن $V(\mu) = \mu^3(1 - \mu)^3$ تابع واریانس است. برای بررسی این موضوع، نمودار چگالی احتمال سیمپلکس و نرمال $N(\mu, \sigma^2 V(\mu))$ را در یک قاب و به ازای μ و σ^2 های مختلف رسم کرده ایم (شکل ۲.۲). در این شکل نیز همانند قبل، ردیف اول و دوم و سوم به ترتیب متناظر با σ^2 برابر است با ۱، ۲، ۴. هم چنین ستون اول و دوم و سوم به ترتیب متناظر با μ برابر با ۰/۱، ۰/۵، ۰/۷ است. همان گونه که مشخص است، به ازای σ^2 کوچک تر، توزیع نرمال تطابق بیشتری با

توزیع سیمپلکس دارد.



شکل ۳.۲: نمودار توابع چگالی سیمپلکس (منحنی توپر) و نرمال (منحنی مقطع)

۱.۳.۲ رابطه بین توزیع گوسی معکوس آمیخته و توزیع سیمپلکس

اگر $X \sim MIG(\xi, \epsilon^2, p)$ و قرار دهیم

$$\mu = \frac{\xi}{1 + \xi}, \quad \sigma^2 = \frac{\epsilon^2}{(1 - p)^2} = \epsilon^2 (1 - \xi)^2$$

آن‌گاه

$$Y = \frac{X}{1 + X} \sim S^-(\mu, \sigma^2).$$

برای اثبات به جرجنسن (۱۹۹۷) و ژو (۲۰۰۱) رجوع کنید. یعنی توزیع سیمپلکس می‌تواند با اعمال تغییر متغیری بر توزیع گوسی معکوس آمیخته، به دست آید. این موضوع، کلیدی‌ترین نکته در تولید اعداد تصادفی سیمپلکس است که در بخش بعد آن را دنبال می‌کنیم.

۲.۳.۲ تولید مقادیر تصادفی از توزیع سیمپلکس

با توجه به این که اگر $X \sim MIG(\xi, \epsilon^2)$: آن گاه $Y = \frac{X}{1+X} \sim S^-(\mu, \sigma^2)$ ، لذا الگوریتم تولید اعداد تصادفی سیمپلکس دقیقاً همان الگوریتم تولید اعداد تصادفی از توزیع MIG است که در انتها، مقادیر Y از رابطه بالا به دست می‌آید. این الگوریتم را می‌توان برای تولید اعداد تصادفی از توزیع سیمپلکس، و به استناد زیر بخش ۱.۲.۲ به صورت زیر خلاصه کرد:

الگوریتم ۱.۳.۲. الگوریتم تولید اعداد تصادفی از توزیع سیمپلکس

$$\text{گام ۱-} \quad \text{قرار دهید } \epsilon^2 = \sigma^2 (1 - \mu)^2, \quad \xi = \frac{\mu}{1 - \mu}, \quad p = \mu.$$

گام ۲- عدد تصادفی X را از توزیع گوسی معکوس آمیخته، $X \sim MIG(\xi, \epsilon^2, p)$ (زیر بخش ۱.۲.۲)، تولید کنید.

$$\text{گام ۳-} \quad \text{برای تولید عدد تصادفی سیمپلکس قرار دهید } Y = \frac{X}{1+X}.$$

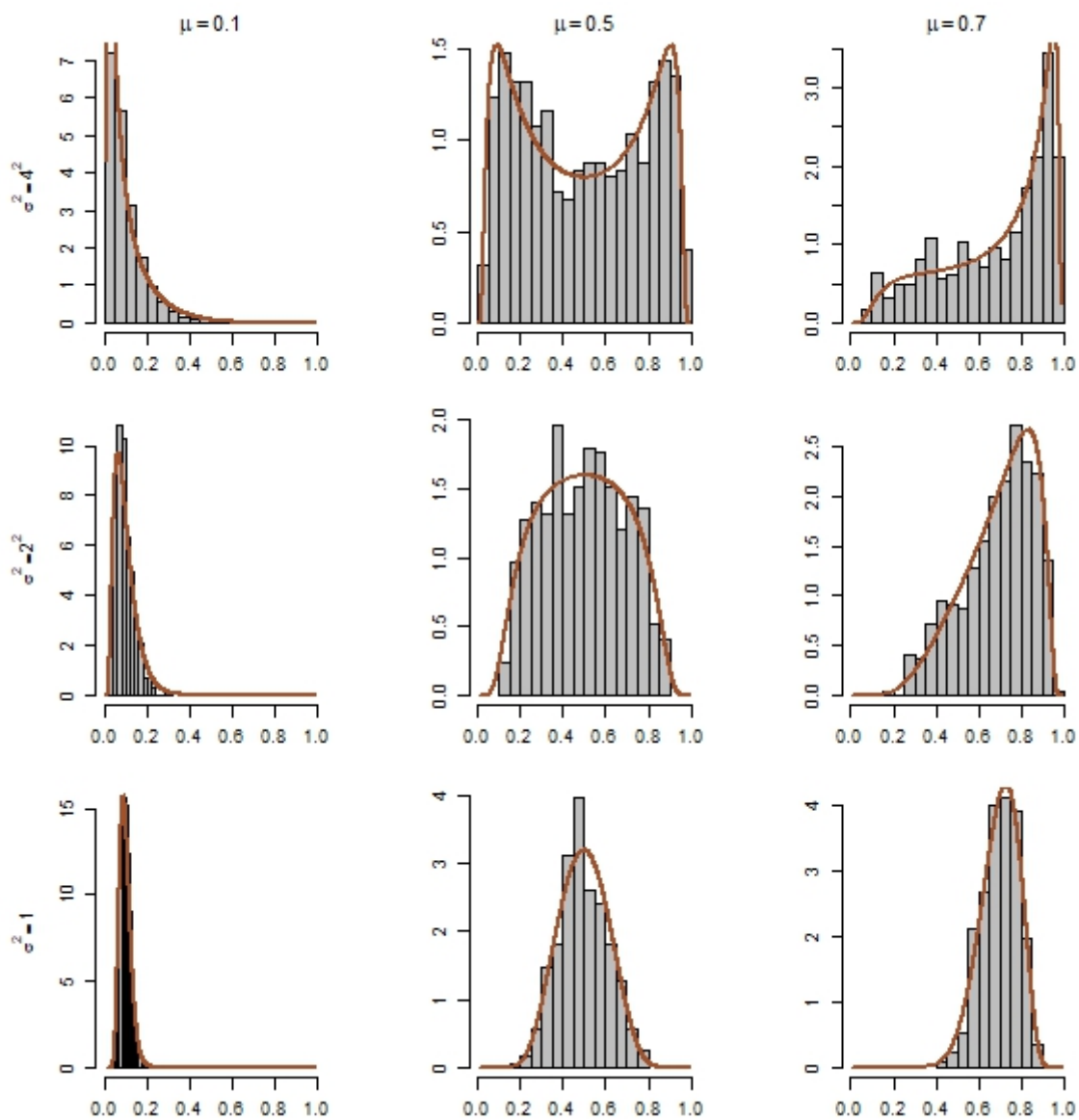
بررسی تطابق داده‌های تولیدشده از توزیع سیمپلکس با چگالی توزیع سیمپلکس

در شکل ۴.۲، هر نمودار بیانگر بافت‌نگار ۵۰۰ عدد تصادفی تولیدشده از توزیع سیمپلکس و هم‌چنین تابع چگالی سیمپلکس $S^-(\mu, \sigma^2)$ است. همان‌طور که ملاحظه می‌شود، تابع چگالی سیمپلکس پوشش مناسبی برای بافت‌نگار ایجاد کرده است. هم‌چنین در شکل ۵.۲ هر نمودار بیانگر هسته چگالی با ۱۵۰۰ عدد تصادفی تولیدشده از توزیع سیمپلکس و تابع چگالی سیمپلکس $S^-(\mu, \sigma^2)$ است.

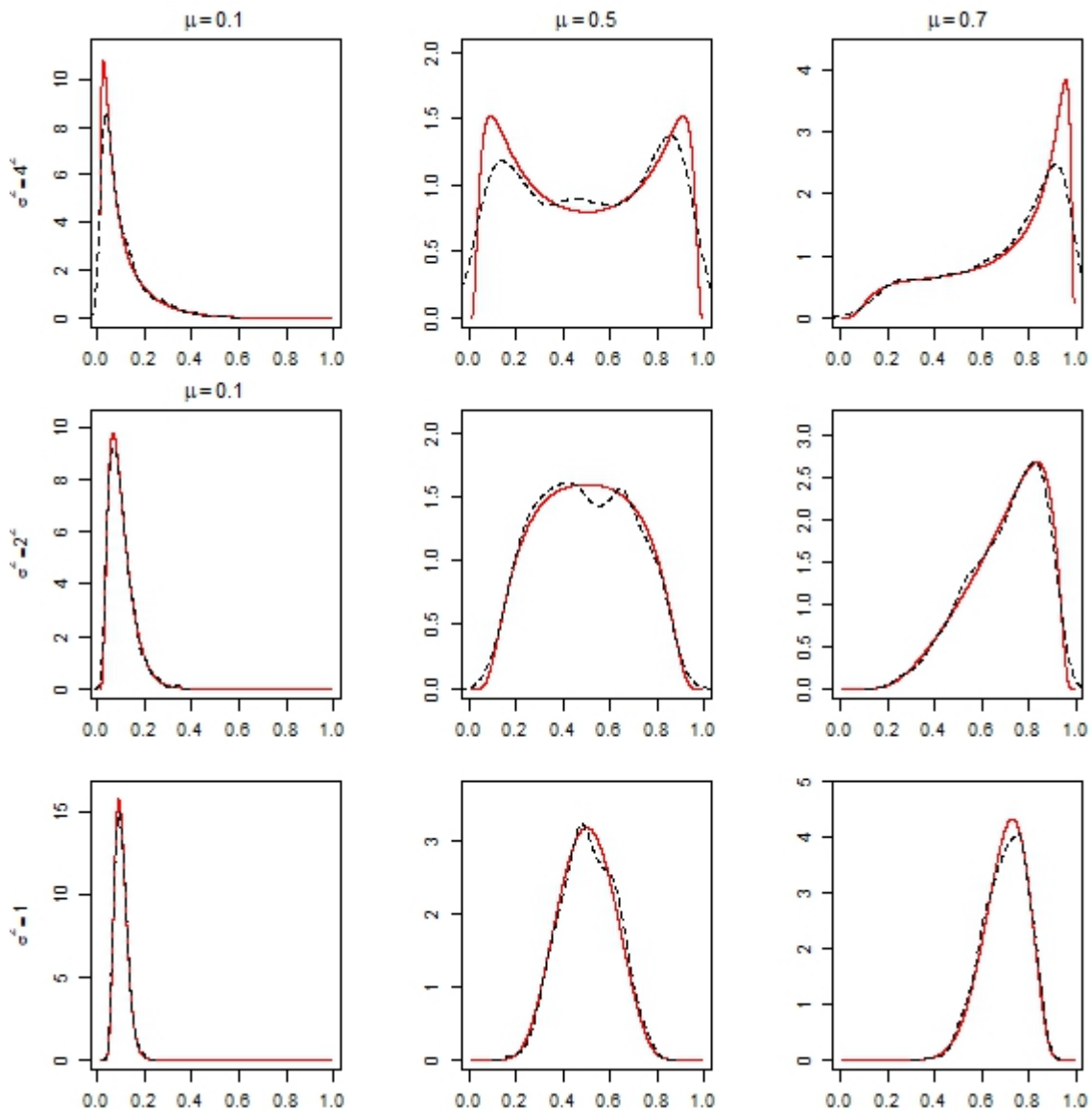
دستورات مربوط به توزیع احتمالی سیمپلکس در بسته‌افزار `simplexreg` در نرم‌افزار R به شرح زیر هستند:

```
dsimplex(x, mu, sig)
psimplex(q, mu, sig)
qsimplex(p, mu, sig)
rsimplex(n, mu, sig)
```


که به ترتیب تابع چگالی احتمال، تابع توزیع تجمعی، تابع چندک و تولید اعداد تصادفی هستند.



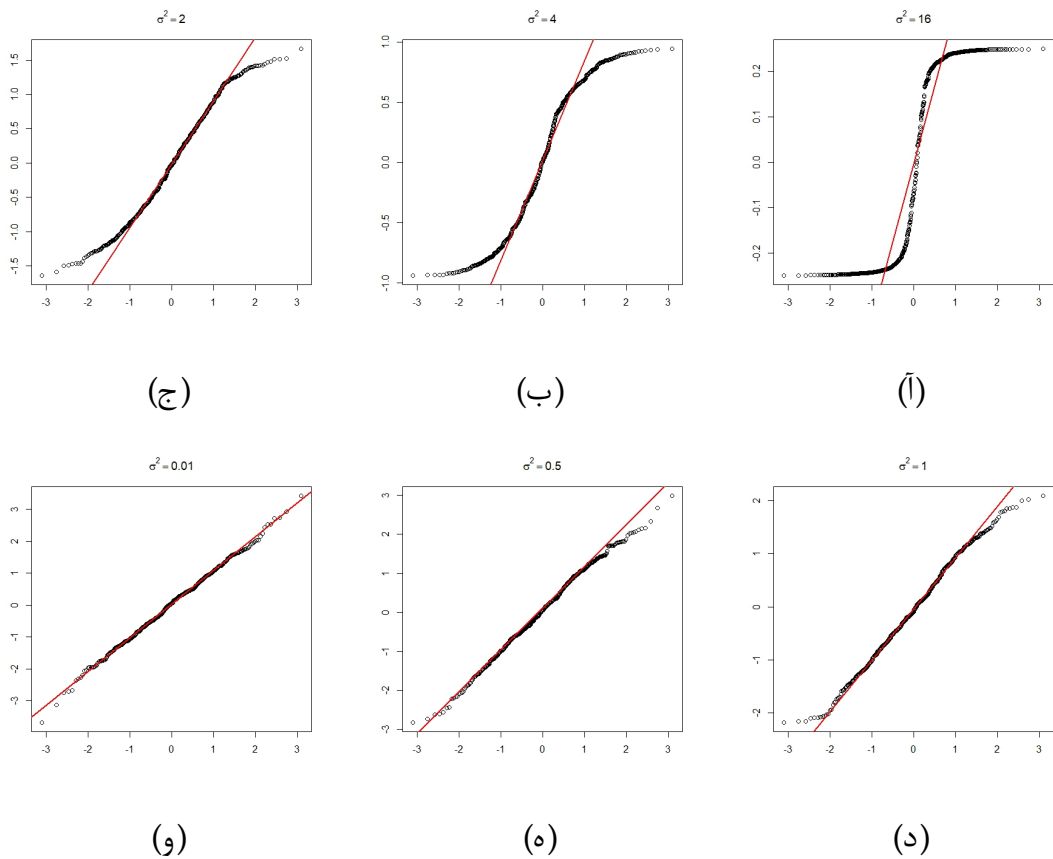
شکل ۴.۲: نمودار چگالی (منحنی توپر) به همراه بافت‌نگار اعداد تصادفی از توزیع سیمپلکس



شکل ۵.۲: نمودار چگالی (منحنی توپر) به همراه هسته چگالی (منحنی مقطع) اعداد تصادفی از توزیع سیمپلکس

بررسی نظریه نرمال مجانبی توسط نمودار Q-Q

برای بررسی تقریب نرمال مندرج در رابطه (۱۲.۲) مبنی بر این که اگر $\sigma^2 \rightarrow \infty$ آن گاه $\frac{Y - \mu}{\sigma\sqrt{V(\mu)}} \xrightarrow{d} N(0, 1)$ ، تعداد ۵۰۰ عدد تصادفی سیمپلکس را با $\mu = 0.5$ و σ^2 های مختلف (۱۶، ۴، ۲، ۱، ۰/۵، ۰/۱) تولید کرده و برای آزمون نرمال بودن، از ابزار بصری نمودار چندک چندک استفاده کرده ایم. همان طور که در شکل ۶.۲ دیده می شود، می توان نرمال مجانبی بودن را نتیجه گرفت.



شکل ۶.۲: نمودار چندک-چندک

۳.۳.۲ برخی خواص توزیع سیمپلکس

اگر $Y \sim S^-(\mu, \sigma^2, p)$ ، با $E[Y] = \mu$ آن گاه

$$(۱) E\{d(Y; \mu)\} = \sigma^2$$

$$(۲) E\{(Y - \mu)d'(Y; \mu)\} = -2\sigma^2$$

$$(۳) E\{(Y - \mu)d(Y; \mu)\} = 0$$

$$(۴) E\{d'(Y; \mu)\} = 0$$

$$(۵) \frac{1}{p} E\{d''(Y; \mu)\} = \frac{3\sigma^2}{\mu(1-\mu)} + \frac{1}{\mu^3(1-\mu)^3}$$

$$(۶) V\{d(Y; \mu)\} = 2(\sigma^2)^2$$

$$(۷) V\{u(Y; \mu)\} = \frac{3\sigma^4}{\mu(1-\mu)} + \frac{\sigma^2}{\mu^3(1-\mu)^3}$$

برهان. با اثبات قسمت (۲) شروع می‌کنیم؛ می‌دانیم که

$$0 = E[(Y - \mu)] = \int_0^1 (y - \mu) f(y; \mu, \sigma^2) dy,$$

حال از طرفین نسبت به μ مشتق گرفته

$$\begin{aligned} 0 &= \frac{d}{d\mu} \int_0^1 (y - \mu) a(y; \sigma^2) \exp \left\{ \frac{-1}{2\sigma^2} d(y; \mu) \right\} dy \\ &= \int_0^1 a(y; \sigma^2) \left[-1 \exp \left\{ \frac{-1}{2\sigma^2} d(y; \mu) \right\} \right] + \int_0^1 a(y; \sigma^2) \left[(y - \mu) \frac{d'(y; \mu)}{2\sigma^2} \exp \left\{ \frac{-1}{2\sigma^2} d(y; \mu) \right\} \right] dy \\ &= -1 - \frac{1}{2\sigma^2} E[(Y - \mu) d'(Y; \mu)] \end{aligned}$$

در نتیجه اثبات کامل می‌شود؛ یعنی $E[(Y - \mu) d'(Y; \mu)] = -2\sigma^2$ برای اثبات قسمت (۴) از قسمت (۳) و رابطه (۵.۳) استفاده کرده و سپس امید ریاضی می‌گیریم.

برای اثبات قسمت (۵) از قسمت‌های (۱)، (۲) و (۳) استفاده می‌کنیم و می‌توان نوشت

$$\begin{aligned} \frac{1}{2} d''(y; \mu) &= \frac{1}{\mu(1-\mu)} d(y; \mu) + \frac{1-2\mu}{\mu^2(1-\mu)^2} (y-\mu) d(y; \mu) \\ &+ \frac{1}{\mu^3(1-\mu)^3} + \frac{1-2\mu}{\mu^4(1-\mu)^4} (y-\mu) \\ &- \frac{1}{\mu(1-\mu)} (y-\mu) d'(y; \mu) - \frac{2(2\mu-1)}{\mu^4(1-\mu)^4} (y-\mu). \end{aligned}$$

حال از رابطه فوق امید ریاضی می‌گیریم و نتیجه می‌شود

$$\begin{aligned} \frac{1}{2} E\{d''(Y; \mu)\} &= \frac{1}{\mu(1-\mu)} [E\{d(Y; \mu)\} - E\{(Y - \mu) d'(Y; \mu)\}] \\ &+ \frac{1-2\mu}{\mu^2(1-\mu)^2} E\{(Y - \mu) d(Y; \mu)\} + \frac{1}{\mu^3(1-\mu)^3} \\ &= \frac{3\sigma^2}{\mu(1-\mu)} + \frac{1}{\mu^3(1-\mu)^3}. \end{aligned}$$

برای اثبات سایر موارد به جرجنسن (۱۹۹۷)، سانگ و تن (۲۰۰۰)، ژو (۲۰۰۱)، سانگ و همکاران (۲۰۰۴)، ژنگ و ژو (۲۰۱۴) رجوع کنید. □

فصل ۳

مدل رگرسیونی سیمپلکس

مقدمه

در این فصل، یک مدل رگرسیونی مناسب را برای مدل‌بندی متغیر پاسخی که به‌طور پیوسته در فاصله واحد، یعنی $y \in (0, 1)$ ، اندازه‌گیری می‌شود، به‌نام رگرسیون سیمپلکس^۱ معرفی می‌کنیم. این مدل رگرسیونی بر اساس پذیره پیروی متغیر پاسخ از توزیع سیمپلکس که در فصل قبل به‌طور مفصل توضیح داده شد، بنا نهاده شده‌است.

از آن‌جایی که توزیع سیمپلکس عضوی از مدل‌های پراکندگی است، روند استنباط برای مدل رگرسیونی سیمپلکس، مشابه مدل‌های پراکندگی و مدل‌های خطی تعمیم‌یافته است. بنابراین، برازش مدل را با روش درست‌نمایی ماکسیمم انجام می‌دهیم. در پایان، با مطالعه شبیه‌سازی، مدل را مورد ارزیابی قرار می‌دهیم.

۱.۳ مدل خطی تعمیم‌یافته سیمپلکس

ما مدل خطی تعمیم‌یافته سیمپلکس^۲ (SGLM) را برای مدل‌بندی متغیرهای وابسته‌ای که به‌طور پیوسته در فاصله واحد مقدار می‌گیرند، معرفی می‌کنیم که مدل بر اساس پذیره پیروی

¹ Simplex Regression

² Simplex Generalized Linear Model

متغیر پاسخ از توزیع سیمپلکس است. فرض بر این است که $(Y_i | x_i \sim S^-(\mu_i, \sigma^2))$ که تابع چگالی احتمال آن در رابطه (۱۱.۲) با میانگین μ_i و پارامتر پراکندگی σ^2 نشان داده شده است. با در نظر گرفتن μ_i به عنوان پارامتر میانگین، پیشگوی خطی مدل به صورت زیر می باشد

$$\eta_i = g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta}$$

این مدل را مدل سیمپلکس همگن^۳ گویند که در آن $g : (0, 1) \rightarrow (-\infty, \infty)$ تابع پیوند بوده و x_i یک بردار p عنصری شامل مشاهده i ام برای p متغیر تبیینی است که $i = 1, 2, \dots, n$ و $\boldsymbol{\beta}$ بردار پارامترهای مجهول رگرسیون است. هدف ما مدل سازی متغیر پاسخ Y بر حسب تابعی از x است. این یک نمونه از مدل های خطی تعمیم یافته است که توسط جرجنسن (۱۹۹۷) معرفی شد؛ در این جا از تابع پیوند لجیت استفاده می کنیم، گرچه سایر تابع های پیوند مانند پروبیت نیز می توانند استفاده شوند.

۱.۱.۳ برازش مدل

برای برآورد پارامترهای $\boldsymbol{\beta}$ در SGLM از روش کمترین توان های دوم بازموزون تکراری^۴ (IRWLS) استفاده می کنیم. رابطه (۱۱.۲) را در نظر بگیرید؛ لگاریتم تابع درستنمایی برای یک نمونه به حجم n به صورت زیر است

$$\begin{aligned} \ell(\boldsymbol{\beta}, \sigma^2) &= \sum_{i=1}^n \ell_i(\beta_i, \sigma^2) \\ &= -\frac{1}{2} \sum_{i=1}^n \left\{ \frac{d(y_i; \mu_i)}{\sigma^2} - \log \sigma^2 + h(y_i) \right\} \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n d(y_i; \mu_i) - \frac{n}{2} \log \sigma^2 - \frac{1}{2} \sum_{i=1}^n h(y_i) \end{aligned}$$

که در آن $h(y_i) = \log \left\{ 2\pi [y_i(1-y_i)]^3 \right\}$. تابع ℓ_i از طریق $\mu_i \in (0, 1)$ توسط تابع لجیت زیر، وابسته به $\boldsymbol{\beta}$ می شود

$$\eta_i = g(\mu_i) = \log \frac{\mu_i}{1 - \mu_i} = \mathbf{x}_i^T \boldsymbol{\beta}.$$

برای یافتن برآورد پارامترهای مجهول، مشتق های جزئی تابع درستنمایی را در نظر بگیرید

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_j} = -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{\partial d(y_i; \mu_i)}{\partial \beta_j}$$

که بنا به قاعده زنجیره ای داریم

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_j} = -\frac{1}{2\sigma^2} \sum_{i=1}^n \frac{\partial d(y_i; \mu_i)}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j}$$

³ Homogeneous Simplex Model

⁴ Iteratively Re-Weighted Least Squares

اکنون هر یک از سه جمله در مجموع فوق را به صورت زیر محاسبه و خلاصه می کنیم

$$\frac{\partial d(y_i; \mu_i)}{\partial \mu_i} = d'(y_i; \mu_i) \quad (1.3)$$

$$\frac{\partial \mu_i}{\partial \eta_i} = \frac{1}{g'(\mu_i)} \quad (2.3)$$

$$\frac{\partial \eta_i}{\partial \beta_j} = \frac{\partial}{\partial \beta_j} \mathbf{x}_i^T \boldsymbol{\beta} = \frac{\partial}{\partial \beta_j} (x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{ij}\beta_j + \dots + x_{ip}\beta_p) = x_{ij} \quad (3.3)$$

در نهایت با توجه به روابط (۱.۳)، (۲.۳) و (۳.۳) داریم

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \beta_j} = -\frac{1}{2\sigma^2} \sum_{i=1}^n d'(y_i; \mu_i) \frac{1}{g'(\mu_i)} x_{ij} = \sum_{i=1}^n c_i x_{ij}$$

و از آن جا که

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \begin{pmatrix} \frac{\partial}{\partial \beta_1} \ell(\boldsymbol{\beta}) \\ \frac{\partial}{\partial \beta_2} \ell(\boldsymbol{\beta}) \\ \vdots \\ \frac{\partial}{\partial \beta_p} \ell(\boldsymbol{\beta}) \end{pmatrix}$$

داریم

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \begin{pmatrix} \sum_{i=1}^n c_i x_{i1} \\ \sum_{i=1}^n c_i x_{i2} \\ \vdots \\ \sum_{i=1}^n c_i x_{ip} \end{pmatrix} = c_1 \begin{pmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1p} \end{pmatrix} + \dots + c_n \begin{pmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{np} \end{pmatrix}.$$

در نتیجه داریم

$$\frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \mathbf{x}_i c_i = -\sum_{i=1}^n \mathbf{x}_i \frac{d'(y_i; \mu_i)}{2\sigma^2 g'(\mu_i)}.$$

برای استفاده از روش IRWLS بایستی تابع اطلاع فیشر را محاسبه کنیم

$$I(\boldsymbol{\beta}) = -E \left(\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \right)$$

$$\begin{aligned} \frac{\partial}{\partial \beta^T} \left(\frac{\partial \ell(\beta)}{\partial \beta} \right) &= \frac{\partial}{\partial \beta^T} \left[- \sum_{i=1}^n \mathbf{x}_i \frac{d'(y_i; \mu_i)}{\sqrt{\sigma^2} g'(\mu_i)} \right] \\ &= - \frac{1}{\sqrt{\sigma^2}} \sum_{i=1}^n \mathbf{x}_i \frac{\partial}{\partial \beta^T} \left[\frac{d'(y_i; \mu_i)}{g'(\mu_i)} \right] \\ &= - \frac{1}{\sqrt{\sigma^2}} \sum_{i=1}^n \mathbf{x}_i \frac{\partial}{\partial \mu_i} \left(\frac{d'(y_i; \mu_i)}{g'(\mu_i)} \right) \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta^T} \\ &= - \frac{1}{\sqrt{\sigma^2}} \sum_{i=1}^n \mathbf{x}_i \left[\frac{d''(y_i; \mu_i) \cdot g'(\mu_i)}{\{g'(\mu_i)\}^2} - \frac{g''(\mu_i) \cdot d'(y_i; \mu_i)}{\{g'(\mu_i)\}^2} \right] \frac{1}{g'(\mu_i)} \mathbf{x}_i^T \end{aligned}$$

بنابراین تابع $I(\beta)$ عبارت است از

$$I(\beta) = \frac{1}{\sqrt{\sigma^2}} \sum_{i=1}^n \mathbf{x}_i \left(\frac{E[d''(y_i; \mu_i)]}{g'(\mu_i)} - \frac{E[d'(y_i; \mu_i)]g''(\mu_i)}{\{g'(\mu_i)\}^2} \right) \frac{1}{g'(\mu_i)} \mathbf{x}_i^T$$

بنا به خاصیت شماره ۴ توزیع سیمپلکس از زیربخش (۳.۳.۲)، $E\{d'(Y; \mu)\} = 0$ و با در نظر گرفتن

$$w_i = \frac{E[d''(y_i; \mu_i)]}{\sqrt{\sigma^2} \{g'(\mu_i)\}^2}$$

در نهایت داریم

$$I(\beta) = \sum_{i=1}^n \mathbf{x}_i w_i \mathbf{x}_i^T \quad (4.3)$$

که با توجه به خاصیت شماره ۵ توزیع سیمپلکس از زیربخش ۳.۳.۲، یعنی

$$\frac{1}{\sqrt{\sigma^2}} E\{d''(Y; \mu)\} = \frac{\sqrt{\sigma^2}}{\mu(1-\mu)} + \frac{1}{\mu^3(1-\mu)^3}$$

$$w_i = \frac{1}{\sigma^2 \{g'(\mu_i)\}^2} \left(\frac{\sqrt{\sigma^2}}{\pi_i} + \frac{1}{\pi_i^3} \right)$$

که در آن $\pi_i = \mu_i(1 - \mu_i)$.

حال تابع امتیاز متناظر با β_j ها را نیز به دست می آوریم

$$\begin{aligned} U(\beta) &= \frac{\partial \ell(\beta)}{\partial \beta} = - \sum_{i=1}^n \mathbf{x}_i^T \frac{d'(y_i; \mu_i)}{\sqrt{\sigma^2} g'(\mu_i)} \\ &= - \sum_{i=1}^n \mathbf{x}_i^T \frac{E[d''(y_i; \mu_i)]}{\sqrt{\sigma^2} \{g'(\mu_i)\}^2} \frac{d'(y_i; \mu_i) g'(\mu_i)}{E[d''(y_i; \mu_i)]} \end{aligned}$$

می دانیم

$$\begin{aligned} u(y; \mu) &= - \frac{1}{\sqrt{\sigma^2}} d'(y; \mu) \\ &= \frac{y - \mu}{\mu(1 - \mu)} \left\{ d(y; \mu) + \frac{1}{\mu^2(1 - \mu)^2} \right\}. \end{aligned} \quad (5.3)$$

$$\begin{aligned} U(\beta) &= \sum_{i=1}^n \mathbf{x}_i^T w_i \frac{\frac{1}{\sigma_i^2} u(y_i; \mu_i) g'(\mu_i)}{\frac{1}{\sigma_i^2} w_i \{g'(\mu_i)\}^2} \\ &= \sum_{i=1}^n \mathbf{x}_i^T w_i \left(\frac{u(y_i; \mu_i)}{\sigma_i^2 w_i g'(\mu_i)} \right) = \sum_{i=1}^n \mathbf{x}_i^T w_i (s_i - g(\mu_i)) \\ &= \sum_{i=1}^n \mathbf{x}_i^T w_i (s_i - \eta_i) \end{aligned} \quad (6.3)$$

که در آن $s_i = g(\mu_i) + \frac{u(y_i, \mu_i)}{\sigma_i^2 w_i g'(\mu_i)}$ این روند، نشان دهنده یک الگوریتم تکراری است، که برآوردگر $(k+1)$ امین β می باشد که به صورت زیر است:

$$\beta^{k+1} = \beta^k + I^{-1}(\beta^k) U(\beta^k) \quad (7.3)$$

حال با توجه به روابط (4.3) و (6.3) در رابطه (7.3) داریم

$$\beta^{(k+1)} = \beta^{(k)} + \left(\sum_{i=1}^m \mathbf{x}_i^T w_i^{(k)} \mathbf{x}_i \right)^{-1} \sum_{i=1}^k \mathbf{x}_i^T w_i^{(k)} (s_i^{(k)} - \eta_i^{(k)}).$$

۲.۱.۳ آماره انحراف

برای مقایسه مجموعه‌ای از مدل‌های رگرسیونی آشیانه‌ای^۵ که متغیر پاسخ آن‌ها دارای توزیع نرمال است، عموماً از روش تحلیل واریانس^۶ استفاده می‌شود که بر اساس مجموع توان دوم مختلف بیان کننده تغییر پذیری داده‌ها هستند. برای GLMs این به تحلیل انحراف تعمیم می‌یابد.

فرض کنید \hat{L}_c و \hat{L}_f به ترتیب نشان دهنده مقدار ماکسیمم تابع درستنمایی برای مدل مورد نظر^۷ و مدل اشباع شده^۸ (مدلی با تمام متغیرهای تبیینی) باشند. آماره انحراف مقیاس بندی شده^۹ برای یک مدل به صورت زیر تعریف می‌شود

$$\begin{aligned} D^* &= -2 \log \left(\frac{\hat{L}_c}{\hat{L}_f} \right) \\ &= -2 \left\{ \log \hat{L}_c - \log \hat{L}_f \right\} \end{aligned} \quad (8.3)$$

که شکل مشابه‌ای با نسبت دو تابع درستنمایی دارد و در احتمال به طور مجانبی دارای توزیع کای دو با $n - p$ درجه آزادی است، که در آن n تعداد کل پارامترها و p تعداد پارامترهای

⁵ Nested

⁶ Analysis of Variance (ANOVA)

⁷ Current Model

⁸ Full Model

⁹ Scaled Deviance

رگرسیونی نامعلومی است که در برازش مدل برآورد شده‌اند. برای داده‌ها با فرض داشتن توزیع نرمال وقتی واریانس معلوم و تابع پیوند همانی باشد این آماره دقیقاً دارای توزیع کای دو است. اما در حالت کلی به دلیل ثابت نبودن تعداد پارامترها، ممکن است تقریب کای دو مفید نباشد. زمانی که یک مدل مناسب برازش شود و تقریب χ^2_{n-p} برقرار باشد، مقدار تقریبی امید ریاضی رابطه (۸.۳) برابر $n - p$ است.

در GLMs با توجه به لگاریتم تابع درست‌نمایی زیر

$$\log L(\beta) = \sum_{i=1}^n \left\{ \frac{1}{\phi} [\theta_i y_i + c(\theta_i)] + d(\phi, y_i) \right\}$$

داریم

$$\log \hat{L}_c = \sum_{i=1}^n \left\{ \frac{1}{\phi} [\hat{\theta}_i y_i + c(\hat{\theta}_i)] + d(\phi, y_i) \right\}$$

که در آن $\hat{\theta}_i$ ، برآورد درست‌نمایی ماکسیمم θ_i است. در مدل اشباع‌شده برآورد $y_i = \mu_i$ است و $\tilde{\theta}_i$ متناظر با برآورد θ_i است؛ حال با استفاده از این که $\mu_i = -c'(\theta_i)$ می‌توان $\tilde{\theta}_i$ را به صورت زیر محاسبه کرد

$$-c'(\tilde{\theta}_i) = y_i$$

سپس

$$\log \hat{L}_f = \sum_{i=1}^n \left\{ \frac{1}{\phi} [\tilde{\theta}_i y_i + c(\tilde{\theta}_i)] + d(\phi, y_i) \right\}$$

و در آخر

$$D^* = \frac{2}{\phi} \sum_{i=1}^n \left\{ (\tilde{\theta}_i - \hat{\theta}_i) y_i + c(\tilde{\theta}_i) - c(\hat{\theta}_i) \right\}$$

حال آماره انحراف را می‌توان به صورت زیر بیان کرد

$$D = \phi D^* \quad (9.3)$$

که در آن D آماره انحراف و ϕ پارامتر پراکندگی است. در حقیقت آماره انحراف برابر با حاصلضرب مقداری ثابت در انحراف مقیاس‌بندی شده است و زمانی که این مقدار ثابت بزرگ‌تر یا کوچک‌تر از یک شود، تمایز بین این دو مقدار نقش مهمی در تحلیل‌های رگرسیونی ایفا می‌کند. از این رو، آماره انحراف مستقل از پارامتر پراکندگی ϕ تعریف می‌شود. جدول ۱.۳ مقدار آماره انحراف مقیاس‌بندی شده را برای برخی از توزیع‌های معروف متعلق به خانواده نمایی نشان می‌دهد. در واقع آماره انحراف، آماره لگاریتم نسبت درست‌نمایی برای مقایسه مدل مورد نظر c ، با مدل اشباع‌شده f است و یا به عبارتی آماره‌ای برای آزمون فرضی است که بیان می‌کند تمام پارامترهایی که در مدل اشباع‌شده وجود دارند و در مدل مورد نظر نیستند، برابر صفر هستند. در حقیقت هدف از محاسبه آماره انحراف تصمیم‌گیری نسبت به سودمندی پارامترهای اضافه

جدول ۱.۳: آماره انحراف مقیاس بندی شده برخی از توزیع های متعلق به خانواده نمایی

توزیع	آماره انحراف مقیاس بندی شده
نرمال	$\sum_i (y_i - \mu_i)^2$
پواسون	$2 \sum_i \left(y_i \ln \left(\frac{y_i}{\mu_i} \right) - (y_i - \mu_i) \right)$
گاما	$2 \sum_i \left(\frac{y_i - \mu_i}{\mu_i} - \ln \left(\frac{y_i}{\mu_i} \right) \right)$
گوسی معکوس	$\sum_i \frac{(y_i - \mu_i)^2}{y_i \mu_i^2}$

شده به مدل مورد نظر است. این آماره برای مدل رگرسیون پواسون همان آماره آزمون نیکویی برازش برای این مدل است که این آزمون به وسیله مقایسه مقادیر مشاهده شده با مقادیر برازش شده، انجام می شود.

۳.۱.۳ برآورد پارامتر پراکندگی

عموما در تحلیل داده ها فرض آن که پارامتر پراکندگی σ^2 معلوم باشد، دور از ذهن است و در اکثر موارد این پارامتر نامعلوم است و محقق ملزم به برآورد آن می باشد. در این حالت پارامتر پراکندگی را می توان توسط روش گشتاوری و با استفاده از عبارت زیر برآورد نمود

$$\widehat{\sigma^2} = \frac{1}{n-p} \sum_{i=1}^n w_i \hat{r}_i \quad (10.3)$$

که در آن نیز n تعداد کل مشاهدات و p تعداد پارامترهای رگرسیونی نامعلومی است که در برازش متداول مدل به کمک باقی مانده های پیرسون به شکل $\hat{r}_i = \frac{Y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i)}}$ برآورد شده اند و w_i وزن معلوم است که محقق برای برخی مقادیر از Y_i یا تمام آن ها، به دلخواه انتخاب می کند.

برای برآورد پارامتر پراکندگی، برآوردگر دیگری مبنی بر آماره انحراف نیز وجود دارد. بر مبنای روابط (۸.۳) و (۹.۳) یک برآورد گشتاوری برای σ^2 به شکل زیر تعریف می شود

$$\widehat{\sigma^2} = \frac{1}{n-p} \sum_{i=1}^n d(y_i; \mu_i) \quad (11.3)$$

این برآوردگر گشتاوری به طور مجانبی نااریب است (کیدیم و فوکیانوز، ۲۰۰۵).

۴.۱.۳ تحلیل باقی مانده‌ها برای مدل رگرسیونی سیمپلکس

در این بخش، به تحلیل باقی مانده‌ها می‌پردازیم و انواع باقی مانده‌ها را معرفی می‌کنیم. هدف از تحلیل باقی مانده‌ها، ارزیابی نیکویی برازش مدل و تشخیص تاثیر متغیرهای تبیینی بر روی متغیر پاسخ است. در حقیقت انحراف يك مقدار مشاهده شده از مقدار برآورد شده آن را باقی مانده می‌نامند. در يك مدل که در آن متغیرهای پاسخ Y به صورت $Y = X(\beta) + \epsilon$ مدل سازی می‌شوند و در آن مولفه‌های ϵ مستقل و هم توزیع، $\epsilon \sim N(0, \sigma^2)$ ، هستند، باقی مانده‌های پاسخ متناظر برای مقادیر $i = 1, \dots, n$ به صورت زیر تعریف می‌شوند:

$$\hat{\epsilon}_i = y_i - \hat{\mu}_i$$

که در آن مقدار برازش شده حاصل از برآورد ML بردار ضرایب رگرسیون β است و ساده ترین نوع باقی مانده است.

یکی دیگر از انواع باقی مانده‌های مورد استفاده، باقی مانده‌های پیرسونی است که فراری و سریباری (۲۰۰۴) آن را باقی مانده‌های معمولی استاندارد شده نامیدند و به صورت زیر بیان می‌شود:

$$r_i^p = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{V}(y_i)}} = \frac{y_i - \hat{\mu}_i}{\hat{\sigma}_i} \quad (12.3)$$

از آن جایی که $V(y_i)$ صورت بسته‌ای ندارد، برای داده‌های بیش پراکنده که پارامتر پراکندگی σ^2 خیلی بزرگ است و واریانس متغیر پاسخ به سمت $\mu(1 - \mu)$ میل می‌کند، باقی مانده پیرسونی مجانبی برابر است با

$$r_i^a = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\mu}_i(1 - \hat{\mu}_i)}} \quad (13.3)$$

حال برای ارزیابی نیکویی برازش مدل مد نظر با استفاده از باقی مانده‌ها، اگر مدل درست انتخاب شده باشد، آن گاه باقی مانده‌ها باید دارای توزیع نرمال استاندارد مجانبی باشند و نقاط در نمودار پراکنش مقادیر برآورد شده در مقابل باقی مانده‌ها، بدون نمایش هیچ الگوی منظمی عمدتاً در بازه $(-2, 2)$ قرار دارند.

۵.۱.۳ مطالعه شبیه سازی

در این بخش، ویژگی‌های برآوردگرهای پارامترهای مدل رگرسیونی سیمپلکس که در این فصل، به طور نظری بیان کردیم، برای حجم‌های نمونه متفاوتی مورد ارزیابی قرار می‌دهیم. شبیه سازی‌ها در محیط نرم افزار R و با کمک بسته `simplexreg` اجرا شده‌اند. نحوه برازش یک مدل رگرسیونی سیمپلکس با استفاده از بسته `simplexreg` را در فصل ۵ به طور کامل تشریح خواهیم کرد.

مثال ۱.۱.۳. مثال اول

برای مثال اول شبیه سازی، در مدل (۱۱.۲)، فرض می کنیم پیشگوی خطی به صورت

$$\text{logit}(\mu_i) = \beta_0 + \beta_1 x_{1i}, \quad i = 1, \dots, n. \quad (14.3)$$

باشد، که در آن $(\beta_0, \beta_1) = (1/5, 0/5)$ و پارامتر پراکندگی $\sigma^2 = 2$ انتخاب شدند. مقادیر متغیر تبیینی x_1 را از توزیع برنولی با احتمال موفقیت $0/6$ تولید کردیم؛ هم چنین برای بررسی عملکرد روش، پنج حجم نمونه $n = 30, 50, 100, 200, 400$ را در نظر گرفتیم. در جدول ۲.۳ برآوردهای نهایی هر پارامتر که با میانگین گیری از برآوردهای حاصله در ۵۰۰ تکرار، محاسبه شده، نشان داده شده اند.

جدول ۲.۳: خلاصه نتایج شبیه سازی برای برآورد ضرایب مدل (۱۱.۲) با پیشگوی خطی (۱۴.۳)

حجم نمونه	پارامترها	مقدار واقعی	برآورد	Stdev	MSE	اریبی
۳۰	β_0	۱/۵	۱/۵۱۴۰	۰/۱۶۳۴	۰/۰۲۶۸	۰/۰۱۴
	β_1	۰/۵	۰/۴۷۴۲	۰/۱۹۶۳	۰/۰۳۹	-۰/۰۲۵۷
	σ^2	۲	۱/۶۰۹۶	۰/۲۶۲۸	۰/۲۲۱۲	-۰/۳۹۰۳
۵۰	β_0	۱/۵	۱/۵۱۳۱	۰/۱۱۳۶	۰/۰۱۳	۰/۰۱۳۱
	β_1	۰/۵	۰/۴۷۲۸	۰/۱۳۹۳	۰/۰۲	-۰/۰۲۷۱
	σ^2	۲	۱/۷۰۸۲	۰/۲۳۰۱	۰/۱۳۷۹	-۰/۲۹۱۷
۱۰۰	β_0	۱/۵	۱/۵۰۵۷	۰/۰۷۶۲	۰/۰۰۵۸	۰/۰۰۵۷
	β_1	۰/۵	۰/۴۸۵۸	۰/۰۹۳۵	۰/۰۰۸۹	-۰/۰۱۴۱
	σ^2	۲	۱/۷۹۳۲	۰/۱۶۲۶	۰/۰۰۶۹	-۰/۲۰۶۷
۲۰۰	β_0	۱/۵	۱/۵۰۲۲	۰/۰۵۹۷	۰/۰۰۳۵	۰/۰۰۲۲
	β_1	۰/۵	۰/۴۹۳۱	۰/۰۷۳۲	۰/۰۰۵۳	-۰/۰۰۶۸
	σ^2	۲	۱/۸۳۹۸	۰/۱۱۳۶	۰/۰۳۸۵	-۰/۱۶۰۱
۴۰۰	β_0	۱/۵	۱/۵۰۱۱	۰/۰۴۰۳	۰/۰۰۱۶	۰/۰۰۱۱
	β_1	۰/۵	۰/۴۹۷۱	۰/۰۴۸۷	۰/۰۰۲۳	-۰/۰۰۲۸
	σ^2	۲	۱/۸۸۲۹	۰/۰۸۶۲	۰/۰۲۱۱	-۰/۱۱۷

همان گونه که مشخص است میانگین مقادیر برآورد پارامترها، بسیار نزدیک به مقادیر واقعی است و با افزایش حجم نمونه، دقت این برآوردها به واسطه کاهش خطای معیار، افزایش می یابد. هم چنین با افزایش حجم نمونه، MSE برآوردها در حال کاهش و همگرا به صفر است و مقادیر اریبی برآوردها نیز ناچیز است. بنابراین، می توان سازگاری برآوردها را برای حجم های نمونه کوچک و متوسط نیز نتیجه گرفت.

مثال ۲.۱.۳. مثال دوم

برای شبیه‌سازی در مدل (۱۱.۲)، فرض می‌کنیم که پیشگوی خطی به صورت

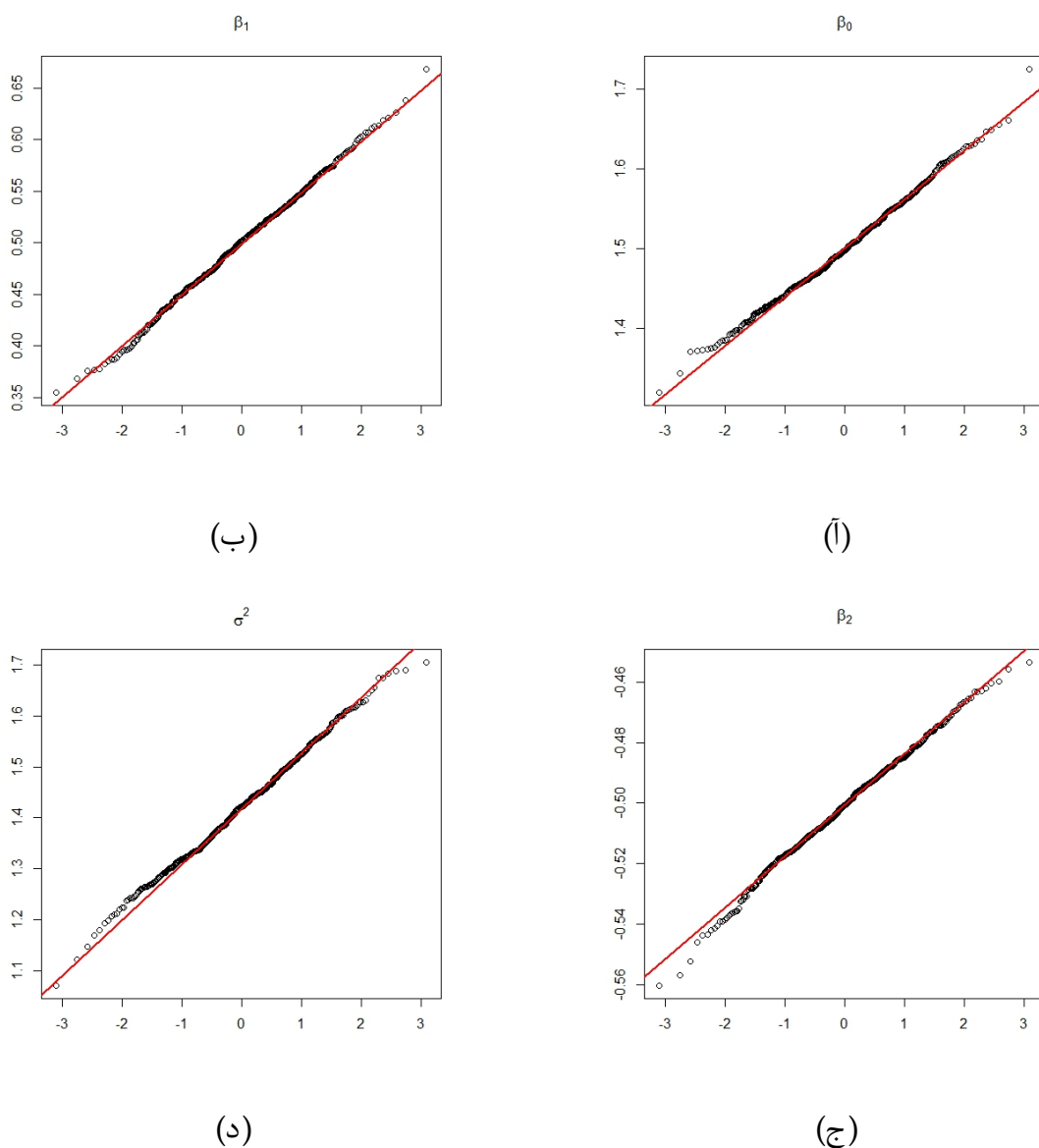
$$\text{logit}(\mu_i) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}, \quad i = 1, \dots, n. \quad (15.3)$$

باشد، که در آن $(\beta_0, \beta_1, \beta_2) = (1/5, 0/5, -0/5)$ و پارامتر پراکندگی $\sigma^2 = 2$ انتخاب شدند. مقادیر متغیر تبیینی x_1 را از توزیع برنولی با احتمال موفقیت $0/6$ و مقادیر x_2 را از توزیع $U(0, 5)$ تولید کردیم. هم‌چنین برای بررسی عملکرد روش، پنج حجم نمونه $300, 500, 1000, 2000, 4000$ را در نظر گرفتیم. در جدول ۳.۳ برآوردهای نهایی هر پارامتر که با میانگین‌گیری از برآوردهای حاصله در 500 تکرار، محاسبه شده، نشان داده شده‌اند.

جدول ۳.۳: خلاصه نتایج شبیه‌سازی برای برآورد ضرایب مدل (۱۱.۲) با پیشگوی خطی (۱۵.۳)

حجم نمونه	پارامترها	مقدار واقعی	برآورد	Stdev	MSE	اریبی
۳۰	β_0	۱/۵	۱/۳۸۴	۰/۲۱۸۷	۰/۰۶۰۷	-۰/۱۱۶
	β_1	۰/۵	۰/۴۷۶۴	۰/۱۷۵۱	۰/۰۳۰۸	-۰/۰۲۳۶
	β_2	-۰/۵	-۰/۴۵۳۷	۰/۰۶۳	۰/۰۰۰۶	۰/۰۴۶۳
	σ^2	۲	۱/۳۷۷۸	۰/۲۸	۰/۰۹۵۷	-۰/۰۶۲۲۲
۵۰	β_0	۱/۵	۱/۴۱۳۴	۰/۱۷۶۷	۰/۰۳۸۱	-۰/۰۸۶۶
	β_1	۰/۵	۰/۴۶۳۲	۰/۱۶۲۲	۰/۰۲۷۱	-۰/۰۳۶۸
	β_2	-۰/۵	-۰/۴۶۲۸	۰/۰۵۴۷	۰/۰۰۴۳	۰/۰۳۷۲
	σ^2	۲	۱/۴۶۶۵	۰/۱۷۵۷	۱/۱۵۹۶	-۰/۰۵۳۳۵
۱۰۰	β_0	۱/۵	۱/۴۰۸۷	۰/۱۴۶۳	۰/۰۲۸۶	-۰/۰۹۱۳
	β_1	۰/۵	۰/۴۸۷۴	۰/۱۳۴۴	۰/۰۱۷۳	-۰/۰۱۲۶
	β_2	-۰/۵	-۰/۴۵۷۸	۰/۰۴۶۵	۰/۰۰۳۸	۰/۰۴۲۲
	σ^2	۲	۱/۵۴۲	۰/۱۵۱۳	۰/۰۵۷۲۵	-۰/۰۴۵۸
۲۰۰	β_0	۱/۵	۱/۴۲۵	۰/۰۲۵۵	۰/۰۰۵۹	-۰/۰۰۷۵
	β_1	۰/۵	۰/۳۹۲۲	۰/۰۲۷	۰/۰۱۱۹	-۰/۰۰۷۸
	β_2	-۰/۵	-۰/۴۴۹۵	۰/۰۰۶۳	۰/۰۰۲۵	۰/۰۰۵۰۵
	σ^2	۲	۱/۶۹۵	۰/۰۲۲۹	۰/۰۲۴۸	-۰/۰۰۳۰۵
۴۰۰	β_0	۱/۵	۱/۵۰۱۶	۰/۰۶۸۳	۰/۰۰۴۶	۰/۰۰۱۶
	β_1	۰/۵	۰/۴۹۸۷	۰/۰۵۸۲	۰/۰۰۳۳	-۰/۰۰۰۱۳
	β_2	-۰/۵	-۰/۵۰۱۵	۰/۰۲	۰/۰۰۰۴	-۰/۰۰۰۱۵
	σ^2	۲	۲/۰۰۶۹	۰/۱۴۶۵	۰/۰۰۶۹	۰/۰۰۰۶۹

روند تغییرات انحراف معیار برآوردگرها، مشابه مثال اول هستند و در این جا نیز می‌توان سازگاری برآوردگرها را نتیجه گرفت. شکل ۱.۳، نمودارهای چندک-چندک پارامترهای برآورد شده در ۵۰۰ مجموعه داده شبیه‌سازی شده را برای پیشگوی خطی (۱۵.۳) نمایش می‌دهد. با توجه به این نمودارها، می‌توان نرمال مجانبی برآوردگرها را نتیجه گرفت، زیرا توزیع نرمال تقریب مناسبی برای توزیع برآوردگرها است. با توجه به نتایج به دست آمده می‌توان گفت که روش پیشنهادی قابل اعتماد است.



شکل ۱.۳: نمودار چندک-چندک پارامترهای مدل (۱۱.۲)، با پیشگوی خطی (۱۵.۳)

۲.۳ مدل رگرسیونی سیمپلکس با پارامتر پراکندگی متغیر

در بخش قبل مدل رگرسیونی سیمپلکس را برای حالتی که پارامتر پراکندگی، σ^2 ، ثابت بود، مورد بررسی قرار دادیم.

در این بخش، مدل را برای حالتی که پارامتر σ_i^2 متغیر و خود تابعی از متغیرهای تبیینی است، تعمیم می‌دهیم. هم‌چنین به برازش مدل با استفاده از روش درست‌نمایی ماکسیمم می‌پردازیم.

۱.۲.۳ معرفی مدل

همان‌طور که در فصل قبل بیان کردیم، متغیر پاسخ y دارای توزیع سیمپلکس به صورت (۱۱.۲) است. اکنون فرض می‌کنیم پارامتر پراکندگی σ^2 نیز تابعی از مشاهدات است، به طوری که متغیرهای تبیینی $z_i, i = 1, \dots, m$ که می‌توانند زیرمجموعه‌ای از x_i ها است، بر روی آن تاثیر دارند. فرض کنید میانگین و پارامتر پراکندگی برای y_i به صورت زیر مدل بندی شوند:

$$g(\mu_i) = \eta_{1i} = \mathbf{x}_i^T \boldsymbol{\beta}, \quad h(\sigma_i^2) = \eta_{2i} = \mathbf{z}_i^T \boldsymbol{\gamma}. \quad (۱۶.۳)$$

که در آن $g : (0, 1) \rightarrow (-\infty, \infty)$ و $h : (0, \infty) \rightarrow (-\infty, \infty)$. این مدل را سیمپلکس ناهمگن^{۱۰} گویند که در آن هر دو تابع پیوند، اکیدا یکنوا و مشتق‌پذیر می‌باشند؛ $\boldsymbol{\beta}$ و $\boldsymbol{\gamma}$ نیز بردارهای پارامترهای نامعلوم برای توابع میانگین و پراکندگی است؛ هم‌چنین η_{1i} و η_{2i} پیشگوه‌های خطی اند و x_i و z_i متغیرهای تبیینی معلوم می‌باشند.

تابع پیوند $g(\cdot)$ همان تابع پیوندی است که در فصل ۲ معرفی شده، $g(\mu) = \log \frac{\mu}{1-\mu}$ ، اما تابع پیوندی که برای $h(\cdot)$ در نظر گرفته می‌شود، تابع پیوند لگاریتمی است، یعنی $h(\sigma_i^2) = \log \sigma_i^2$.

۲.۲.۳ برازش مدل

تابع لگاریتم درست‌نمایی برای یک نمونه به حجم n ، برای این مدل رگرسیون سیمپلکس به صورت زیر است

$$\ell(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \sum_{i=1}^n \ell_i(\mu_i, \sigma_i^2) \quad (۱۷.۳)$$

که در آن

$$\ell_i(\boldsymbol{\beta}_i, \boldsymbol{\gamma}_i) = -\frac{1}{2} \sum_{i=1}^m \left\{ \frac{d(y_i; \mu_i)}{\sigma_i^2} - \log \sigma_i^2 \right\}.$$

توابع $\mu_i = g^{-1}(\eta_{1i})$ و $\sigma_i^2 = h^{-1}(\eta_{2i})$ که در رابطه (۱۶.۳) تعریف شده‌اند، به ترتیب توابعی از $\boldsymbol{\beta}$ و $\boldsymbol{\gamma}$ می‌باشند. برای برآورد پارامترها، مشابه زیربخش ۱.۱.۳، باید تابع امتیاز متناظر با $\boldsymbol{\gamma}$ را

¹⁰ Simplex Heterogeneous Model

نیز به دست آوریم.

با استفاده از قاعده زنجیره‌ای، برای محاسبه تابع امتیاز U ، با مشتق‌گیری از تابع لگاریتم درستنمایی (۱۷.۳)، نسبت به γ_j ، $j = 1, \dots, n$ داریم

$$U(\gamma) = \frac{\partial \ell(\beta, \gamma)}{\partial \gamma_j} = -\frac{1}{2} \sum_{i=1}^n \left\{ \frac{\partial d(y_i; \mu_i) - \log(\sigma_i^2)}{\partial \sigma_i^2} \frac{\partial \sigma_i^2}{d\eta_{\gamma_i}} \frac{d\eta_{\gamma_i}}{\partial \gamma_j} \right\}. \quad (18.3)$$

بنابراین برای محاسبه (۱۸.۳)، باید مشتق‌های جزئی مجموع سمت راست آن را به ترتیب محاسبه و جایگذاری کنیم (ژانگ و ژو، ۲۰۱۴). می‌توان نوشت

$$\frac{\partial d(y_i; \mu_i)}{\partial \sigma_i^2} = \left(\frac{-d(y_i; \mu_i)}{(\sigma_i^2)^2} + \frac{1}{\sigma_i^2} \right) \quad (19.3)$$

حال با توجه به این که $\sigma_i^2 = h^{-1}(\eta_{\gamma_i})$ داریم

$$\frac{\partial \sigma_i^2}{d\eta_{\gamma_i}} = \frac{1}{h'(\sigma_i^2)}, \quad (20.3)$$

و

$$\frac{d\eta_{\gamma_i}}{\partial \gamma_j} = z_{ij} \quad (21.3)$$

در نهایت با جایگذاری (۱۹.۳)، (۲۰.۳) و (۲۱.۳) در رابطه (۱۸.۳)، نتیجه می‌شود

$$U(\gamma) = -\frac{1}{2} \sum_{i=1}^n \left(\frac{-d(y_i; \mu_i)}{(\sigma_i^2)^2} + \frac{1}{\sigma_i^2} \right) \frac{1}{h'(\sigma_i^2)} z_{ij}$$

قرار می‌دهیم

$$v_i = \frac{d(y_i; \mu_i) - \sigma_i^2}{2(\sigma_i^2)^2 h'(\sigma_i^2)}, \quad i = 1, \dots, n.$$

حال می‌توان رابطه را به صورت زیر بازنویسی کرد

$$U(\gamma) = \sum_{i=1}^n z_i^T v_i \quad (22.3)$$

برای حل معادله‌های تابع درستنمایی، باید از مشتق دوم تابع لگاریتم درستنمایی، امید ریاضی بگیریم. به همین منظور می‌توان نوشت:

$$I(\gamma) = -E \left(\frac{\partial^2 \ell(\gamma)}{\partial \gamma \partial \gamma^T} \right).$$

ابتدا عبارت داخل امید ریاضی را محاسبه کرده سپس امید ریاضی می‌گیریم.

$$\begin{aligned} \frac{\partial^2 \ell(\gamma)}{\partial \gamma \partial \gamma^T} &= \frac{\partial}{\partial \gamma} \left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{-d(y_i; \mu_i)}{(\sigma_i^2)^2} + \frac{1}{\sigma_i^2} \right) \frac{1}{h'(\sigma_i^2)} \mathbf{z}_i \right] \\ &= -\frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \sigma_i^2} \left[\frac{1}{h'(\sigma_i^2)} \left(\frac{-d(y_i; \mu_i)}{(\sigma_i^2)^2} + \frac{1}{\sigma_i^2} \right) \right] \frac{\partial \sigma_i^2}{\partial \eta_{\gamma_i}} \frac{\partial \eta_{\gamma_i}}{\partial \gamma^T} \mathbf{z}_i \end{aligned}$$

با استفاده از قاعده زنجیره‌ای فرمول بالا را باز نویسی می‌کنیم

$$\begin{aligned} \frac{\partial^2 \ell(\gamma)}{\partial \gamma \partial \gamma^T} &= -\frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \sigma_i^2} \left[\frac{-d(y_i; \mu_i)}{h'(\sigma_i^2)(\sigma_i^2)^2} + \frac{1}{h'(\sigma_i^2)\sigma_i^2} \right] \frac{1}{h'(\sigma_i^2)} \mathbf{z}_i^T \mathbf{z}_i \\ &= -\frac{1}{2} \sum_{i=1}^n \left[\frac{(h''(\sigma_i^2)(\sigma_i^2)^2 + 2\sigma_i^2 h'(\sigma_i^2))d(y_i; \mu_i)}{h'^3(\sigma_i^2)(\sigma_i^2)^4} - \frac{h''(\sigma_i^2)\sigma_i^2 + h'(\sigma_i^2)}{h'^3(\sigma_i^2)(\sigma_i^2)^2} \right] \frac{1}{h'(\sigma_i^2)} \mathbf{z}_i^T \mathbf{z}_i \\ &= -\sum_{i=1}^n \left[\frac{\mathbf{z}_i^T \mathbf{z}_i (h''(\sigma_i^2)\sigma_i^2 d(y_i; \mu_i) + 2h'(\sigma_i^2)d(y_i; \mu_i))}{2h'^3(\sigma_i^2)(\sigma_i^2)^3} \right] \\ &\quad + \sum_{i=1}^n \left[\frac{\mathbf{z}_i^T \mathbf{z}_i (h''(\sigma_i^2)\sigma_i^2 + h'(\sigma_i^2))}{2h'^3(\sigma_i^2)(\sigma_i^2)^2} \right] \\ &= -\sum_{i=1}^n \left[\frac{\mathbf{z}_i^T \mathbf{z}_i h''(\sigma_i^2)\sigma_i^2 d(y_i; \mu_i)}{2h'^3(\sigma_i^2)(\sigma_i^2)^3} \right] - \sum_{i=1}^n \left[\frac{2\mathbf{z}_i^T \mathbf{z}_i d(y_i; \mu_i)}{2h'^3(\sigma_i^2)(\sigma_i^2)^3} \right] + \sum_{i=1}^n \left[\frac{\mathbf{z}_i^T \mathbf{z}_i h''(\sigma_i^2)}{2h'^3(\sigma_i^2)\sigma_i^2} \right] \\ &\quad + \sum_{i=1}^n \left[\frac{\mathbf{z}_i^T \mathbf{z}_i}{2h'^3(\sigma_i^2)(\sigma_i^2)^2} \right] \\ &= \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T h''(\sigma_i^2)}{2h'^3(\sigma_i^2)} \left(-\frac{d(y_i; \mu_i)}{(\sigma_i^2)^2} + \frac{1}{\sigma_i^2} \right) - \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{2h'^3(\sigma_i^2)} \left(\frac{2d(y_i; \mu_i)}{(\sigma_i^2)^3} - \frac{1}{(\sigma_i^2)^2} \right) \end{aligned}$$

حال از طرفین امید ریاضی می‌گیریم؛ با توجه به ویژگی ۴ در زیربخش ۳.۳.۲ داریم

$$E\left[-\frac{\partial^2 \ell(\gamma)}{\partial \gamma \partial \gamma^T}\right] = \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{2h'^3(\sigma_i^2)} \left(\frac{1}{(\sigma_i^2)^2} \right)$$

می‌توان رابطه بالا را به صورت زیر بازنویسی می‌کنیم

$$= E\left[-\frac{\partial^2 \ell(\gamma)}{\partial \gamma \partial \gamma^T}\right] = \sum_{i=1}^n \mathbf{z}_i^T q_i \mathbf{z}_i \quad (23.3)$$

که در آن $q_i = \frac{1}{2(\sigma_i^2)^2 h'^3(\sigma_i^2)}$. این روند نشان‌دهنده یک الگوریتم تکراری است، که برآوردگر $(k+1)$ امین γ با توجه به رابطه (۷.۳) و روابط (۲۲.۳) و (۲۳.۳) به صورت زیر است:

$$\gamma^{(k+1)} = \gamma^{(k)} + \left(\sum_{i=1}^n \mathbf{z}_i^T q_i \mathbf{z}_i \right)^{-1} \sum_{i=1}^n \mathbf{z}_i^T v_i^{(k)}$$

۳.۲.۳ مطالعه شبیه‌سازی

روند کار به همان صورتی است که در بخش قبل انجام شد، با این تفاوت که در این بخش پارامتر پراکندگی مدل به صورت یک تابع است. یعنی باید برای پارامتر پراکندگی، همانند پارامتر میانگین، μ ، عمل کنیم. در نتیجه مطالعه شبیه‌سازی را با استفاده از تابع پیوند لجیت برای میانگین و تابع پیوند لگاریتمی برای پارامتر پراکندگی، انجام دادیم. مدل مورد نظر را به صورت

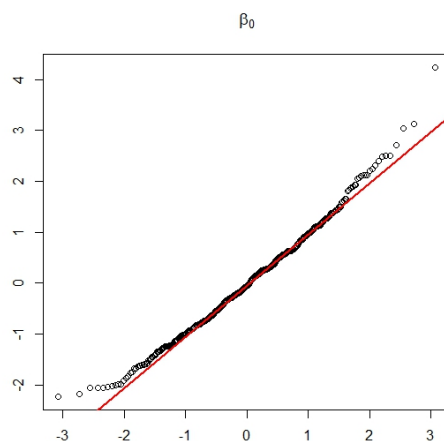
$$\text{logit}(\mu_i) = \eta_{1i} = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}, \quad \log(\sigma_i^2) = \eta_{2i} = \gamma_0 + \gamma_1 z_{1i}, \quad i = 1, \dots, n \quad (24.3)$$

در نظر گرفتیم، که مقادیر واقعی ضرایب عبارتند از $(\beta_0, \beta_1, \beta_2) = (1, 0.5, -1)$ و $(\gamma_0, \gamma_1) = (5, 0.3)$. مقادیر متغیرهای تبیینی x_1 و x_2 را به ترتیب از توزیع $N(0, 1)$ و توزیع برنولی با احتمال موفقیت ۰/۶، تولید کردیم. مقادیر متغیر تبیینی z_1 را نیز از توزیع $N(0, 1)$ تولید کردیم. همچنین برای بررسی عملکرد روش، پنج حجم نمونه $n = 30, 50, 100, 200, 400$ را در نظر گرفتیم. در جدول ۴.۳ برآوردهای نهایی پارامترها که با میانگین‌گیری از برآوردهای حاصله در ۵۰۰ تکرار، محاسبه شده نشان داده شده‌اند که با افزایش حجم نمونه MSE کاهش یافته و همچنین آریبی برآوردها ناچیز است.

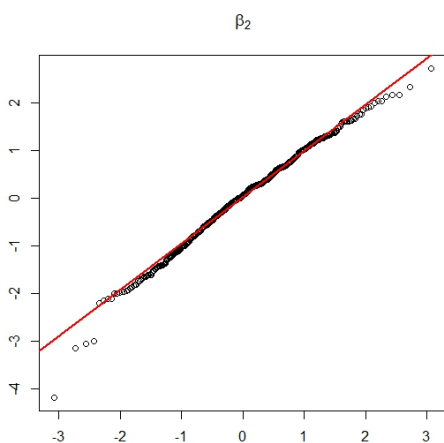
شکل ۲.۳ نیز منحنی چگالی مقادیر برآورده شده از مجموعه داده شبیه‌سازی شده را برای ضرایب مدل (۲۴.۳) با $n = 400$ نمایش می‌دهد. این شکل نیز تایید دیگری بر خوب بودن تقریب نرمال برای توزیع برآوردهای پارامتر مدل است. با توجه به نتایج به دست آمده، می‌توان عملکرد مناسب روش برازش پیشنهادی را نتیجه گرفت.

جدول ۴.۳: خلاصه نتایج شبیه‌سازی برای برآورد ضرایب مدل (۱۱.۲) با پیشگوی خطی (۲۴.۳)

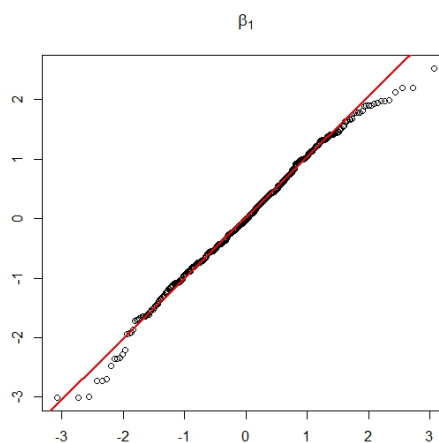
حجم نمونه	پارامترها	مقدار واقعی	میانگین برآورد	Stdev	MSE	اریبی
۳۰	β_0	۱	۱/۰۶۹۳	۰/۴۴۲۸	۰/۲۰۰۵	۰/۰۶۹۳
	β_1	۰/۵	۰/۵۳۲	۰/۳۱۴۳	۰/۰۹۹۶	۰/۰۳۲
	β_2	-۱	-۱/۰۴۷۶	۰/۵۳۸۵	۰/۲۹۱۶	-۰/۰۴۷۶
	γ_0	۵	۴/۸۱۰۶	۰/۲۷۲۹	۰/۱۱۰۲	-۰/۱۸۹۴
	γ_1	۰/۳	۰/۲۷۳۳	۰/۳۲۱۵	۰/۱۰۳۸	-۰/۰۲۶۷
۵۰	β_0	۱	۱/۰۵۳۱	۰/۳۳۳۳	۰/۱۱۳۷	۰/۰۵۳۱
	β_1	۰/۵	۰/۵۵۶	۰/۲۳۲۳	۰/۰۵۷	۰/۰۵۶
	β_2	-۱	-۱/۰۵۹۹	۰/۴۲۴۹	۰/۱۸۳۷	-۰/۰۵۹۹
	γ_0	۵	۴/۸۶۹۳	۰/۲۰۸۱	۰/۰۶۰۲	-۰/۱۳۰۶
	γ_1	۰/۳	۰/۲۹۸۹	۰/۲۱۱۴	۰/۰۴۴۶	-۰/۰۰۱
۱۰۰	β_0	۱	۱/۰۶۲	۰/۲۶۴۴	۰/۰۷۳۶	۰/۰۶۲
	β_1	۰/۵	۰/۵۲۶۵	۰/۱۴۸۹	۰/۰۲۲۸	۰/۰۲۶۵
	β_2	-۱	-۱/۰۵۸۹	۰/۳۰۷۷	۰/۰۹۷۹	-۰/۰۵۸۹
	γ_0	۵	۴/۹۴۰۳	۰/۱۴۳۴	۰/۰۲۴	-۰/۰۵۹۶
	γ_1	۰/۳	۰/۲۸۴۵	۰/۲۸۴۵	۰/۰۲۲۵	-۰/۰۱۵۴
۲۰۰	β_0	۱	۱/۰۰۷۹	۰/۱۵۷۳	۰/۰۲۴۷	۰/۰۰۷۹
	β_1	۰/۵	۰/۵۱۴۹	۰/۰۹۳۵	۰/۰۰۸۹	۰/۰۱۴۹
	β_2	-۱	-۱/۰۰۶	۰/۱۹۸۴	۰/۰۳۹۳	-۰/۰۰۶
	γ_0	۵	۴/۹۷۰۶	۰/۰۹۹	۰/۰۱۰۶	-۰/۰۲۹۳
	γ_1	۰/۳	۳/۰۲۶	۰/۱۰۵۵	۰/۰۱۱۱	۰/۰۰۲۶
۴۰۰	β_0	۱	۱/۰۰۱۳	۰/۱۰۶۸	۰/۰۱۱۴	۰/۰۰۱۳
	β_1	۰/۵	۰/۵۰۷۹	۰/۰۶۶۵	۰/۰۰۴۴	۰/۰۰۷۹
	β_2	-۱	-۰/۹۹۳	۰/۱۳۱۶	۰/۰۱۷۳	۰/۰۰۶۹
	γ_0	۵	۴/۹۸۵۹	۰/۰۷۳۷	۰/۰۰۵۶	-۰/۰۱۴۱
	γ_1	۰/۳	۰/۲۹۳۶	۰/۰۶۸	۰/۰۰۴۶	-۰/۰۰۶۳



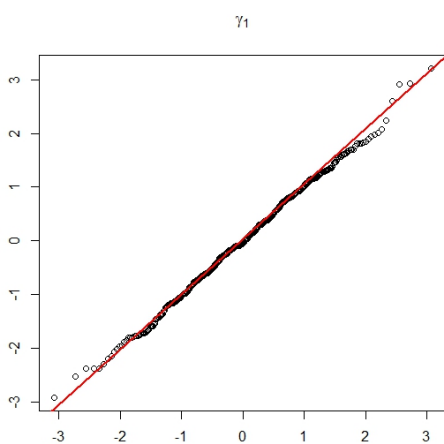
(ا)



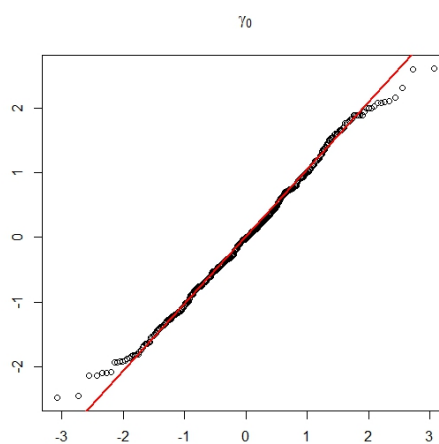
(ج)



(ب)



(ه)



(د)

شکل ۲.۳: نمودار چندک-چندک برآوردهای ضرایب رگرسیونی میانگین و پارامتر پراکندگی در مدل (۲۴.۳) برای $n = 400$ و تعداد تکرار 500 .

فصل ۴

مدل حاشیه‌ای برای تحلیل داده‌های نسبت طولی

مقدمه

امروزه استفاده از مطالعات طولی در بسیاری از علوم، نظیر علوم اجتماعی و اقتصادی به خوبی رواج یافته است. مطالعات طولی در قیاس با مطالعات مقطعی و سری زمانی، به دلیل این که واحدهای مورد نظر را در طول زمان بررسی می‌نمایند، از قابلیت بیشتری جهت مطالعه پدیده‌های مختلف برخوردارند.

در موقعیت‌هایی که با اندازه‌های تکراری در داده‌ها مواجه هستیم، به دلیل وجود همبستگی برای هر واحد بین اندازه‌ها، نمی‌توانیم از مدل‌های خطی استفاده نماییم و بایستی معادلات برآوردیابی دیگری را به کار ببریم، زیرا در یک مدل رگرسیون خطی کلاسیک پذیره‌های استقلال مشاهدات، ثابت بودن واریانس و وجود ارتباط خطی بین مشاهدات و متغیرهای تصادفی از الزامات تعریف مدل، برای صحیح بودن نتایج تحلیل است و در GLM تنها دانستن واریانس مشاهدات، به عنوان تابعی از میانگین و وابسته بودن میانگین مشاهدات به تابعی خطی یا غیرخطی از متغیرهای تبیینی توسط تابع پیوند، کافی است. در حالی که استفاده از تحلیل‌هایی بدون در نظر گرفتن همبستگی موجب صحیح نبودن برآورد مولفه‌های واریانس و نیز استنباط اشتباه نسبت به بردار ضرایب رگرسیونی می‌شود. به همین دلیل در فصل ۱ به معرفی انواع

مدل‌های خطی و GLM و مباحث مربوط به آن‌ها پرداختیم و به منظور بررسی و تحلیل این نوع داده‌های تکراری که همبستگی بین مشاهدات آن‌ها وجود دارد، روش شبه‌درست‌نمایی^۱ بر مبنای GLM با عنوان GEEs در بخش‌های بعدی معرفی می‌شود که این روش توسط GLM برآوردهایی کارا و ناریب برای پارامترهای مدل ارائه می‌دهد.

۱.۴ معرفی داده‌های طولی

داده‌های طولی که در اقتصاد و جامعه‌شناسی با عنوان داده‌های پانلی^۲ شناخته می‌شوند، به مجموعه داده‌هایی گفته می‌شود که در زمان‌های متوالی از واحدهای مشخص به دست می‌آیند. بنابراین این داده‌ها دارای دو بعد هستند: بعد مقطعی (واحدها) و بعد زمانی. در نتیجه داده‌های مقطعی و داده‌های سری زمانی حالت‌های خاصی از داده‌های طولی هستند.

۱.۱.۴ ساختار داده‌های طولی

مطالعات طولی، واحدهای مشخصی را در زمان‌های متوالی مورد بررسی قرار داده و در آن متغیر پاسخ و متغیرهای تبیینی اندازه‌گیری می‌شوند. چنین داده‌هایی را می‌توان به صورت جدول ۱.۴ نمایش داد: در این جدول y_{it} مقدار متغیر پاسخ مربوط به واحد i ام در زمان t ام

جدول ۱.۴: ساختار داده‌ها در مطالعات طولی

	واحد				
	۱	...	i	...	n
	$(y_{i1}, x_{i11}, \dots, x_{i1k})$...	$(y_{i1}, x_{i11}, \dots, x_{i1k})$...	$(y_{in}, \dots, x_{in1}, x_{ink})$
	⋮	...	⋮	...	⋮
زمان t	$(y_{t1}, x_{t11}, \dots, x_{t1k})$...	$(y_{it}, x_{it1}, \dots, x_{itk})$...	$(y_{tn}, \dots, x_{tn1}, x_{tnk})$
	⋮	...	⋮	...	⋮
	$(y_{T1}, x_{T11}, \dots, x_{T1k})$...	$(y_{iT}, x_{iT1}, \dots, x_{iTk})$...	$(y_{nT}, \dots, x_{nT1}, x_{nTk})$

و x_{itr} مقدار r امین متغیر تبیینی مربوط به زمان t ام و واحد i ام می‌باشد. T و n نیز به ترتیب

^۱ Quasi-Likelihood

^۲ Panel Data

معرف تعداد دوره‌های زمانی و تعداد واحدهایی هستند که بررسی شوند. جدول ۱.۴ مربوط به یک مطالعه طولی متعادل می‌باشد زیرا تمام واحدها به‌طور مساوی و در T دوره زمانی، مورد ارزیابی قرار گرفته‌اند. همچنین مواردی وجود دارند که در آن واحد i ام در T_i دوره زمانی اندازه‌گیری می‌شود. چنین مطالعه‌ای را نامتعادل گویند.

در جدول ۱.۴ سطرها (زمان) ماهیت مقطعی داده‌ها و ستون‌ها (واحدها) ماهیت زمانی داده‌ها را نشان می‌دهند. به‌عنوان مثال، اگر داده‌هایی را که در زمان t ام جمع‌آوری شده‌اند در نظر بگیریم در واقع تمام n واحد را تنها در یک مقطع از زمان مورد بررسی قرار داده‌ایم. بنابراین برای هر دوره زمانی، داده‌ها مقطعی می‌باشند. اما اگر مثلاً داده‌های مربوط به واحد i ام را در نظر بگیریم، تنها یک واحد را در زمان‌های متوالی مورد بررسی قرار داده‌ایم؛ بنابراین داده‌ها سری زمانی می‌باشند.

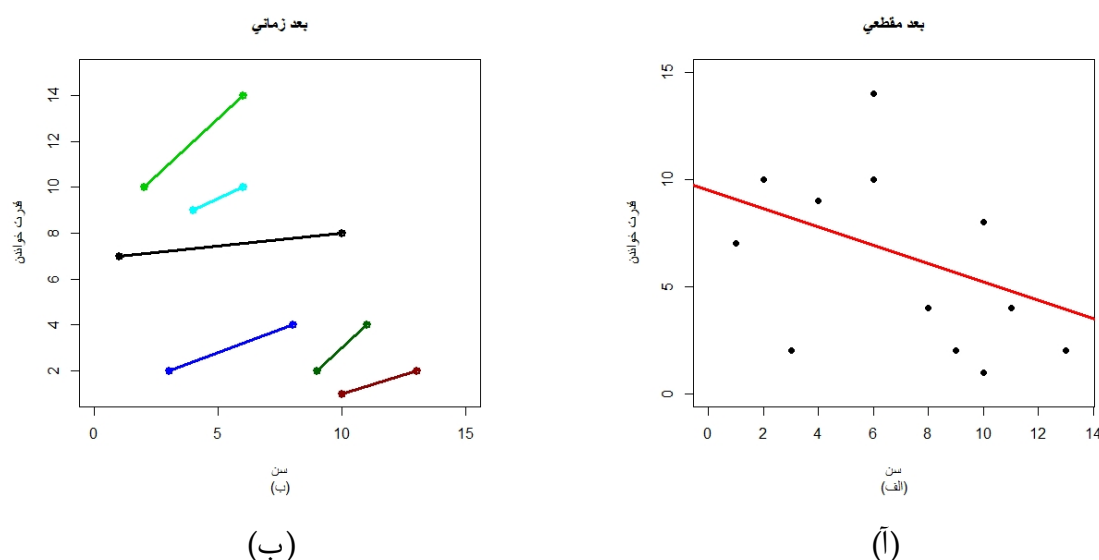
۲.۱.۴ اهداف و ضرورت تحلیل داده‌های طولی

با جمع‌آوری مشاهدات واحدها در طول زمان، علاوه بر دستیابی به رابطه بین متغیرها، می‌توان تغییرات واحدها در طول زمان را نیز مورد مطالعه قرار داد. در نتیجه در تحلیل داده‌های طولی با سه هدف زیر روبرو هستیم:

۱. مطالعه تغییرات واحدها در طول زمان.
۲. مطالعه رابطه بین متغیرها (اثر متغیرهای پیشگو بر روی متغیر پاسخ).
۳. مطالعه تاثیر زمان بر رابطه بین متغیرها.

برای درک بهتر به مثال زیر از دیگل و همکاران (۲۰۰۲) اشاره می‌شود. در شکل ۱.۴ داده‌ها از یک مطالعه طولی مربوط به قدرت خواندن کودکان به‌دست آمده‌اند که در آن شاخص توانایی مطالعه هر کودک، دو بار اندازه‌گیری شده است. بعد مقطعی داده‌ها در قسمت (آ)، کاهش قدرت خواندن کودکان با گذشت زمان را نشان می‌دهد؛ در حالی که در قسمت (ب)، بیان‌کننده افزایش توانایی هر کودک با گذشت زمان است. در مجموع می‌توان گفت در حالی که توانایی خواندن در کودکان جوان‌تر سطح بالاتری دارد؛ همه کودکان 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 برخوردارند.

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. بر خلاف پذیره‌های روش‌های آماری متداول مبتنی بر پایه استقلال مشاهدات مربوط به هر واحد دارای همبستگی درونی بوده و از مشاهدات سایر واحدها مستقل می‌باشند. واضح است که جهت رسیدن به استنباط علمی، در نظر گرفتن این همبستگی امری ضروری



شکل ۱.۴: نمودار قدرت خواندن کودکان

است.

دلایل بالا از جمله مهم‌ترین مواردی است که روش‌های آماری کلاسیک از قبیل رگرسیون و سری‌های زمانی به تنهایی قادر به تحلیل این داده‌ها نیستند.

۲.۴ روش معادلات برآوردیابی تعمیم‌یافته

روش معادلات برآوردیابی تعمیم‌یافته توسط لیانگ و زیگر (۱۹۸۶) به‌عنوان یک روش برآوردیابی مناسب در GLM با مشاهدات وابسته، معرفی شد. معادله برآوردیابی تعمیم‌یافته، تعمیمی از روش‌های برآوردیابی متداول در مدل‌های رگرسیونی کلاسیک بر اساس شبه‌درست‌نمایی است. این روش در تحلیل داده‌های طولی بسیار کارا تر از روش‌های متداول در مدل‌های رگرسیونی کلاسیک است.

در این دیدگاه همبستگی بین مشاهدات با در نظر گرفتن ماتریس‌های همبستگی مبنای^۳ مختلف مدل‌سازی می‌شود که تعیین درستی از این ساختارهای همبستگی در بهبود کارایی برآوردیابی ضرایب رگرسیونی موثر خواهد بود. برآورد ضرایب رگرسیونی و مولفه‌های واریانس در این روش با توجه به فرض ساختارهای متفاوت از ماتریس همبستگی مبنای، توسط امید ریاضی و واریانس توزیع فرضی جامعه به‌دست می‌آیند و لذا نیازی به تعیین توزیع احتمال حاشیه‌ای کامل مشاهدات یا تابع درست‌نمایی نیست.

³ Working Correlation Matrix

۱.۲.۴ خانواده توزیع‌های نمایی

خانواده‌ای از توزیع‌ها را عضو خانواده توزیع‌های نمایی گوییم، اگر بتوان تابع چگالی احتمال آن را به صورت زیر نوشت:

$$f_Y(y) = \exp \left\{ \frac{1}{\phi} [\theta y + c(\theta)] + d(\phi, y) \right\}$$

که در آن θ پارامتر طبیعی است که تابعی از امید ریاضی متغیر پاسخ است و ϕ پارامتر پراکندگی یا مقیاس است که برای تمامی متغیرهای پاسخ یکسان است و $c(\theta)$ و $d(\phi, y)$ توابع معلوم متناظر با نوع چگالی مربوط به خانواده توزیع‌های نمایی می‌باشند. برای بسیاری از اعضای این خانواده ϕ برابر با σ^2 یا ۱ است.

اگر متغیر تصادفی Y داری توزیعی از خانواده توزیع‌های نمایی باشد، امید ریاضی و واریانس آن به صورت زیر محاسبه می‌شوند:

$$E(Y) = \mu = -c'(\theta) \quad (1.4)$$

$$V(Y) = \sigma^2 = -c''(\theta) \phi. \quad (2.4)$$

از جمله توزیع‌های نمایی معروف می‌توان به توزیع نرمال، پواسون، دوجمله‌ای، نمایی، گاما، بتا و گوسی معکوس اشاره کرد. به عنوان مثال، توزیع پواسون را در نظر بگیرید:

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

که می‌توان آن را به شکل خانواده توزیع‌های نمایی نوشت:

$$\exp \{ [y \log \mu - \mu] - \log y! \}$$

که در آن با توجه به روابط (۱.۴) و (۲.۴) داریم:

$$\phi = 1, \quad \theta = \log \mu, \quad c(\theta) = -e^\theta, \quad d(\phi, y) = -\log y!$$

$$E(Y) = \mu = -c'(\theta) = e^\theta, \quad V(Y) = \mu = -c''(\theta) \phi = e^\theta$$

۲.۲.۴ تابع شبه درست‌نمایی

هدف از برازش مدل‌های مختلف رگرسیون بررسی تأثیر متغیرهای تبیینی بر متغیر پاسخ است که در آن نتایج تحلیل داده‌ها انجام شود. استنباط در مورد پارامترهای مدل به محقق کمک می‌کند تا تشخیص دهد کدام متغیر تبیینی روی متغیر پاسخ اثر گذار می‌باشد، در حالی که سایر متغیرهای تبیینی که می‌توانستند اثر تداخلی داشته باشند، ثابت در نظر گرفته شده‌اند. در حقیقت اندازه پارامترهای برآورد شده مدل، توانایی و اهمیت اثر را مشخص می‌سازند. دو دیدگاه متمایز برای استنباط پارامترها وجود دارند.

۱- کمترین توان‌های دوم کلاسیک: این دیدگاه برآورد بهینه از نقطه نظر توابع برآورد است که توسط گودامب (۱۹۶۰) معرفی شده است.

۲- روشی مشابه با ML: دیدگاهی که توسط ودربرن (۱۹۷۴) معرفی شده و به‌عنوان پایه‌ای برای تحلیل GLMs است که روش شبه‌درست‌نمایی نامیده شد (هید ۱۹۹۷؛ هاردین و هیلب، ۲۰۰۲).

در برخی از تحلیل‌ها از جمله زمانی که مقادیر پاسخ وابسته‌اند و به خانواده توزیع‌های نمایی تعلق دارند، دسترسی به شکل صحیح و کاملی از تابع درست‌نمایی غیرممکن است. به این دلیل دیدگاه ودربرن (۱۹۷۴) می‌تواند به‌عنوان حالتی ویژه از دیدگاه تابع برآورد بهینه، در جایی که فضای تابع برآورد به رده‌ای خاص محدود شده، در نظر گرفته شود. وی پذیره‌های روی تابع احتمال را با تعریف تابعی تنها براساس رابطه بین میانگین و واریانس جامعه، جایگزین کرد که معادلات و ویژگی‌هایی مشابه با تابع لگاریتم درست‌نمایی دارند. برای معرفی تابع شبه‌درست‌نمایی در یک مدل رگرسیونی تابع امتیاز $U = \frac{\partial \ell}{\partial \mu}$ را به‌وسیله معادله مشتق در نظر بگیرید که با استفاده از قاعده زنجیره‌ای و روابط (۱.۴) و (۲.۴) به شکل زیر ساده می‌شود:

$$\frac{\partial \ell}{\partial \mu} = \frac{\partial \ell}{\partial \theta} \frac{\partial \theta}{\partial \mu} = \frac{y + c'(\theta) - 1}{\phi} \frac{-1}{c''(\theta)} = \frac{y - \mu}{V(y)}$$

یادآوری می‌کنیم که روابط $E \left[\frac{\partial \ell}{\partial \theta_i} \right] = 0$ و $E \left[\frac{\partial^2 \ell}{\partial \theta_i^2} + \left(\frac{\partial \ell}{\partial \theta_i} \right)^2 \right] = 0$ برای این تابع برقرار است. از این رو معادله $U = 0$ مانند مشتق یک تابع امتیاز رفتار می‌کند و این یک تابع برآورد شبه‌امتیاز^۴ نامیده می‌شود. از این رو شبه‌درست‌نمایی به کمک رابطه بین گشتاورهای مرتبه اول و دوم توزیع فرضی جامعه به برآورد پارامترها می‌پردازد که به‌صورت زیر تعریف شده است:

$$Q(y; \mu) = \int_y^\mu \frac{y - \mu}{V(\mu)} d\mu$$

هم‌چنین واریانس پاسخ به‌عنوان تابعی از میانگین به‌صورت $V(Y_i) = \phi V(\mu_i)$ تعریف می‌شود و ϕ پارامتر پراکندگی توزیع مورد نظر است. وی هم‌چنین اشاره می‌کند که این معادلات برای هر انتخابی از تابع توزیع، حتی زمانی که انتگرال معادلات متناظر با عضوی از خانواده توزیع‌های نمایی نباشند یا حتی زمانی که توزیع پارامتر معلوم نباشد، حل خواهد شد. ضوابط تابع واریانس و شبه‌درست‌نمایی برای چند توزیع از خانواده توزیع‌های نمایی در جدول ۲.۴ فهرست شده‌اند.

⁴ Quasi Score

جدول ۲.۴: ضوابط تابع واریانس و شبه‌درست‌نمایی برای چند توزیع از خانواده توزیع‌های نمایی

توزیع	$V(\mu)$	$Q(\mu)$
نرمال	۱	$-\frac{1}{\mu} \sum (y - \mu)^2$
برنولی	$\mu(1 - \mu)$	$y \ln(\mu/(1 - \mu)) + \ln(1 - \mu)$
پواسون	μ	$y \ln(\mu) - \mu$
گاما	μ^2	$-(y/\mu + \ln(\mu))$
گوسی معکوس	$\mu + \mu^2$	$y(\ln(\mu) - 2 \ln(\mu + 1))$

۳.۴ مدل حاشیه‌ای سیمپلکس برای تحلیل داده‌های نسبت طولی

حالت خاصی از روش GEEs که توسط لیانگ و زیگر (۱۹۸۶) و پرنیس (۱۹۸۸) برای مشاهدات وابسته ارائه شده است؛ روش شبه‌درست‌نمایی و دربرن (۱۹۷۴) است. به عبارت دیگر، در این روش وابستگی بین مشاهدات نیز در تحلیل داده‌ها، توسط یک ساختار ماتریسی به نام ساختار همبستگی مبنا منظور می‌شود. این روش برآورد پارامترها را توسط گشتاورهای مرتبه اول و دوم توزیع فرضی جامعه و نه بر اساس توزیع احتمال حاشیه‌ای کامل مشاهدات یا روش درست‌نمایی کامل انجام می‌دهد که از ویژگی‌های آن در کاربرد برای GLM است و به‌طور مستقیم از گشتاور مرتبه اول جامعه برای برآورد ضرایب رگرسیونی و از گشتاور مرتبه دوم برای برآورد مولفه‌های واریانس و کوواریانس و در واقع برای برآورد پارامترهای همبستگی بین واحدها در یک گروه، بهره می‌گیرد و به این دلیل در برآوردیابی ضرایب رگرسیونی، کارایی بیشتری دارد.

۱.۳.۴ مدل حاشیه‌ای

فرض کنید y_{ij} مقادیر متغیر پاسخ برای $(j = 1, \dots, n_i, i = 1, \dots, m)$ و نشان‌دهنده j امین اندازه روی i امین واحد باشد و همچنین $t_{ij}, j = 1, \dots, n_i$ زمان‌های متوالی متناظر با هر

واحد است که در آن زمان‌ها اندازه‌گیری انجام شده است. علاوه بر این $k = 1, \dots, p$ x_{ijk} متشکل از متغیرهای تبیینی است که با هر مقدار y_{ij} (متغیر پاسخ) متناظر است. حال فرض کنید y_{ij} که مصداق متغیرهای تصادفی Y_{ij} هستند از توزیع سیمپلکس ($Y_{ij} \sim S^-(\mu_{ij}, \sigma_{ij}^2)$) که تابع چگالی آن در رابطه (۱۱.۲) آمده است، پیروی کنند. در توزیع مذکور $\mu_{ij} \in (0, 1)$ پارامتر میانگین و $\sigma_{ij}^2 > 0$ پارامتر پراکندگی هستند و هر دو پارامتر می‌توانند تابعی از متغیرها باشند. قرار می‌دهیم $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{im})^T$ ، $\mathbf{x}_{ij} = (x_{ij1}, \dots, x_{ijp})^T$ که در آن Y_1, \dots, Y_m از هم مستقل‌اند.

مشخصات کلی مدل حاشیه‌ای سیمپلکس به‌صورت زیر است:

۱- مولفه اول میانگین حاشیه‌ای یا احتمال پاسخ برای نمونه i ام به‌صورت زیر که مدلی برای توصیف اثرات متوسط جامعه است، تعریف می‌شود

$$g(\mu_{ij}) = \mathbf{x}_{ij}^T \boldsymbol{\beta}$$

که در آن $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})^T$ بردار p بعدی از ضرایب رگرسیونی و $g : (0, 1) \rightarrow (-\infty, \infty)$ تابع پیوند لجیت است.

۲- مولفه دوم، مدلی برای توصیف پارامتر پراکندگی است؛ با در نظر گرفتن $\sigma_{ij}^2 > 0$ به‌عنوان پارامتر پراکندگی، پیشگوی خطی که تابعی از متغیرهای z_{ij} است (که ممکن است z_{ij} زیرمجموعه‌ای از x_{ij} ها باشند) به‌صورت زیر تعریف می‌شود:

$$h(\sigma_{ij}^2) = \mathbf{z}_{ij}^T \boldsymbol{\gamma}$$

که در آن $\boldsymbol{\gamma} = (\gamma_0, \dots, \gamma_{r-1})^T$ که γ_0 عرض از مبدا است و $h : (0, \infty) \rightarrow (-\infty, \infty)$ تابع پیوند لگاریتمی است. همچنین σ_{ij}^2 را می‌توان به‌صورت زیر نوشت:

$$\sigma_{ij}^2 = \exp(\mathbf{z}_{ij}^T \boldsymbol{\gamma}) = \prod_{k=0}^{r-1} (e^{\gamma_k})^{z_{ijk}} = e^{\gamma_0} \prod_{k=1}^{r-1} (e^{\gamma_k})^{z_{ijk}}.$$

جایی که z_{ij} داده‌های رسته‌ای^۵ هستند و ساختار جمعی^۶ می‌تواند مورد توجه قرار گیرد. برای مثال، در مثال داده‌های جراحی چشم در فصل ۵، مدل‌سازی پراکندگی متغیر در ۳ تیمار متفاوت (تأثیر سطح غلظت گاز در داده‌ها) به‌صورت زیر است

$$\sigma_{ij}^2 = \sigma_0 + \sigma_1 z_{ij1} + \sigma_2 z_{ij2}$$

که در آن z_{ijk} ، متغیرهای ظاهری مربوط به تیمار ۱ و ۲ هستند؛ به‌طوری‌که σ_0 نشان‌دهنده پراکندگی تیمار ۳ است.

⁵ Categorical Covariates

⁶ Additive Structure

۳- مولفه سوم برای مدل‌سازی ساختار همبستگی است؛ همبستگی بین Y_{ik} و Y_{ij} که تابعی از میانگین‌های حاشیه‌ای و احتمالاً پارامترهای فرعی $\alpha = (\alpha_1, \dots, \alpha_q)^T$ است، به صورت زیر تعریف می‌شود

$$\text{Corr}(Y_{ij}, Y_{ik}) = \rho(\mu_{ij}, \mu_{ik}, \alpha)$$

که در آن ρ تابعی معلوم و α برداری از پارامترهای نامعلوم است. انواع مختلفی از ساختارهای همبستگی ممکن است برای تابع ρ مورد استفاده قرار گیرند که سه مورد رایج در تحلیل داده‌های طولی، تبادلی پذیر^۷، اتورگرسیو مرتبه اول (AR(1))^۸ و ساختار همبستگی وابسته مرتبه m ^۹ هستند.

۲.۳.۴ روش GEEs برای برآورد پارامترها

با توجه به بردار میانگین برای هر i ، $\mu_i = (\mu_{i1}, \dots, \mu_{in_i})^T$ ، بردار امتیاز به صورت زیر تعریف می‌شود

$$u_i = (u_{i1}, \dots, u_{in_i})^T, \quad u_{ij} = -\frac{1}{\psi} d'(y_{ij}; \mu_{ij})$$

که تحت شرایط نظم، $E(u_{ij}) = 0$ که می‌توان نتیجه گرفت $E(u_i) = 0$ (سانگ و تان، ۲۰۰۰)؛ حال واریانس u_{ij} طبق رابطه (۱۱.۱) برابر است با

$$\text{var}(u_{ij}) = \frac{\sigma_{ij}^2}{\psi} E\{d''(Y_{ij}; \mu_{ij})\} = \frac{\psi \sigma_{ij}^4}{\mu_{ij}(1 - \mu_{ij})} + \frac{\sigma_{ij}^2}{v(\mu_{ij})}$$

ماتریس کوواریانس بردار پاسخ i ام برای مقادیر $i = 1, \dots, n$ ، در روش GEEs به صورت زیر تعریف می‌شود

$$V_i = \text{diag}\{\text{var}(s_{ij})\} R(\alpha) \text{diag}\{\text{var}(s_{ij})\}$$

که در آن $\text{diag}\{\text{var}(s_{ij})\}$ برای مقادیر $i = 1, \dots, n_i$ یک ماتریس قطری $(n_i \times n_i)$ است. همچنین، کمیت $R(\alpha)$ یک ماتریس $(n_i \times n_i)$ است که ماتریس همبستگی مبنا نامیده می‌شود و به طور کامل با بردار پارامتر فرعی α که برداری $(q \times 1)$ است، مشخص می‌شود که این بردار را پارامتر همبستگی می‌نامند و می‌بایست در روند برازش مدل برآورد شود. بردار مبنا را به صورت $s_i = \text{diag}\{v(\mu_{ij})\} u_i$ تعریف می‌کنیم.

روش GEEs برای برآورد بردار ضرایب رگرسیون β ، شامل یک سیستم غیرخطی از معادلات بر حسب β است که عموماً به طور تکراری با الگوریتم نیوتن-رافسون حل می‌شود؛ این معادلات به صورت زیر هستند

$$\Psi_1(\beta, \gamma, \alpha) = \sum_{i=1}^m D_i^T A_i V_i^{-1} s_i = 0 \quad (3.4)$$

⁷ Exchangeable Correlation

⁸ Autoregressive Correlation (AR)

⁹ M-Dependent Correlation

که در آن $A_i = \text{diag}\{\sigma_{ij}^{-2} v(\mu_{ij}) \text{var}(u_{ij})\}$ و $D_i^T = \frac{\partial \mu_i^T}{\partial \beta}$ ماتریس $(p \times n_i)$ از مشتقات جزئی میانگین i امین واحد نسبت به ضرایب رگرسیون و V_i ماتریس کوواریانس بردار پاسخ Y_i است. در این روش تنها به پذیره‌هایی از گشتاورهای مرتبه اول نیاز است که در صورت برقراری آن‌ها، برآوردهای ضرایب رگرسیون سازگار خواهند بود. مزیت روش GEEs در آن است که تنها لازم است پارامترهای میانگین به‌طور صحیح مشخص شده باشند که در این صورت برآورد ضرایب رگرسیونی سازگارند و این موضوع حتی اگر ساختار ماتریس همبستگی مبنا به‌صورت صحیح تعیین نشده باشد، نیز برقرار است (پارک و شین، ۱۹۹۹).

چنان‌چه گفته شد، از ویژگی‌های روش GEEs مدل‌سازی همبستگی بین مشاهدات در یک واحد با ماتریس همبستگی مبنا است و انتخاب صحیح ساختار همبستگی در بهبود کارآمدی برآوردیابی ضرایب رگرسیونی موثر است. به این منظور لازم است ماتریس همبستگی صحیحی از بردار پاسخ i ام داشته باشیم تا بتوانیم ماتریس کوواریانس درستی از بردار Y_i محاسبه کنیم. در واقع باید کمیت $R_i(\alpha)$ صحیح انتخاب شود تا V_i مناسبی برای تحلیل به‌دست آید (هورتن و لیپسیتز، ۱۹۹۹ و هاردین و هیلب، ۲۰۰۳).

ماتریس همبستگی مبنا با بردار پارامتر فرعی α تعیین می‌شود که عموماً نامعلوم است و لازم است در روند برازش تکراری به‌همراه ضرایب رگرسیونی β برآورد شود.

در برآورد بردار پارامتر همبستگی α برای مقادیر $i = 1, \dots, m$ و $j = 1, \dots, n_i$ یک معادله برآوردیابی به‌صورت زیر تعریف می‌شود

$$\Psi_3(\beta, \gamma, \alpha) = \sum_{i=1}^m \left(\frac{\partial \xi_i^T}{\partial \alpha} \right) H_i^{-1} (r_i - \xi_i) = [0]_{q \times 1} \quad (4.4)$$

که در آن $\xi_i = E(r_i)_{q \times 1}$ و H_i ، ماتریس کواریانس که ماتریسی است قطری $(q \times q)$ از مقادیر واریانس و $r_i = (r_{i1}, r_{i2}, \dots, r_{in_i-1}, r_{in_i})$ است که در آن r_{ij} ، i, j امین مانده پیرسون و $q = (n_i)$ است. این معادله برآوردیابی نتیجه می‌دهد که پارامتری کردن ماتریس همبستگی مبنا، از طریق معادله $\xi_i = E(r_i)$ است و برای برآورد پارامترهای همبستگی α بایستی مقادیر \hat{r}_{ij} به‌دست آمده از برآورد متداول $\hat{\beta}$ ، جانشین مانده‌های پیرسون شوند.

دیدگاه‌های یکسانی در ارتباط با ساختارهای ماتریس متفاوت همبستگی مبنا وجود دارند. ممکن است مشاهدات درون یک واحد از مرتبه مشخصی پیروی نکنند یا به‌طور یکسان همبسته باشند یا فرض‌های پیچیده‌تری بر روی مشاهدات درون یک واحد در نظر گرفته شوند. در نتیجه در هر حالتی ساختاری متفاوت مد نظر محقق است و شاید یک بردار از پارامترهای اضافی تا حدود $(\max n_i) - 1$ پارامتر یا ماتریس کاملی از پارامترهای با $(\max n_i) - n_i$ پارامتر اضافی لازم باشد.

بر این اساس محققین انتخاب‌هایی را از ساختارهای ماتریس همبستگی مبنا با توجه به موضوع مورد بررسی و مدل‌های مد نظر فراهم کرده‌اند که در قسمت‌های بعدی تعدادی از

آن‌ها را همراه با فرمول‌هایی ساده برای برآورد مولفه‌های همبستگی α مربوط به آن‌ها، معرفی می‌کنیم. در تمام حالات فرمول‌ها به‌طور مستقیم از معادله برآوردیابی ارائه‌شده در عبارت (۴.۴) نتیجه می‌شود.

ساختار همبستگی اتورگرسیو مرتبه اول

این ساختار يك همبستگی زمانی را برای مشاهدات مختلف درون واحدها در نظر می‌گیرد. در واقع اگر همبستگی مشاهدات تکراری در درون واحدها، مرتبه‌ای طبیعی داشته باشند این ساختار به خوبی می‌تواند آن‌را نشان دهد. برای مثال، در مطالعات پزشکی، واحدها می‌توانند بیمارانی با اندازه‌های تکراری باشند که در ماه‌های متوالی مقدار ویژگی مشخصی از آن‌ها اندازه‌گیری شده‌است. در حالت کلی ساختار همبستگی اتورگرسیو مرتبه اول برای مقادیر $i = 1, \dots, n$ و زمان‌های $t = 1, \dots, n_i$ به‌صورت زیر تعریف می‌شود

$$\text{Corr}(Y_{it}, Y_{it'}) = \begin{cases} \alpha^{|t-t'|} & t \neq t' \\ 1 & t = t' \end{cases} \quad (5.4)$$

در این تعریف مقدار برآورد پارامتر همبستگی α توسط معادله برآوردیابی (۴.۴)، مانده‌های پیرسون \hat{r}_{it} و رابطه (۱۰.۳) به‌صورت زیر است

$$\hat{\alpha} = \frac{1}{\hat{\phi}} \sum_{i=1}^n \left[\frac{\sum_{t=1}^{n_t-1} \hat{r}_{it} \hat{r}_{i,t+1}}{n_i}, \dots, \frac{\sum_{t=1}^{n_t-l} \hat{r}_{it} \hat{r}_{i,t+l}}{n_i} \right]$$

نکته قابل توجه آن است که يك فرآیند اتورگرسیو مرتبه l ، خود همبستگی غیر صفر، برای بیش از l تاخیر دارد و به‌عنوان نمونه برای زمانی که چهار متغیر تبیینی مورد توجه محقق باشد، ساختار (۵.۴) به‌شکل ماتریسی زیر باز نویسی می‌شود

$$\text{AR}(1) = \begin{bmatrix} 1 & \alpha & \alpha^2 & \alpha^3 \\ \alpha & 1 & \alpha & \alpha^2 \\ \alpha^2 & \alpha & 1 & \alpha \\ \alpha^3 & \alpha^2 & \alpha & 1 \end{bmatrix}$$

ساختار همبستگی تبادل‌پذیر

ساده‌ترین شکل ساختار ماتریس همبستگی يك ماتریس همانی است. در تعمیم ساده‌ای از این ساختار همبستگی مستقل می‌توان فرض کرد که برخی مشاهدات در درون واحدها همبستگی مشترك داشته باشند و در حقیقت يك پارامتر فرعی به ماتریس اضافه می‌شود و مقدار (j, l) امین عنصر از این ماتریس همبستگی مبنا برای مقادیر $i = 1, \dots, n$ و $j = 1, \dots, n_i$ در حالت کلی

به شکل زیر تعریف می‌شود

$$\text{Corr}(Y_{ij}, Y_{il}) = \begin{cases} 1 & j = l \\ \alpha & j \neq l \end{cases} \Leftrightarrow \mathbf{R}(\alpha) = \begin{bmatrix} 1 & \alpha & \cdots & \alpha \\ \alpha & 1 & \cdots & \alpha \\ \vdots & \vdots & \ddots & \vdots \\ \alpha & \alpha & \cdots & 1 \end{bmatrix}$$

این ساختار همبستگی برای مجموعه داده‌هایی که در آن‌ها اندازه‌های تکراری وابستگی زمانی ندارند و همچنین برای هر جایگشتی از آن‌ها، معتبر است. در این حالت پارامتر فرعی α يك اسکالر است که به کمک معادله برآوردیابی (۴.۴) و توسط معادله زیر برآورد می‌شود

$$\hat{\alpha} = \frac{1}{\hat{\phi}} \sum_{i=1}^n \left[\frac{\sum_{j=1}^{n_i} \sum_{l=1}^{n_i} \hat{r}_{ij} \hat{r}_{il} - \sum_{j=1}^{n_i} \hat{r}_{ij}^2}{n_i(n_i - 1)} \right]$$

که در آن از برآورد مانده‌های پیرسون \hat{r}_{ij} و رابطه (۱۰.۳) برای برآورد پارامتر پراکندگی ϕ استفاده می‌شود. این ساختار تحت چندین نام دیگر از جمله ساختار همبستگی مساوی^{۱۰} و همبستگی مشترك^{۱۱} نیز در مباحث مختلف آماری شناخته می‌شود.

ساختار همبستگی وابسته مرتبه m

در مباحث سری زمانی، ممکن است همبستگی بین مشاهدات برای تعداد معینی از واحدهای زمانی و نه برای تمامی آن‌ها مد نظر محقق باشد. در این حالت ساختار ماتریس همبستگی برای مقادیر $i = 1, \dots, n$ و $j = 1, \dots, n_i$ به شکل زیر تعریف می‌شود

$$\text{Corr}(Y_{ij}, Y_{i,j+t}) = \begin{cases} 1 & t = 0 \\ \alpha_t & t = 1, 2, \dots, m \\ 0 & t > m \end{cases}$$

مشابه قبل برای مقادیر $t = 1, \dots, m$ می‌توان پارامترهای همبستگی $\hat{\alpha}_t$ را توسط مانده‌های پیرسون برآورد شده \hat{r}_{ij} ، معادله برآوردیابی (۴.۴) و رابطه (۱۰.۳) به صورت زیر برآورد نمود

$$\hat{\alpha}_t = \frac{1}{(N_t - p)\hat{\phi}} \sum_{i=1}^n \sum_{j \leq n_i - t} \hat{r}_{ij} \hat{r}_{i,j+t}$$

که در آن p بعد بردار ضرایب رگرسیونی و $N_t = \sum_{i=1}^n (n_i - t)$ است. برای مثال، اگر در مدل‌سازی ماتریسی با بعد ۳ مورد استفاده قرار گیرد، آن‌گاه آن ماتریس

¹⁰ Equal Correlation

¹¹ Common Correlation

ساختار زیر را دارد

$$R(\alpha) = \begin{bmatrix} 1 & \alpha_1 & \alpha_2 & \alpha_3 \\ \alpha_1 & 1 & \alpha_1 & \alpha_2 \\ \alpha_2 & \alpha_1 & 1 & \alpha_1 \\ \alpha_3 & \alpha_2 & \alpha_1 & 1 \end{bmatrix}$$

۳.۳.۴ روش معادلات برآوردیابی تعمیم‌یافته مرتبه اول (GEE1) و مرتبه دوم (GEE2)

در مدل‌سازی به روش GEEs اگر مقدار ماتریس همبستگی مبنا انتخاب شده در تحلیل داده‌ها، بسیار متفاوت از ماتریس همبستگی حقیقی باشد ممکن است این روش در برآورد پارامترهای مدل ناکارآمد باشد. به این منظور محققان مدل‌سازی ساختاری دیگر را توصیه کرده‌اند که در آن یک ساختار همبستگی توسط گشتاورهای مرتبه بالاتر نیز در تحلیل اعمال می‌شود. این تغییرپذیری شامل در نظر گرفتن یک معادله برآوردیابی اضافی برای پارامتر همبستگی α است که آن را معادله برآوردیابی تعمیم‌یافته مرتبه دوم (GEE2) می‌نامند و شکل ساده‌تر آن، معادله برآوردیابی تعمیم‌یافته مرتبه اول (GEE1) است. در واقع این دو روش تعمیم‌هایی از روش GEEs هستند (بالمی و لی، ۲۰۰۹؛ زورن، ۲۰۰۱).

روش GEE1 تقریباً به همان کارایی روش GEE2 برای برآورد بردار ضرایب رگرسیونی β است و تنها زمانی که بهبود کارایی قابل توجه‌ای مد نظر محقق باشد، استفاده از روش GEE2 توصیه می‌شود. در حالت کلی مجموعه معادلات ساده GEE1 ترجیح داده می‌شوند و در تحلیل داده‌های وابسته کاربرد بیشتری دارند.

با توجه به پرنتمی و ژائو (۱۹۹۱)، در روش GEE2، با اضافه کردن یک مجموعه از معادلات برآورد شده برای پارامترهای همبستگی براساس باقی‌مانده‌های امتیاز استاندارد شده^{۱۲}، به صورت زیر تعریف می‌شود:

$$r_{ij} = \frac{u_{ij}}{\sqrt{\text{var}(u_{ij})}} = \frac{u_{ij}}{\sqrt{\frac{\sigma_{ij}^2}{\gamma} E\{d''(Y_{ij}; \mu_{ij})\}}} \quad (6.4)$$

به راحتی می‌توان دید که خواص گشتاوری $E(r_{ij}) = 0$ ، $\text{var}(r_{ij}) = 1$ و $\text{Corr}(u_{ij}, u_{ij'}) = \text{Corr}(r_{ij}, r_{ij'})$ برای باقی‌مانده‌های امتیاز برقرار است.

حال برای برآورد ضرایب رگرسیون γ ، معادله برآوردیابی به صورت زیر تعریف می‌شود

$$\Psi_{\gamma}(\beta, \gamma, \alpha) = \sum_{i=1}^m \left(\frac{\partial \sigma_i^T}{\partial \gamma} \right) \Sigma_i^{-1} (d_i - \sigma_i) = 0 \quad (7.4)$$

¹² Standardised Score Residuals

که در آن $\mathbf{d}_i = (d(y_{i1}; \mu_{i1}), \dots, d(y_{in_i}; \mu_{in_i}))^T$ و ماتریس کواریانس و $\boldsymbol{\sigma}_i = E(\mathbf{d}_i) = (\sigma_{i1}^2, \dots, \sigma_{in_i}^2)^T$.

با تعمیم روش GEEs توسط بسط معادلات برآوردیابی برای پارامترهای $\boldsymbol{\theta} = (\beta, \gamma, \alpha)$ معادلات GEE2 به صورت زیر حاصل می‌شوند

$$\Upsilon(\boldsymbol{\theta}) = \Upsilon(\beta, \gamma, \alpha) = \begin{bmatrix} \Psi_1(\beta, \gamma, \alpha) \\ \Psi_2(\beta, \gamma, \alpha) \\ \Psi_3(\beta, \gamma, \alpha) \end{bmatrix} = \mathbf{0}. \quad (۸.۴)$$

معادله برآوردیابی $\Upsilon(\boldsymbol{\theta}) = \mathbf{0}$ نااریب است، زیرا $E(\Upsilon(\boldsymbol{\theta})) = \mathbf{0}$. از این رو طبق نظریه استاندارد^{۱۳} معادلات برآوردیابی تحت شرایط نظم، برآوردگر $\hat{\boldsymbol{\theta}} = (\hat{\beta}, \hat{\gamma}, \hat{\alpha})$ ، برآوردگری سازگار است و $mJ^{-1}(\boldsymbol{\theta})$ به طور مجانبی گاوسی چند متغیره با میانگین صفر و ماتریس کواریانس $mJ^{-1}(\boldsymbol{\theta})$ است که $J(\boldsymbol{\theta})$ ماتریس اطلاع گودامب^{۱۴} به صورت $J(\boldsymbol{\theta}) = S^T R^{-1} S$ است؛ که در آن ماتریس حساسیت $S = E\{\partial \Upsilon(\boldsymbol{\theta}) / \partial \boldsymbol{\theta}^T\}$ و ماتریس تغییرپذیری برابر با $R = E\{\Upsilon(\boldsymbol{\theta}) \Upsilon^T(\boldsymbol{\theta})\}$ است. برای مشاهده جزئیات بیشتر به سانگ و همکاران (۲۰۰۴) رجوع شود (ضمیمه را ۲.آ ببینید). با استفاده از الگوریتم نیوتن-رافسون، حل $\hat{\boldsymbol{\theta}}$ برای معادله (۸.۴) را می‌توان به صورت زیر به دست آورد

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - S^{-1} \Upsilon(\boldsymbol{\theta}^{(k)}).$$

۴.۳.۴ انواع معادلات برآوردیابی تعمیم‌یافته

برخی محققان دو رده‌بندی از مدل‌ها را برای ساختار واحدهای درون داده‌ها پیشنهاد کرده‌اند که در حقیقت این تقسیم‌بندی بر مبنای نحوه محاسبه میانگین و به کارگیری اثر تصادفی در مدل است که در این بخش به معرفی آن‌ها به عنوان دو نوع گوناگون از مجموعه روش‌های GEEs می‌پردازیم (هاردین و هیلب، ۲۰۰۳).

مدل متوسط جامعه^{۱۵}

مدلی که مقدار همبستگی درون واحدها را به واسطه اثرات متوسط روی تمام واحدها در نظر می‌گیرد و به عنوان یک مدل حاشیه‌ای مطرح است را مدل متوسط جمعیت می‌نامند. در حقیقت این مدل یک مشخصه‌سازی را برای ساختار کوواریانس سطوح واحدها فرض می‌کند و آن‌گاه مقدار کوواریانس یا همبستگی سطوح را با میانگین گرفتن در میان اطلاعات تمام واحدها برآورد می‌کند. مدل‌های انتقالی^{۱۶} و پاسخ‌های شرطی از نوع مدل‌های متوسط جامعه

¹³ Standard Theory

¹⁴ Godambe

¹⁵ Population Averaged

¹⁶ Transitional

هستند و در تحلیل مطالعات طولی زمانی که بایستی مقدار همبستگی متداول مقادیر پاسخ روی پاسخ‌های قبلی نشان داده شود، استفاده می‌شوند.

مدل واحد ویژه^{۱۷}

این مدل همبستگی داخل واحدها را توسط مطرح کردن مولفه‌های تصادفی سطوح آن‌ها بیان می‌کند و به واسطه معرفی اثر یک واحد، در حالی که به یک کوواریانس سطح اشاره می‌کند، به دست می‌آید. در این مدل هر اثر با استفاده از اطلاعات تنها یک واحد خاص، برآورد می‌شود. مدل‌ها با اثرات ثابت و تصادفی از این نوع مدل هستند.

۵.۳.۴ تحلیل باقی‌مانده‌ها

پیشنهاد می‌شود از دو نوع باقی‌مانده استفاده شود که تشخیص برای پذیره‌های مدل کلیدی، یعنی

(۱) توزیع‌های حاشیه‌ای

(۲) توابع پیوند

(۳) ساختار همبستگی

ایجاد شود. اولی، باقی‌مانده‌های امتیاز استاندارد شده r_{ij} در (۶.۴) است و دیگری باقی‌مانده‌های معمولی استاندارد شده پیرسون $e_{ij} = (y_{ij} - \mu_{ij}) / \sqrt{\text{var}(Y_{ij})}$ است. از آنجایی که $\text{var}(Y_{ij})$ شامل تابع گامای ناکامل است، صورت بسته‌ای ندارد. هم‌آورد نمونه از r_{ij} یا e_{ij} با جایگزین کردن پارامترها با تخمین‌های متناظر به دست می‌آید که آن‌ها را با \hat{r}_{ij} یا \hat{e}_{ij} نشان می‌دهند. تحلیل باقی‌مانده‌ها برای شناسایی سیگنال‌های قوی مربوط به نقض فرض مدل خاص در زیر آورده شده است.

پذیره توزیع سیمپلکس می‌تواند با نمودار \hat{e}_{ij} در مقابل \hat{m}_{ij} که هدف آن بررسی رابطه واریانس- میانگین است، کنترل شود. اگر این پذیره درست باشد، $\text{var}(e_{ij}) = 1$ "مستقل از میانگین μ_{ij} است. بنابراین، نقاط موجود در نمودار باید به‌طور تصادفی در اطراف خط افقی در صفر (امید باقی‌مانده‌ها) با تقریب ۹۵٪ در نوار افقی بین ۲ و -۲ پخش شوند. هر گونه واکنش ظاهری از این امر تخطی از پذیره توزیع یا احتمالاً برازش یک مدل ضعیف را نشان می‌دهد. تحقیقات بیشتری برای شناسایی عامل چنین خروجی مورد نیاز است. از آنجایی که σ^2 بزرگ است این روش قابل اعتماد است زیرا رابطه واریانس میانگین، تحت کنترل $\mu(1 - \mu)$ یک حالت مشابه با توزیع دوجمله‌ای دارد.

با توجه به مک‌کالا و نلدر (۱۹۸۹)، نمودار متغیر وابسته تعدیل یافته a_{ij} در مقابل پیشگوی خطی $\hat{\eta}_{ij}$ برای بررسی تابع پیوند انتخاب شده، استفاده می‌کنیم. حال تعریف می‌کنیم:

$$a_i = g(\mu_i) + \left\{ \frac{3\sigma_{ij}^4}{\mu_{ij}(1-\mu_{ij})} + \frac{\sigma_{ij}^2}{v(\mu_{ij})} \right\}^{-\frac{1}{3}} u(y_i; \mu_i) \quad (9.4)$$

واضح است که $E(a_i) = g(\mu_{ij})$ زیرا $E(u_{ij}) = 0$ و $\text{var}(a_i) = E\{a_i - g(\mu_{ij})\}^2 = 1$ است. اگر تابع پیوند مناسب باشد، نمودار برآورد \hat{a}_i در مقابل $\hat{\eta}_{ij} = \mathbf{x}_{ij}^T \hat{\beta}$ باید یک خط مستقیم را که با تقریب ۹۵٪ نقاط در یک باند با محدودیت‌های بالا و پایین حدود $\hat{\eta}_{ij} \pm 2$ قرار می‌گیرند، نشان دهد. این نمودار مانند مدل‌های خطی تعمیم یافته بهترین تابع پیوند را برای مدل پیشنهاد نمی‌دهد بلکه فقط به بررسی غیر رسمی برای هر گونه نقض شدید پیوند مورد استفاده، می‌پردازد.

اگر چه مدل‌سازی حقیقی ساختار همبستگی داده‌های طولی دشوار است اما ساختارهای همبستگی صحیح به ضرایب رگرسیون، این امکان را می‌دهند که به‌طور موثر تخمین زده شوند. بنابراین، ارزیابی تناسب همبستگی استفاده شده در GEEs از طریق آنالیز باقی مانده‌ها مهم است. توجه داشته باشید که

$$\text{Corr}(r_{ij}, r_{ij'}) = \text{Corr}(a_i, a_{i'})$$

این امر حاکی از آن است که همبستگی واقعی متغیر a_i برابر با باقی مانده امتیاز استاندارد شده r_{ij} است. برای ارزیابی همبستگی داده‌ها، می‌توان برخی از روش‌های اکتشافی ارائه شده در بخش ۳.۴ دیگل و همکاران (۲۰۰۲) را به کار برد.

فصل ۵

کاربرد مدل رگرسیونی سیمپلکس و برازش آن با R

در این فصل، نحوه استخراج استنباط‌های مطرح‌شده در فصل‌های قبلی را برای مدل‌های رگرسیونی سیمپلکس با استفاده از بسته `simplexreg` شرح می‌دهیم. سپس کاربرد مدل را با دو مثال واقعی شامل داده‌های سلول بنیادی خون محیطی و داده‌های جراحی چشم نشان می‌دهیم.

۱.۵ مقدمه

تابع `brzsh` مدل رگرسیونی سیمپلکس، `simplexreg()`، طوری طراحی شده‌است که تا حد امکان مشابه تابع `glm()` استاندارد (تیم مرکزی توسعه R^۱، ۲۰۰۹) که برای برازش GLM به کار می‌رود، باشد. یک تفاوت اساسی بین این دو تابع آن است که در `simplexreg()` می‌توان معادله‌هایی که معرف میانگین و پارامتر پراکندگی هستند را در این تابع لحاظ کرد و در نتیجه دو ماتریس متفاوت از متغیرهای تبیینی، دو پیشگوی خطی و دو مجموعه از ضرایب به دست

^۱R Development Core Team

آورد. روش استفاده از `simplexreg()` در R به صورت زیر است:

```
simplexreg(formula, data, subset, na.action,
link = c("logit", "probit", "cloglog", "neglog"), corr = "Ind", id = NULL,
control = simplexreg.control(...), model = TRUE, y = TRUE, x = FALSE, ...)
```

چهار آرگومان اول، شامل مشخصات مدل استاندارد (چمبرز و هستی، ۱۹۹۲) هستند. از `formula` برای مشخص کردن مدل رگرسیون و پیشگوه‌های خطی، مشابه تابع GLM، استفاده می‌شود. داده‌ها توسط `data` مشخص می‌شوند. آرگومان `link` معرف نوع تابع پیوند و `corr` ساختار همبستگی را تعیین می‌کند. آرگومان `model` معرف یک مقدار منطقی است و تعیین می‌کند که از مقادیر خروجی تابع `simplexreg.fit` مشخصات مدل ظاهر شود. در بسته `simplexreg` و بسته‌های مرتبط با آن می‌توان به توابع زیر که برای استخراج استنباط‌ها در مدل رگرسیونی سیمپلکس مفید هستند، اشاره کرد:

- `summary()` خلاصه‌ای از مدل رگرسیونی برازش داده‌شده را ارائه می‌دهد.
- `coef()` برای استخراج ضرایب برآوردشده مدل، مورد استفاده قرار می‌گیرد.
- `vcov()` ماتریس کواریانس برآوردگرها را ارائه می‌دهد.
- `predict()` برای پیش‌گویی داده‌های جدید می‌توان از آن بهره برد.
- `fitted()` برازش میانگین‌ها برای داده‌های مشاهده‌شده را انجام می‌دهد.
- `residuals()` برای استخراج باقی‌مانده‌ها مورد استفاده قرار می‌گیرد.
- `estfun()` تخمین مشاهدات و برآورد پارامترها را نمایش می‌دهد.
- استخراج مولفه‌های مدل، مولفه‌های ماتریس مدل و مولفه‌های اصلی مدل به ترتیب با توابع `terms()`، `model.matrix()` و `model.frame()` صورت می‌گیرد.
- `loglik()` لگاریتم درست‌نمایی برازش شده را ارائه می‌دهد.
- `plot()` برای رسم نمودارهای مناسب مدل می‌توان از آن بهره برد.

تابع `simplexreg.fit` یک ابزار دیگری برای محاسبات مدل است که به صورت زیر است:

```
simplexreg.fit(y, x, z = NULL, t = NULL, link = "logit", corr = "Ind",
id = NULL, control = simplexreg.control())
```


آرگومان *corr* ساختار همبستگی را در مدل حاشیه‌ای مشخص می‌کند و شامل سه گزینه مستقل، تبادل‌پذیر و اتورگرسیو مرتبه اول است که آن‌ها را به ترتیب با "Ind"، "Exc" و "AR(1)" نشان می‌دهند. به طور پیش فرض، آرگومان *corr* مقدار "Ind" را می‌گیرد که در این صورت مدل حاشیه‌ای معادل مدل خطی تعمیم‌یافته، SGLM، است. در تابع *simplexreg()* اگر به عنوان مثال *formula* به صورت $y \sim x_1 + x_2$ باشد، آن‌گاه متغیرهای تبیینی x_1 و x_2 تنها برای مدل بندی میانگین رگرسیون سیمپلکس در نظر گرفته می‌شوند. در این حالت σ^2 ثابت فرض می‌شود. به عبارت دیگر، $z_i = 1$ و $h(\cdot)$ تابع پیوند همانی است که همان مدل معرفی شده در فصل ۳ بخش ۱.۳ است. اما چنانچه *formula* شامل دو بخش به صورت $y \sim x_1 + x_2 | z_1 + z_2$ باشد، در این حالت σ^2 نیز تابعی از متغیرهای تبیینی z_1 و z_2 است که مدل معرفی شده در فصل ۳ بخش ۲.۳ می‌باشد. برای داده‌های نسبت طولی، مدل‌های سیمپلکس حاشیه‌ای از سه جز، اثرات میانگین جامعه، الگوی پراکنش^۲ و ساختار همبستگی تشکیل می‌شود. حال اگر *formula* به صورت سه قسمتی $y \sim x_1 + x_2 || t$ باشد، معرف مدل با پارامتر پراکندگی ثابت است و اگر به صورت $y \sim x_1 + x_2 | z_1 + z_2 | t$ باشد برای مدل‌های حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر ("AR(1)" یا "Exc" = corr) به کار می‌رود. پارامتر *t* در تابع *simplexreg.fit* و در *formula* مربوط به متغیر کمکی زمان است. تابع پیوند پیش فرض برای میانگین $g(\mu) = \text{logit}$ و برای پارامتر پراکندگی $h(\sigma^2) = \log$ است. تابع پیوند برای میانگین μ را می‌توان با انتخاب *link*، تغییر داد. توابع مجاز برای μ همان توابع پیوند مجاز در *glm()* شامل "logit"، "probit"، "cloglog"، "neglog" و برای σ^2 ، توابع همانی، لگاریتمی و ریشه دوم هستند.

```
plot(x, type = c("residuals", "corr", "GOF"), res = "adjvar", lag = 1, ...)
```

آرگومان *type* نوع نمودار را مشخص می‌کند. اگر *type = "residuals"* تحلیل باقی‌مانده‌ها صورت می‌گیرد که چهار انتخاب برای آن وجود دارد

۱. اگر *res : stdPerr* آن‌گاه باقی‌مانده استاندارد شده پیرسون محاسبه می‌شود که در رابطه (۱۲.۳) آمده است.

۲. اگر *res : appstdPerr* آن‌گاه تقریبی از باقی‌مانده پیرسون محاسبه می‌شود که در رابطه (۱۳.۳) آمده است.

۳. اگر *res : stdscor* آن‌گاه باقی‌مانده استاندارد شده امتیاز محاسبه می‌شود که در رابطه (۶.۴) آمده است.

۴. اگر *res : adjvar* آن‌گاه متغیر وابسته تعدیل یافته، محاسبه می‌شود که در رابطه (۹.۴) آمده است.

^۲Pattern of Dispersion

سه مورد اول به منظور بررسی رابطه میانگین واریانس و نیز تشخیص خطای مدل، در مقابل میانگین μ رسم می‌شوند. در حالی که با استفاده از مورد چهارم، متغیر وابسته تعدیل یافته در مقابل پیشگوی خطی η برای بررسی تابع پیوند رسم می‌شود.

آرگومان `type = "corr"` ساختار همبستگی را نشان می‌دهد. در `lag`، برای بررسی خودهمبستگی از باقی مانده‌های استاندارد شده امتیاز استفاده می‌شود و هرگاه آرگومان `type = "GOF"` باشد انحراف جزئی در مقابل متغیرهای زمان رسم می‌شود.

۲.۵ کاربرد مدل در دو مثال واقعی

این بخش شامل مثال‌های کاربردی است که توسط آن‌ها مباحث نظری بیان شده در فصل‌های سوم و چهارم را به صورت کاربردی تشریح می‌کنیم.

۱.۲.۵ داده‌های سلول بنیادی خون محیطی

پیوندهای سلول بنیادی خون محیطی^۳ (PBSC) به طور گسترده‌ای برای بهبود سریع هماتولوژیک پس از میلوپلیتو تراپی^۴ برای انواع اختلالات خونی بدخیم استفاده می‌شود. بازسازی مجدد هماتوپوئیتیک^۵، به طور گسترده‌ای به تزریق تعداد کافی سلول‌های بنیادی برای پیوند در مغز استخوان محیط میکرو^۶ بستگی دارد. همان‌طور که در مطالعه آلان و همکاران (۲۰۰۲) نشان داده شد، مقدار سلول‌های زنده CD34+ به عنوان یک نشانه مهم برداشت کافی PBSC و همچنین به عنوان یک پیش‌بینی کننده پیوند هماتوپوئیتیک در نظر گرفته می‌شود. در این مطالعه ۲۴۲ بیمار موافقت کردند که پیوند اتولوگ با PBSC را پس از دوزهای شیمی‌درمانی در سال‌های ۲۰۰۳ و ۲۰۰۸ در آزمایشگاه سلول‌های بنیادی هماتوپوئید در مؤسسه سرطان - خدمات درمانی آلبرتا، دریافت و اطلاعات را ثبت کنند. این داده‌ها شامل اطلاعاتی در مورد سن، جنس و ویژگی‌های بالینی بیماران است. داده‌های مربوط به پیوند PBSC، با نام `sdac`، در بسته `simplexreg` گنجانده شده است. این مجموعه داده شامل ۲۳۹ مشاهده (بیمار) و ۵ متغیر تبیینی به شرح زیر است:

³ Peripheral Blood Stem Cell

⁴ Myeloablative Therapy

⁵ Hematopoietic

⁶ Micro-Environment

- age: سن بیماران
- gender: جنسیت بیماران که اکثر آنها یعنی ۱۷۱ نفر، مرد هستند
- rcd: میزان سلول‌های زنده CD34+
- chemo: متغیر ظاهری؛ $\text{chemo} = 0$ معرف پروتکل یک روزه و $\text{chemo} = 1$ معرف پروتکل سه روزه است
- ageadj: متغیر سن تعدیل یافته؛ $\text{age} < 40$ به‌عنوان سن پایه در نظر گرفته و سایر سنین از عدد ۴۰ کم می‌شوند

بخشی از داده‌های sdac موجود در بسته simplexreg به شرح زیر است.

	age	gender	rcd	ageadj	chemo
1	62	M	0.75	22	0
2	39	M	0.83	0	1
3	43	M	0.94	3	1
4	58	M	0.86	18	0
5	43	M	0.54	3	0
6	51	M	0.7	11	1
7	64	F	0.59	24	0
8	64	M	0.82	24	0
9	48	M	0.67	8	1
10	51	M	0.57	11	0
11	57	M	0.75	17	0
12	60	M	0.74	20	1
13	53	M	0.75	13	1
14	47	M	0.94	7	0
15	51	M	0.87	11	1
16	29	M	0.99	0	1
17	43	M	0.83	3	1
18	57	M	0.82	17	1
19	62	F	0.81	22	0

پارامتر پراکندگی ثابت

ابتدا مدل را با فرض ثابت بودن پارامتر پراکندگی σ^2 در نظر گرفتیم. پیشگوی خطی برای μ_i به صورت

$$\text{logit}(\mu_i) = \beta_0 + \beta_1 \text{ageadj} + \beta_2 \text{chemo} \quad (1.5)$$

است. با دستورات زیر می‌توان مدل (۱.۵) را بر روی داده‌ها برازش داد

```
library(simplexreg)
data("sdac", package="simplexreg")
sim.glm1 <- simplexreg(rcd ~ ageadj + chemo, data = sdac)
summary(sim.glm1)
```

نتایج برازش مدل (۱.۵) در جدول ۲.۵ گزارش شده‌اند. نتایج مذکور با ۶ تکرار در مرحله بهینه‌سازی و با انحراف ۲۳۵ به دست آمده که لگاریتم درست‌نمایی آن $۱۰۵/۳$ و p - مقدار مربوطه برابر $۰/۴۹$ است. همان‌طور که ملاحظه می‌شود عرض از مبدا و ضرایب متغیرهای ageadj و chemo معنی‌دار هستند.

جدول ۱.۵: باقی‌مانده استاندارد پیرسون مدل (۱.۵)

کوچکترین مقدار	چارک اول	میانه	چارک سوم	بزرگترین مقدار
-۲/۸۲۵	-۰/۵۸۵	-۰/۰۰۸	۰/۴۹۷	۱/۳۹۶

جدول ۲.۵: برآورد پارامترهای مدل (۱.۵)

پارامترها	برآورد	Std.Err	z-value	Pr(> z)
عرض از مبدا	۱/۱۰۰۲	۰/۱۴۱	۷/۸۲۱	$۵/۲۶e^{-۱۵}***$
ageadj	۰/۰۱۳	۰/۰۰۶	۲/۰۸۲	۰/۰۳۷۳*
chemo	۰/۲۶۶	۰/۱۲۴	۲/۱۲۹	۰/۰۳۳۳*

'***' سطح ۰، '**' سطح ۰/۰۰۱، '*' سطح ۰/۰۵ و '.' سطح ۰/۱ را نشان می‌دهد.

پارامتر پراکندگی متغیر

در این قسمت، مدل را برای حالتی که σ^2 تابعی از متغیرهای تبیینی باشد نیز برازش دادیم که در آن پیشگوه‌های خطی به صورت زیر هستند:

$$\text{logit}(\mu_i) = \beta_0 + \beta_1 \text{ageadj} + \beta_2 \text{chemo} \quad \log(\sigma_i^2) = \gamma_0 + \gamma_1 \text{age} \quad (۲.۵)$$

کد برازش این مدل در R به صورت زیر است

```
data("sdac", package="simplexreg")
sim.glm2 <- simplexreg(rcd~ageadj+chemo|age, link = "logit", data = sdac)
summary(sim.glm2)
```

نتایج برازش پارامترهای مدل (۲.۵) با ۶ تکرار و با انحراف ۲۳۹ در جدول ۴.۵ گزارش شده است. لگاریتم درست‌نمایی برابر $۹۹/۵۱$ و p - مقدار مربوطه معادل $۰/۴۳$ است. p - مقدار آزمون کای دو برای هر دو مدل به ترتیب $۰/۴۹$ و $۰/۴۳$ است و این به معنای نقص برازش نیست.

جدول ۳.۵: باقی مانده استاندارد پیرسون مدل (۲.۵)

کوچکترین مقدار	چارک اول	میانه	چارک سوم	بزرگترین مقدار
-۳/۰۸۲	-۰/۵۳۹	-۰/۰۰۳	۰/۴۹۶	۱/۴۴۸

جدول ۴.۵: برآورد پارامترهای مدل (۲.۵)

پارامترها	برآورد	Std.Err	z-value	Pr(> z)
عرض از مبدا	۱/۱۱۵	۰/۱۴۱	۷/۸۹	$۳/۰۳e^{-۱۵}***$
ageadj	۰/۰۱۳	۰/۰۰۶۴	۲/۰۱۷	۰/۰۴۳۷*
chemo	۰/۲۵۱	۰/۱۲۱	۲/۰۶۸	۰/۰۳۸۶*
عرض از مبدا	۲/۶۰۷	۰/۳۶۶	۷/۱۰۷	$۱/۱۸e^{-۱۲}***$
age	-۰/۰۱۵	۰/۰۰۰۶	-۲/۱۸۱	۰/۰۲۹*

*** سطح ۰، ** سطح ۰/۰۰۱، * سطح ۰/۰۵ و ' ' سطح ۰/۱ را نشان می دهد.

جدول ۵.۵ و ۶.۵، ملاک های AIC و BIC برای دو مدل با σ^2 ثابت و متغیر را نشان می دهند. با توجه به این دو مقدار، مدل (۱.۵)، یعنی مدل با پارامتر پراکندگی ثابت بهتر عمل کرده است.

جدول ۵.۵: مقایسه AIC برای دو مدل

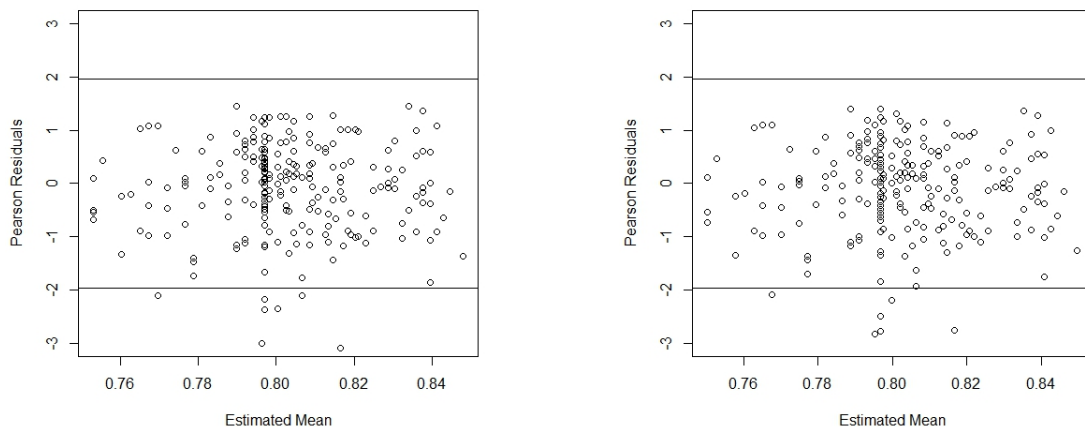
مدل	درجه آزادی	AIC
با پراکندگی ثابت (۱.۵)	۴	-۲۰۲/۵۴۷
با پراکندگی متغیر (۲.۵)	۵	-۱۸۹/۰۳

جدول ۶.۵: مقایسه BIC برای دو مدل

مدل	درجه آزادی	BIC
با پراکندگی ثابت (۱.۵)	۴	-۱۸۸/۶۴۰۸
با پراکندگی متغیر (۲.۵)	۵	-۱۷۱/۶۴۷۷

بررسی فرض ثابت بودن واریانس باقی مانده ها توسط رسم نمودار باقی مانده ها در برابر میانگین مقادیر برازش داده شده، صورت می گیرد. چنانچه این نمودار، ساختار خاصی را نشان ندهد، مثلا به شکل قیفی نباشد، می توان فرض ثابت بودن واریانس باقی مانده ها

را پذیرفت. همان‌طور که در شکل ۱.۵ ملاحظه می‌شود این نمودارها که برای دو مدل با پراکندگی ثابت و متغیر هستند، ساختار خاصی را نشان نمی‌دهند و حدود ۹۷٪ از نقاط در بازه (۱/۹۶, -۱/۹۶) قرار دارند. بنابراین می‌توان فرض ثابت بودن واریانس باقی‌مانده‌ها را پذیرفت.



(آ) بررسی توزیع دقیق باقی‌مانده‌های پیرسون برای مدل با پارامتر پراکندگی ثابت
(ب) بررسی توزیع دقیق باقی‌مانده‌های پیرسون برای مدل با پارامتر پراکندگی متغیر

شکل ۱.۵: بررسی نیکویی برازش مدل‌های رگرسیون سیمپلکس با باقی‌مانده‌های پیرسون، در مثال PBSC

۲.۲.۵ داده‌های جراحی چشم

سانگ و همکاران (۲۰۰۲) داده‌های مربوط به گاز C_3F_8 را برای بررسی تاثیر استفاده از تزریق این گاز درون چشم، در جراحی‌های ترمیم شبکیه که توسط میرز و همکاران (۱۹۹۲) را مجدداً بررسی کردند. داده‌های مربوط به شبکیه، retinal، در بسته simplexreg در نرم‌افزار R گنجانده شده‌است. متغیر پاسخ، درصد حجم گاز باقی‌مانده در چشم است. گاز با سه سطح غلظت متفاوت، ۱۵٪، ۲۰٪ و ۲۵٪ به داخل چشم ۳۱ بیمار قبل از عمل جراحی تزریق شده است. سپس در یک دوره سه ماهه، ۳ تا ۸ بار (به‌طور متوسط ۵ بار) پیگیری شدند و حجم گاز موجود در چشم در دفعات بعدی، به‌عنوان درصد ابتدایی حجم گاز ثبت می‌شود. هدف اولیه، تاثیرگذاری سطح غلظت گاز تزریق‌شده در چشم بیماران بر نرخ فاسد شدن گاز داخل چشم است. این مجموعه داده شامل ۱۸۱ مشاهده (بیماران) و ۶ متغیر تبیینی به شرح زیر است:

- Gas: درصد حجم گاز باقی مانده در چشم
 - Time: زمان بندی روزهای پس از تزریق گاز
 - LogT: لگاریتم زمان ($\log(\text{Time})$)
 - LogT2: توان دوم LogT
 - Level: متغیر ظاهری؛ سطوح غلظت گاز اولیه داخل چشم؛ اگر $\text{Level} = -1$ یعنی سطح غلظت گاز ۱۵٪، $\text{Level} = 0$ سطح غلظت گاز ۲۰٪ و $\text{Level} = 1$ سطح غلظت گاز ۲۵٪
 - ID: یک عامل نشان دهنده بیماران
- بخشی از داده‌های retinal موجود در بسته simplexreg به شرح زیر است. قبل از تعیین مدل، ابتدا ساختار همبستگی را از طریق نمودار خودهمبستگی در تاخیرهای متفاوت بررسی می‌کنیم.

	Gas	Time	LogT	LogT2	Level	ID
1	0.99	1	0	0	0	1
2	0.95	2	0.6931472	0.480453	0	1
3	0.95	3	1.098612	1.206949	0	1
4	0.5	8	2.079442	4.324077	0	1
5	0.375	16	2.772589	7.687248	0	1
6	0.225	33	3.496508	12.22557	0	1
7	0.05	61	4.110874	16.89928	0	1
8	0.999	1	0	0	-1	2
9	0.95	2	0.6931472	0.480453	-1	2
10	0.9	3	1.098612	1.206949	-1	2
11	0.9	4	1.386294	1.921812	-1	2
12	0.55	12	2.484907	6.174761	-1	2
13	0.475	21	3.044522	9.269117	-1	2
14	0.25	29	3.367296	11.33868	-1	2
15	0.05	44	3.78419	14.32009	-1	2
16	0.99	1	0	0	1	3
17	0.99	2	0.6931472	0.480453	1	3
18	0.95	3	1.098612	1.206949	1	3
19	0.9	4	1.386294	1.921812	1	3

پارامتر پراکندگی ثابت

ابتدا مدل خطی تعمیم یافته سیمپلکس بدون همبستگی پیاپی^۷ را برای داده‌های جراحی چشم در نظر می‌گیریم؛ زامین سطح غلظت گاز برای i امین بیمار در روز t_{ij} را با نماد y_{ij} نشان می‌دهیم. متغیرهای ظاهری به صورت زیر بیان می‌شود:

$$x_{ij} = \frac{\text{gas}_{ij} - 20}{5} \begin{cases} 1 & \text{گاز با سطح غلظت } 15\% \\ 0 & \text{گاز با سطح غلظت } 20\% \\ -1 & \text{گاز با سطح غلظت } 25\% \end{cases}$$

⁷ Serial Correlation

با جایگزینی ij توسط k ، متغیر پاسخ ما به صورت $y_k \sim S^-(\mu_k, \sigma^2)$ نشان داده می‌شود که در آن $k = 1, \dots, N$ و $N = \sum_{i=1}^n n_i = 181$ تعداد کل مشاهدات موجود است. پیشگوی خطی برای μ_k به صورت زیر است

$$\text{logit}(\mu_k) = \beta_0 + \beta_1 \log(t_k) + \beta_2 \log^2(t_k) + \beta_3 x_k \quad (3.5)$$

که t_k تعداد روزها پس از تزریق گاز تا زمان مشاهده است. با دستورات زیر می‌توان مدل (3.5) را بر روی این داده‌ها برازش داد؛ در این جا Level همان x_k است.

```
data("retinal", package = "simplexreg")
sim.glm3 <- simplexreg(Gas ~ LogT + LogT2 + Level,
data = retinal, id = ID)
summary(sim.glm3)
```

نتایج برازش مدل (3.5) در جدول 8.5 گزارش شده‌اند. نتایج مذکور با 23 تکرار و با انحراف 177 به دست آمده‌اند که لگاریتم درست‌نمایی آن برابر با 391/2 و p -مقدار 0/486 است. همان‌طور که ملاحظه می‌شود، عرض از مبدا و ضرایب متغیرهای LogT2 و Level معنی‌دار هستند.

جدول 7.5: باقی‌مانده استاندارد پیرسون مدل (3.5)

کوچکترین مقدار	چارک اول	میانه	چارک سوم	بزرگترین مقدار
-4/013	-0/327	0/07	0/441	3/659

جدول 8.5: برآورد پارامترهای مدل (3.5)

پارامترها	برآورد	Std.Err	z-value	Pr(> z)
عرض از مبدا	2/685	0/326	8/231	< 2e-16***
LogT	0/065	0/346	0/188	0/851
LogT2	-0/335	0/0847	-3/959	7/53e-05***
Level	0/325	0/147	2/208	0/027*

***' سطح 0، '**' سطح 0/001، '* ' سطح 0/05 و '.' سطح 0/1 را نشان می‌دهد.

پارامتر پراکندگی متغیر

در این قسمت، مدل را برای حالتی که σ^2 تابعی از متغیرهای ظاهری است برازش می‌دهیم. پیشگوه‌های خطی به صورت زیر هستند:

$$\text{logit}(\mu_k) = \beta_0 + \beta_1 \log(t_k) + \beta_2 \log^2(t_k) + \beta_3 x_k \quad \sigma_k^2 = \sigma_0 + \sigma_1 \log(t_k) + \sigma_2 x_k \quad (4.5)$$

برازش این مدل در R با کد زیر انجام می‌شود

```
sim.glm4 <- simplexreg(Gas ~ LogT + LogT2 + Level | LogT + Level,
  data = retinal, id = ID)
summary(sim.glm4)
```

نتایج مذکور با ۲۰ تکرار و با انحراف ۱۸۱ به دست آمده‌اند که لگاریتم درست‌نمایی آن ۳۷٪ و p -مقدارش ۰/۴۰۲۶ است. در جدول ۱۰.۵ برآورد هر پارامتر، برای مدل حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر نشان داده شده‌است. نتایج نشان می‌دهد که ضریب LogT در ساختار پیشگوی میانگین معنی‌دار نبوده و بایستی حذف شود.

جدول ۹.۵: باقی‌مانده استاندارد پیرسون مدل (۴.۵)

کوچکترین مقدار	چارک اول	میانه	چارک سوم	بزرگترین مقدار
-۴/۳۶۷	-۰/۳۱۹	۰/۰۶۵	۰/۳۹۱	۴/۴۲

جدول ۱۰.۵: برآورد پارامترهای مدل (۴.۵)

پارامترها	برآورد	Std.Err	z-value	Pr(> z)
عرض از مبدا	۲/۷۱۷	۰/۳۵	۷/۷۶۳	$1.5e-15$ ***
LogT	-۰/۰۲	۰/۳۴۶	-۰/۰۵۹	۰/۹۵۲
LogT2	-۰/۳۱۲۶۳	۰/۰۸۱۱۷	-۳/۸۵۱	۰/۰۰۰۱***
Level	۰/۴۰۵	۰/۱۴۵	۲/۷۹۲	۰/۰۰۵**
عرض از مبدا	۶/۱۵۵	۰/۱۹۸	۳۰/۹۵۷	$< 2e-16$ ***
LogT	-۰/۴۵۸	۰/۰۸۰۴	-۵/۷۰۱	$1.19e-08$ ***
Level	-۰/۴۹۳	۰/۱۴۲	-۳/۴۵۶	۰/۰۰۰۵***

*** سطح ۰، ** سطح ۰/۰۰۱، * سطح ۰/۰۵، و . سطح ۰/۱ را نشان می‌دهد.

مدل حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر و ساختار همبستگی

Exc

در این قسمت، پارامتر پراکندگی مدل حاشیه‌ای سیمپلکس را به صورت متغیر در نظر می‌گیریم. همچنین با توجه به ساختار همبستگی تبادل پذیر (Exc) برای بررسی تاثیر زمان بر پارامتر پراکندگی، پیشگوهای خطی به صورت زیر معرفی می‌شوند:

(۵.۵)

$$\text{logit}(\mu_{ij}) = \beta_0 + \beta_1 \log(t_{ij}) + \beta_2 \log^2(t_{ij}) + \beta_3 x_{ij} \quad \sigma_{ij}^2 = \sigma_0 + \sigma_1 \log(t_{ij}) + \sigma_2 \log^2(t_{ij})$$

با اجرای کد زیر، نتایج برازش مدل در ۱۲.۵ آمده است. نتایج مذکور با ۱۷ تکرار و با انحراف ۱۷۷ به دست آمده است.

```
sim.gee1 <- simplexreg(Gas ~ LogT + LogT2 + Level | LogT + Level | Time,
corr = "Exc", id = ID, data = retinal)
summary(sim.gee1)
```

جدول ۱۱.۵: باقی مانده استاندارد پیرسون مدل (۵.۵)

کوچکترین مقدار	چارک اول	میانه	چارک سوم	بزرگترین مقدار
-۳/۸	-۰/۲۵۱	۰/۱۳۶	۰/۴۷۲	۵/۰۴۲

جدول ۱۲.۵: برآورد پارامترها با ساختار همبستگی Exc در مدل (۵.۵)

پارامترها	برآورد	Std.Err	z-value	Pr(> z)
عرض از مبدا	۲/۶۵۴	۰/۳۰۳	۸/۷۶	< ۲e - ۱۶***
LogT	۰/۱۶	۰/۲۴۴	۰/۶۵۷	۰/۵۱۱
LogT2	-۰/۳۷۹	۰/۰۶۴	-۵/۹۱۵	۳/۳۱e - ۰۹***
Level	۰/۲۵	۰/۱۹۹	۱/۲۵۳	۰/۲۱
α	-۱/۳۸۸	۰/۴۸۹	-۲/۸۳۴	۰/۰۰۴۶**
ρ	۰/۲۴۹	۰/۱۲۲	۲/۰۴۱	۰/۰۴۱*

*** سطح ۰، ** سطح ۰/۰۰۱، * سطح ۰/۰۵ و . سطح ۰/۱ را نشان می‌دهد.

نتایج جدول ۱۲.۵ نشان می‌دهد که ضرایب LogT و Level معنی دار نیستند و بایستی از مدل حذف شده و مدل مذکور مجدداً برازش شود.

مدل حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر و ساختار همبستگی

AR(1)

اکنون مدل (۵.۵) را مجدداً با استفاده از ساختار همبستگی اتورگرسیو مرتبه اول در داده‌های طولی، توسط کد زیر برآزش می‌دهیم

```
sim.gee2 <- simplexreg(Gas ~ LogT + LogT2 + Level | LogT + Level | Time,
corr = "AR1", id = ID, data = retinal)
summary(sim.gee2)
```

نتایج مذکور با ۲۵ تکرار و با انحراف ۱۸۱ به‌دست آمده‌اند. در جدول ۱۴.۵ برآورد هر پارامتر، برای مدل حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر و ساختار همبستگی AR(1) نشان داده شده‌است. در این‌جا همانند قبل ضرایب LogT و Level در پیشگوی میانگین و Level در پیشگوی پراکندگی معنی‌دار نیستند و این متغیرها باید از مدل مذکور حذف شوند.

جدول ۱۳.۵: باقی‌مانده استاندارد پیرسون مدل (۵.۵) با ساختار همبستگی AR(1)

کوچکترین مقدار	چارک اول	میانه	چارک سوم	بزرگترین مقدار
-۴/۵۸	-۰/۳۴۵	۰/۰۵۹	۰/۳۹۱	۴/۶۳۷

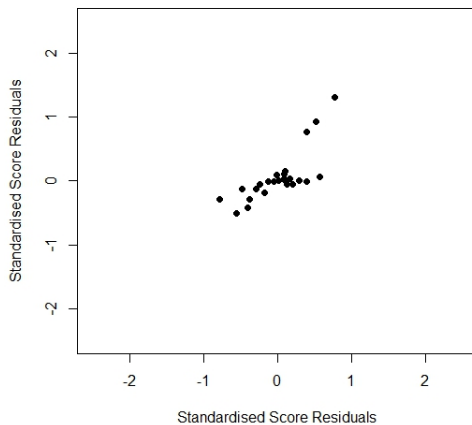
جدول ۱۴.۵: برآورد پارامترها با ساختار همبستگی AR(1) در مدل (۵.۵)

پارامترها	برآورد	Std.Err	z-value	Pr(> z)
عرض از مبدا	۲/۷۲۱	۰/۲۰۳	۱۳/۴۲۵	< ۲e - ۱۶***
LogT	۰/۰۳۴	۰/۳۱۲	۰/۱۰۹	۰/۹۱۳
LogT2	-۰/۳۲۹	۰/۰۸۵	-۳/۸۶۹	۰/۰۰۰۱***
Level	۰/۴۰۹	۰/۲۱۶	۱/۸۸۷	۰/۰۵۹۰
عرض از مبدا	۶/۱۵۳	۰/۳۵۱	۱۷/۵۱۱	< ۲e - ۱۶***
LogT	-۰/۴۵۷	۰/۱۶۹	-۲/۶۹۹	۰/۰۰۷**
Level	-۰/۴۹۲	۰/۳۵۶	-۱/۳۸۱	۰/۱۶۷
α	-۰/۳۴۹	۰/۱۸۶	-۱/۸۷۲	۰/۰۶۱۰
ρ	۰/۷۰۵	۰/۱۳۱	۵/۳۶۳	۸/۱۷e - ۰۸***

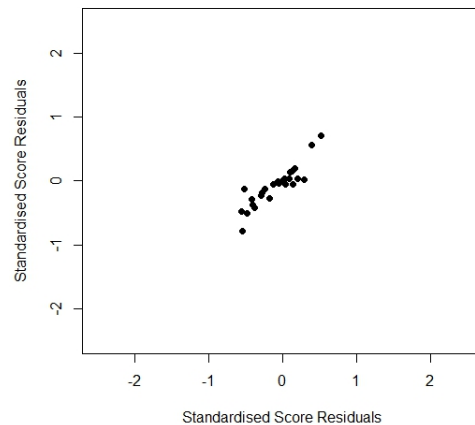
*** سطح ۰، ** سطح ۰/۰۰۱، * سطح ۰/۰۵ و ! سطح ۰/۱ را نشان می‌دهد.

برای بررسی صحت ساختار همبستگی، میزان خودهمبستگی‌ها در چهار تاخیر محاسبه و نتایج

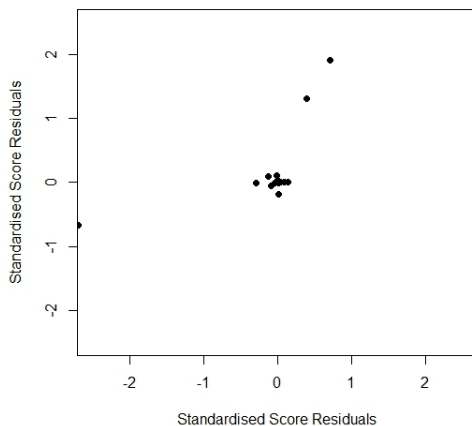
آن در نمودار شکل ۲.۵ رسم شده‌اند؛ که می‌توان دریافت که خودهمبستگی در تاخیر اول، قوی‌تر است و با افزایش تاخیر، خودهمبستگی کاهش می‌یابد، به طوری که در تاخیر چهارم، خودهمبستگی ناچیز می‌شود. در نتیجه مشخص است که ساختار همبستگی $AR(1)$ برای داده‌ها مناسب است. در شکل (آ) ۳.۵، نمودار متغیر \hat{a}_{ij} در برابر پیشگوی خطی η_i نشان



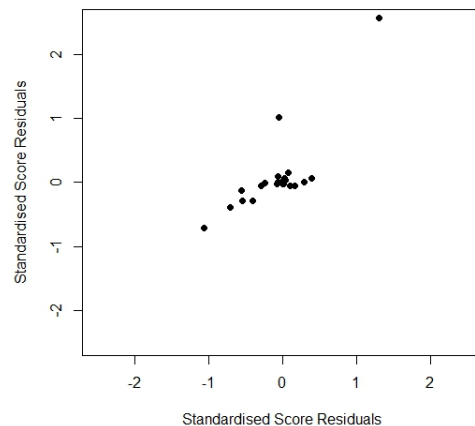
(ب) بررسی خودهمبستگی $lag = 2$



(آ) بررسی خودهمبستگی $lag = 1$



(د) بررسی خودهمبستگی $lag = 4$

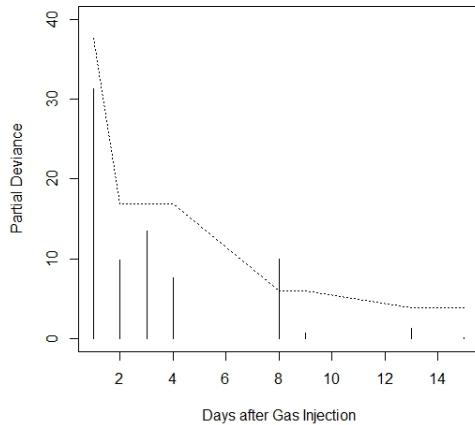


(ج) بررسی خودهمبستگی $lag = 3$

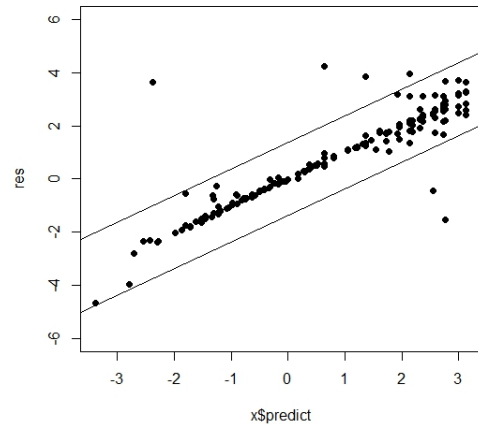
شکل ۲.۵: بررسی خودهمبستگی در چهار تاخیر مختلف، در مثال جراحی چشم.

داده شده است و ۹۷٪ از نقاط در نوار اطمینان ۹۵٪ قرار دارند که نشان‌دهنده یک روند خطی است. هم‌چنین گویای عملکرد تابع پیوند لجیت است که یک انتخاب مناسب برای مدل‌بندی داده‌هاست. شکل (ب) ۳.۵، نشان‌دهنده آزمون نیکویی برازش مدل حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر و ساختار همبستگی اتورگرسیو مرتبه اول است که بیانگر خوب بودن و قابل اعتماد بودن مدل است؛ حتی با توجه به این که در تاخیر هشتم مدل خوب عمل نکرده

است. نتیجه این که مدل حاشیه‌ای سیمپلکس با پارامتر پراکندگی متغیر و ساختار همبستگی اتورگرسیو مرتبه اول عملکرد مناسبی داشته است.



(ب) آزمون نیکویی برازش



(آ) بررسی تابع پیوند

شکل ۳.۵: بررسی تابع پیوند و نیکویی برازش مدل در مثال جراحی چشم

۳.۵ نتیجه‌گیری

در این پایان‌نامه، یک مدل رگرسیونی مناسب برای پاسخ‌هایی که به‌طور پیوسته در فاصله واحد تغییر می‌کنند، ارائه شد. این مدل رگرسیونی بر اساس پذیره پیروی متغیر پاسخ از توزیع سیمپلکس است.

طبق مطالعات شبیه‌سازی انجام‌شده، برآورد پارامترهای مدل در هر دو حالتی که پارامتر پراکندگی ثابت و متغیر باشد، برآوردی سازگار و قابل اعتماد است؛ دقت این برآوردها به واسطه کاهش خطای معیار، افزایش می‌یابد.

ما مدل را علاوه بر مدل بندی پارامتر میانگین، برای حالتی که پارامتر پراکندگی نیز مدل بندی می‌شود، بررسی کردیم. در این پایان‌نامه، مدل‌های متنوعی مورد بررسی قرار گرفتند. در برخی از آن‌ها، پذیره استقلال بین مشاهدات مد نظر قرار گرفت و در برخی دیگر به تحلیل داده‌های طولی که همبستگی جز لاینفک ساختار آن‌ها است، پرداختیم.

به کمک بسته simplexreg در R، تابع چگالی، تابع توزیع تجمعی، تابع چندک و تولید اعداد تصادفی را مورد بررسی قرار داده و به تحلیل باقی‌مانده‌های رگرسیون سیمپلکس با استفاده از روش درست‌نمایی ماکسیمم و معادلات برآوردیابی تعمیم‌یافته پرداختیم.

پیوست آ

۱.آ تولید عدد تصادفی از توزیع MIG

برای تولید متغیر تصادفی از توزیع گوسی معکوس آمیخته از روش پیشنهادی جرجنسن (۱۹۹۱) استفاده می‌کنیم که در آن تابع مولد گشتاور X به صورت زیر می‌باشد:

$$m(s; \xi, \epsilon^2, p) = (1 - p)m_1(s; \xi, \epsilon^2) + pm_2(s; \xi, \epsilon^2)$$

که در آن m_1 و m_2 به ترتیب تابع مولد گشتاورهای X_1 و X_2 می‌باشند. حال از سشادری (۱۹۹۳) داریم

$$m_1(s; \xi, \epsilon^2) = \exp \left\{ \frac{1}{\epsilon} \left[\frac{1}{\xi \epsilon} - \left(\frac{1}{\xi^2 \epsilon^2} - 2s \right)^{\frac{1}{2}} \right] \right\}.$$

و از جانسون و کاتز (۱۹۷۰) داریم

$$m_2(s; \xi, \epsilon^2) = m_1(s; \xi, \epsilon^2)m_3(s; \xi, \epsilon^2),$$

که در آن m_3 برابر است با

$$m_3(s; \xi, \epsilon^2) = (1 - 2\xi^2 \epsilon^2 s)^{-\frac{1}{2}}.$$

m_3 همان تابع مولد گشتاور $\xi^2 \epsilon^2 \chi_{(1)}^2$ است که $X_3 \sim \xi^2 \epsilon^2 \chi_{(1)}^2$ و X_1 مستقل از هم به صورت زیر می‌باشند

$$m(s; \xi, \epsilon^2, p) = m_1(s; \xi, \epsilon^2)[1 - p + pm_3(s; \xi^2 \epsilon^2)].$$

که می‌توان نوشت $X = X_1 + V$ ، که در آن

$$V = \begin{cases} \circ & \text{با احتمال } 1-p \\ X_3 & \text{با احتمال } p \end{cases}$$

و هم‌چنین X_1 و V از هم مستقل و $X_1 \sim IG(\xi, \epsilon^2)$.

۲.۱ ماتریس اطلاع گودامب

ماتریس حساسیت، یک ماتریس (3×3) به صورت زیر است

$$S = E\left\{\frac{\partial \Upsilon(\theta)}{\partial \theta^T}\right\} = \begin{pmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{pmatrix}$$

که در آن $S_{12} = \circ$ ، $S_{13} = \circ$ و $S_{23} = \circ$. مولفه چهارم، S_{21} ، نیز برابر با صفر است زیرا $E[u_{ij}] = \circ$ ؛ در نتیجه می‌توان ماتریس S را به شکل زیر بازنویسی کرد

$$S = \begin{pmatrix} S_{11} & \circ & \circ \\ \circ & S_{22} & \circ \\ S_{31} & S_{32} & S_{33} \end{pmatrix}$$

حال وارون ماتریس فوق به صورت زیر است

$$S^{-1} = \begin{pmatrix} S_{11}^{-1} & \circ & \circ \\ \circ & S_{22}^{-1} & \circ \\ -S_{33}S_{31}S_{11}^{-1} & -S_{33}^{-1}S_{32}S_{22}^{-1} & S_{33}^{-1} \end{pmatrix}$$

باید توجه داشت که تمام مولفه‌های روی قطری اصلی، وارون هستند و زمانی که توزیع $r_{ij}r_{ij}'$ مستقل از پارامترهای میانگین و پراکندگی باشند، مولفه‌های S_{31} و S_{32} نیز برابر با صفرند. مولفه‌های روی قطر اصلی ماتریس S برابرند با

$$S_{11} = -\sum_{i=1}^n D_i^T A_i V_i^{-1} A_i D_i$$

$$S_{22} = -\sum_{i=1}^n \left(\frac{\partial \sigma_i^T}{\partial \gamma}\right) \Sigma_i^{-1} \left(\frac{\partial \sigma_i}{\partial \gamma^T}\right)$$

$$S_{33} = -\sum_{i=1}^n \left(\frac{\partial \xi_i^T}{\partial \alpha}\right) H_i^{-1} \left(\frac{\partial \xi_i}{\partial \alpha^T}\right)$$

ماتریس متغیر R یک ماتریس 3×3 به صورت زیر است

$$R = E\{\Upsilon(\theta)\Upsilon^T(\theta)\} \begin{pmatrix} V_{11} & V_{12} & V_{13} \\ V_{21} & V_{22} & V_{23} \\ V_{31} & V_{32} & V_{33} \end{pmatrix}$$

که مولفه‌های آن در زیر آمده است

$$V_{11} = E\{\Psi_1\Psi_1^T\} = \sum_{i=1}^n D_i^T A_i V_i^{-1} \text{cov}(s_i) V_i^{-1} A_i D_i$$

$$V_{12} = E\{\Psi_1\Psi_2^T\} = \sum_{i=1}^n D_i^T A_i V_i^{-1} \text{cov}(s_i, d_i) \Sigma_i^{-1} \left(\frac{\partial \sigma_i}{\partial \gamma^T}\right)$$

$$V_{13} = E\{\Psi_1\Psi_3^T\} = \sum_{i=1}^n D_i^T A_i V_i^{-1} \text{cov}(s_i, r_i) H_i^{-1} \left(\frac{\partial \xi_i}{\partial \alpha^T}\right)$$

$$V_{22} = E\{\Psi_2\Psi_2^T\} = \sum_{i=1}^n \left(\frac{\partial \sigma_i^T}{\partial \gamma}\right) \Sigma_i^{-1} \text{cov}(d_i) \Sigma_i^{-1} \left(\frac{\partial \sigma_i}{\partial \gamma^T}\right)$$

$$V_{23} = E\{\Psi_2\Psi_3^T\} = \sum_{i=1}^n \left(\frac{\partial \sigma_i^T}{\partial \gamma}\right) \Sigma_i^{-1} \text{cov}(d_i, r_i) H_i^{-1} \left(\frac{\partial \xi_i}{\partial \alpha^T}\right)$$

$$V_{33} = E\{\Psi_3\Psi_3^T\} = \sum_{i=1}^n \left(\frac{\partial \xi_i^T}{\partial \alpha}\right) H_i^{-1} \text{cov}(r_i) H_i^{-1} \left(\frac{\partial \xi_i}{\partial \alpha^T}\right).$$

می‌دانیم که $V_{32} = V_{23}^T$ و $V_{31} = V_{13}^T$ ، $V_{21} = V_{12}^T$

با توجه به این که $\text{cov}(w_i) = \text{diag}\{v(\mu_{ij})\} \text{cov}(u_i) \text{diag}\{v(\mu_{ij})\}$ که برآورد $\text{cov}(u_i)$ از طریق برآورد $\hat{\mu}_{ij}$ ها به دست می‌آید و در این عبارت می‌توان $\text{cov}(u_i)$ را با $\hat{u}_i \hat{u}_i^T$ جایگذاری کرد.

۳.آ دستورات نرم افزار R

شکل ۱.۲

```
x<-seq(0,3,0.01)
f<-(2*pi*x^3)^(-1/2)*(exp((-x-1)^2)/(2*x))
f1<-(2*pi*x^3)^(-1/2)*(exp((-x-2)^2)/(8*x))
f2<-(2*pi*x^3)^(-1/2)*(exp((-x-5)^2)/(50*x))
plot(x, f, type="l", ylim=c(0, 3), col=3, lwd=2)
lines(x, f1, lty=2, col=4, lwd=2)
lines(x, f2, lty=4, col=6, lwd=2)
legend(2.3,3, c(expression(xi==1), expression(xi==2), expression(xi==5)),
```

```

col=c(3,4,6), text.col="black", lty=c(1,2,4) ,lwd=2,
merge=TRUE, bg="white")
xlab='x'
f<-(2*pi*x^3)^(-1/2)*(exp((-x-1)^2)/(2*x))
f1<-(2*(2*pi*x^3))^(-1/2)*(exp((-x-1)^2)/(4*x))
f2<-(5*(2*pi*x^3))^(-1/2)*(exp((-x-1)^2)/(10*x))
plot(x, f, type="l", ylim=c(0,3), col=3, lwd=2)
lines(x, f1, lty=2, col=4, lwd=2)
lines(x, f2, lty=4, col=6, lwd=2)
legend(2.3,3, c(expression(epsilon^2==1), expression(epsilon^2==2),
expression(epsilon^2==5)), col=c(3,4,6), text.col="black",
lty=c(1,2,4), merge=TRUE, bg="white", lwd=2)
xlab='x'

```

شکل ۲.۲

```

library(simplexreg)
f=function(x,a,b){
dsimplex(x,a,b)
}
curve(f(x,0.1,4), c(0,1), col="red", xlim=c(0,1), ylim=c(0,18),
main=expression(mu==0.1), xlab=' ', ylab=expression(sigma^2==4^2))

```

شکل ۳.۲

```

library(simplexreg)
f=function(x,a,b){
dsimplex(x,a,b)
}
curve(f(x,0.1,4), c(0,1),col="red", xlim=c(0,1), ylim=c(0,15),
main=expression(mu==0.1), xlab=' ', ylab=expression(sigma^2==4^2))
par(new=TRUE)
v=sqrt((0.1*0.9)^3)
curve(dnorm(x,0.1,v*4), xlim=c(0,1), lty=2, ylim=c(0,15),
main=' ', xlab=' ', ylab=' ')

```

شکل ۴.۲

```
## or plotting (histogram):
rsimplexplot <- function(n, mu, sigma) {
  x <- seq(0.01, 0.99, 0.01)
  for (j in 1:length(sigma)) {
    for (i in 1:length(mu)) {
      hist(rsimplex(n, mu[i], sigma[j]), density = 50,
           xlim = c(0, 1), ylim = c(0, max(dsimplex(x, mu[i], sigma[j])) + 0.5),
           freq = F, breaks = 20, xlab = " ", ylab = " ", main = NULL)
      lines(x, dsimplex(x, mu[i], sigma[j]), type = "l", col = "red")
    }
  }
  par(mfrow = c(3, 3), mar = c(2, 2, 2, 2), cex = 0.55)
  rsimplexplot(5000, c(0.1, 0.5, 0.7), c(4, 2, 1))
}
```

شکل ۵.۲

```
## plotting (Kernel density):
rsimplexplot <- function(n, mu, sigma) {
  x <- seq(0.01, 0.99, 0.01)
  for (j in 1:length(sigma)) {
    for (i in 1:length(mu)) {
      plot(x, dsimplex(x, mu[i], sigma[j]),
           ylim = c(0, max(dsimplex(x, mu[i], sigma[j])) + 0.5),
           type = "l", col = "red")
      lines(density(rsimplex(n, mu[i], sigma[j])), lty = 2)
    }
  }
  par(mfrow = c(3, 3), mar = c(2, 2, 2, 2), cex = 0.55)
  rsimplexplot(5000, c(0.1, 0.5, 0.7), c(4, 2, 1))
}
```

شکل ۶.۲

```
library(simplexreg)
z<-array(0,dim=c(500,1))
mu=0.5
sigma=0.01
xx=rsimplex(500,mu,sigma)
for (i in 1:500){
z[i]=(xx[i] -mu)/(sigma*sqrt((mu^3)*(1-mu)^3))
}
qqnorm(z,main=expression(sigma^2==0.01),xlab="",ylab="")
qqline(z, col = 2,lwd=2)
```

شکل ۱.۵

```
plot(sim.glm1, type="residuals", res="stdPerr", ylim=c(-3,3))
plot(sim.glm2, type="residuals", res="stdPerr", ylim = c(-3,3))
```

شکل ۲.۵

```
plot(sim.glm3, type="corr", xlim=c(-2.5,2.5), ylim=c(-2.5,2.5), pch=16)
plot(sim.glm3, type="corr", lag=2, xlim=c(-2.5,2.5),
ylim=c(-2.5,2.5), pch=16)
plot(sim.glm3, type="corr", lag=3, xlim=c(-2.5,2.5),
ylim=c(-2.5,2.5), pch=16)
plot(sim.glm3, type="corr", lag=4, xlim=c(-2.5,2.5),
ylim=c(-2.5,2.5), pch=16)
```

شکل ۳.۵

```
plot(sim.gee2, type="residuals", ylim=c(-6,6), pch=16)
plot(sim.gee2, type="GOF", xlab="Days after Gas Injection", ylim=c(0,40))
```

کد مربوط به مثال دوم شبیه‌سازی فصل ۳

```
library(Formula)
```

```

library(simplexreg)
remove(list=ls())
M=500
param = matrix (NA, nrow = M, ncol = 3)
param1 = matrix (NA, nrow = M, ncol = 1)
par = matrix (NA, nrow = M, ncol = 4)
beta0=rep(0,M)
beta1=rep(0,M)
beta2=rep(0,M)
sigma1=rep(0,M)
for (j in 1:M)
{
#Data generation from differnt seeds
set.seed(12345+j)
print(j)
n = 200
beta=c(1.5,0.5,-0.5)
sigma=sqrt(4)
x1 = rep(1,length(n))
x2 = rbinom(n,1,0.6)
x3=runif(n, 0, 5)
X = cbind(x1,x2,x3)
eta = X%*%beta
mu = exp(eta)/(1+exp(eta))
y <- rep(0, n)
for (i in 1:n){
y[i] <- rsimplex(1, mu[i],sqrt(sigma))
}
data=data.frame(cbind(y,X))
try({
# Error in svd(w) : infinite or missing values in 'x'
fit=simplexreg(y ~ x2+x3 , data=data, link="logit")
param[j,] = as.vector(unlist(fit$fixef[,1]))
param1[j,] = as.vector(unlist(fit$Dispersion))

```

```
})
summary(fit)
beta0[j]=fit$fixef[1]
beta1[j]=fit$fixef[2]
beta2[j]=fit$fixef[3]
sigma1[j]=fit$Dispersion
}
#All coefficients in one matrix
A<-matrix(NA,nrow=M,ncol=4)
for(j in 1:3){
A[,1]<-beta0
A[,2]<-beta1
A[,3]<-beta2
A[,4]<-sigma1
}
index<-which((rowSums(abs(A))>4))
A<-A[-index,]
A
#Final estimation of coefficients
beta.and.gamma.hat<-colMeans(A)
beta.and.gamma.hat
#Standard deviation of all coefficients after remove outliers
sd<-GMCM:::colSds(A)
sd
# Computing bias
b0=beta.and.gamma.hat[1]-1.5
b1=beta.and.gamma.hat[2]-0.5
b2=beta.and.gamma.hat[3]+0.5
s=beta.and.gamma.hat[4]-2
#Computing MSE
betaa=c(1.5,0.5,-0.5,2)
Mse=numeric(4)
for (k in 1:4){
mse=function(x) mean((x-betaa[k])^2)
```

```
Mse[k]=mse(A[,k])
}
Mse
#Computing variances
var1=mean((beta0-mean(beta0))^2)
var2=mean((beta1-mean(beta1))^2)
var3=mean((beta2-mean(beta2))^2)
var4=mean((sigma1-mean(sigma1))^2)
f1=(beta0-mean(beta0))/sqrt(var1)
f2=(beta1-mean(beta1))/sqrt(var2)
f3=(beta2-mean(beta1))/sqrt(var3)
f4=(sigma1-mean(sigma1))/sqrt(var4)
#Computing qq-plot
qqnorm(f1,main=expression(beta[0]), xlab="", ylab="")
qqline(f1, col=2, lwd=2)
qqnorm(f2, main=expression(beta[1]), xlab="", ylab="")
qqline(f2, col = 2, lwd=2)
qqnorm(f3, main=expression(beta[2]), xlab="", ylab="")
qqline(f3, col = 2, lwd=2)
qqnorm(f4, main=expression(sigma^2), xlab="", ylab="")
qqline(f4, col = 2, lwd=2)
```

کد مربوط به مثال سوم شبیه‌سازی فصل ۳

```
library(Formula)
library(simplexreg)
remove(list=ls())
M=500
param = matrix (NA, nrow = M, ncol = 3)
param1 = matrix (NA, nrow = M, ncol = 2)
par = matrix (NA, nrow = M, ncol = 5)
beta0=rep(0,M)
beta1=rep(0,M)
beta2=rep(0,M)
```

```
gamma0=rep(0,M)
gamma1=rep(0,M)
for (j in 1:M)
{
#Data generation from differnt seeds
set.seed(12345+j)
print(j)
n = 400
beta=c(1,0.5,-1)
sigma=sqrt(2)
x1 = rep(1,length(n))
x2 = rnorm(n)
x3 = rbinom(n,1,0.6)
X = cbind(x1,x2,x3)
eta = X%*%beta
mu = exp(eta)/(1+exp(eta))
gamma=c(5,0.3)
z1=rep(1,length(n))
z2=rnorm(n)
Z = cbind(z1,z2)
dispersion =Z%*%gamma
sigma1 = exp(dispersion)
y <- rep(0, n)
for (i in 1:n){
y[i] <- rsimplex(1, mu[i], sqrt(sigma1[i]))
}
data=data.frame(cbind(y,X,Z))
try({
#Error in svd(w) : infinite or missing values in 'x'
fit=simplexreg(y ~ x2+x3 | z2, data=data, link="logit")
param[j,] = as.vector(unlist(fit$fixef[,1]))
param1[j,] = as.vector(unlist(fit$dispar[,1]))
})
summary(fit)
```



```

beta0[j]=fit$fixef[1]
beta1[j]=fit$fixef[2]
beta2[j]=fit$fixef[3]
gamma0[j]=fit$dispar[1]
gamma1[j]=fit$dispar[2]
}
#All coefficients in one matrix
A<-matrix(NA,nrow=M,ncol=5)
for(j in 1:5){
A[,1]<-beta0
A[,2]<-beta1
A[,3]<-beta2
A[,4]<-gamma0
A[,5]<-gamma1
}
index<-which((rowSums(abs(A))>10))
A<-A[-index,]
A
#Final estimation of coefficients
beta.and.gamma.hat<-colMeans(A)
beta.and.gamma.hat
# Standard deviation of all coefficients after remove outliers
sd<-GMCM:::colSds(A)
sd
# Computing bias
b0=beta.and.gamma.hat[1]-1
b1=beta.and.gamma.hat[2]-0.5
b2=beta.and.gamma.hat[3]+1
g0=beta.and.gamma.hat[4]-5
g1=beta.and.gamma.hat[5]-0.3
#Computing MSE
betaa=c(1, 0.5,- 1, 5, 0.3)
Mse=numeric(5)
for (k in 1:5){

```

```
mse=function(x) mean((x-betaa[k])^2)
Mse[k]=mse(A[,k])
}
Mse
#Computing variances
var1=mean((A[,1]-mean(A[,1]))^2)
var2=mean((A[,2]-mean(A[,2]))^2)
var3=mean((A[,3]-mean(A[,3]))^2)
var4=mean((A[,4]-mean(A[,4]))^2)
var5=mean((A[,5]-mean(A[,5]))^2)
f1=(A[,1]-mean(A[,1]))/sqrt(var1)
f2=(A[,2]-mean(A[,2]))/sqrt(var2)
f3=(A[,3]-mean(A[,3]))/sqrt(var3)
f4=(A[,4]-mean(A[,4]))/sqrt(var4)
f5=(A[,5]-mean(A[,5]))/sqrt(var5)
#Computing qq-plot
qqnorm(f1, main=expression(beta[0]), xlab="", ylab="")
qqline(f1, col=2, lwd=2)
qqnorm(f2, main=expression(beta[1]), xlab="", ylab="")
qqline(f2, col=2, lwd=2)
qqnorm(f3, main=expression(beta[2]), xlab="", ylab="")
qqline(f3, col=2, lwd=2)
qqnorm(f4, main=expression(gamma[0]), xlab="", ylab="")
qqline(f4, col=2, lwd=2)
qqnorm(f5, main=expression(gamma[1]), xlab="", ylab="")
qqline(f5, col=2, lwd=2)
```

مراجع

- [1] Allan, D. S., Keeney, M., Howson-Jan, K., Popma, J., Weir, K., Bhatia, M., Sutherland, D. & Chin-Yee, I. H. (2002). Number of viable CD34+ cells reinfused predicts engraftment in autologous hematopoietic stem cell transplantation. *Bone marrow transplantation*, 29(12), 967.
- [2] Balakrishnan, N., Leiva, V., Sanhueza, A., & Cabrera, E. (2009). Mixture inverse Gaussian distributions and its transformations, moments and applications. *Statistics*, 43(1), 91-104.
- [3] Balemi, A., & Lee, A. (2009). Comparison of GEE1 and GEE2 estimation applied to clustered logistic regression. *Journal of Statistical Computation and Simulation*, 79(4), 361-378.
- [4] Barndorff-Nielsen, O. E., & Jørgensen, B. (1991). Some parametric models on the simplex. *Journal of Multivariate Analysis*, 39(1), 106-116.
- [5] Chambers, J. M., & Hastie, T. J. (Eds.). (1992). *Statistical models in S*, (251). Pacific Grove, CA: Wadsworth & Brooks/Cole Advanced Books & Software.
- [6] Chhikara, R. (1988). *The Inverse Gaussian Distribution: Theory: Methodology, and Applications*, (95). CRC Press.
- [7] Crowder, M. (1987). On linear and quadratic estimating functions. *Biometrika*, 74(3), 591-597.
- [8] Diggle, P., Diggle, P. J., Heagerty, P., Heagerty, P. J., Liang, K. Y., & Zeger, S. (2002). *Analysis of Longitudinal Data*. Oxford University Press.
- [9] Ferrari, S., & Cribari-Neto, F. (2004). Beta regression for modelling rates and proportions. *Journal of Applied Statistics*, 31(7), 799-815.

-
- [10] Fisher, R. A. (1953). Dispersion on a sphere. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 217(1130), 295-305.
- [11] Godambe, V. P. (1960). An optimum property of regular maximum likelihood estimation. *The Annals of Mathematical Statistics*, 31(4), 1208-1211.
- [12] Hardin, J. W., & Hilbe, J. M. (2002). *Generalized Estimating Equations*. Chapman and Hall/CRC, London.
- [13] Heyde, C. C.(1997). *Quasi-likelihood and its Application: A General Approach to Optimal Parameter Estimation*. Springer, New York.
- [14] Horton, N. J., & Lipsitz, S. R. (1999). Review of software to fit generalized estimating equation regression models. *The American Statistician*, 53(2), 160-169.
- [15] Johnson, N. L., & Kotz, S. (1970). *Distributions in Statistics: Continuous Univariate Distributions-1*. Houghton Mifflin.
- [16] Jørgensen, B., Seshadri, V., & Whitmore, G. A. (1991). On the mixture of the inverse Gaussian distribution with its complementary reciprocal. *Scandinavian Journal of Statistics*, 77-89.
- [17] Jørgensen, B. (1992). Exponential dispersion models and extensions: A review. *International Statistical Review/Revue Internationale de Statistique*, 5-20.
- [18] Jorgensen, B. (1997). *The Theory of Dispersion Models*. CRC Press.
- [19] Kedem, B., & Fokianos, K. (2005). *Regression Models for Time Series Analysis* (488). John Wiley & Sons.
- [20] Liang, K. Y., & Zeger, S. L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika*, 73(1), 13-22.
- [21] McCullagh, P., Nelder, JA. (1989),. *Generalized Linear Models*. Chapman and Hall, London.
- [22] Meyers, S. M., Ambler, J. S., Tan, M. I. N. G., Werner, J. C., & Huang, S. S. (1992). Variation of perfluoropropane disappearance after vitrectomy. *Retina (Philadelphia, Pa.)*, 12(4), 359-363.
- [23] Michael, J. R., Schucany, W. R., & Haas, R. W. (1976). Generating random variates using transformations with multiple roots. *The American Statistician*, 30(2), 88-90.

- [24] Mudholkar, G. S., & Natarajan, R. (2002). The inverse Gaussian models: analogues of symmetry, skewness and kurtosis. *Annals of the Institute of Statistical Mathematics*, 54(1), 138-154.
- [25] Mudholkar, G. S., & Wang, H. (2007). IG-symmetry and R-symmetry: interrelations and applications to the inverse Gaussian theory. *Journal of Statistical Planning and Inference*, 137(11), 3655-3671.
- [26] Nelder, J. A., & Wedderburn, R. W. (1972). Generalized linear models. *Journal of the Royal Statistical Society: Series A (General)*, 135(3), 370-384.
- [27] Paik, M. C. (1992). Parametric variance function estimation for nonnormal repeated measurement data. *Biometrics*, 19-30.
- [28] Park, T., & Shin, D. Y. (1999). On the use of working correlation matrices in the GEE approach for longitudinal data. *Communications in Statistics-Simulation and Computation*, 28(4), 1011-1029.
- [29] Prentice, R. L. (1988). Correlated binary regression with covariates specific to each binary observation. *Biometrics*, 1033-1048.
- [30] Prentice, R. L., & Zhao, L. P. (1991). Estimating equations for parameters in means and covariances of multivariate discrete and continuous responses. *Biometrics*, 825-839.
- [31] Qiu, Z. (2001). *Simplex Mixed Models for Longitudinal Proportional Data*, Ph.D. Dissertation, York University
- [32] Qiu, Z., Song, P. X. K., & Tan, M. (2008). Simplex mixed-effects models for longitudinal proportional data. *Scandinavian Journal of Statistics*, 35(4), 577-596.
- [33] Seshadri, V. (1993),. *The Inverse Gaussian Distribution*. New York.
- [34] Seshadri, V. (1999),. *The Inverse Gaussian Distribution: Statistical Theory and Applications*(137). Springer Science& Business Media New York
- [35] Shuster, J. (1968). On the inverse Gaussian distribution function. *Journal of the American Statistical Association*, 63(324), 1514-1516.
- [36] Song, P. X. K., & Tan, M. (2000). Marginal models for longitudinal continuous proportional data. *Biometrics*, 56(2), 496-502.

-
-
- [37] Song, P. X. K., Qiu, Z., & Tan, M. (2004). Modelling heterogeneous dispersion in marginal models for longitudinal proportional data. *Biometrical Journal: Journal of Mathematical Methods in Biosciences*, 46(5), 540-553.
- [38] Wedderburn, R. W. (1974). Quasi-likelihood functions, generalized linear models, and the Gauss—Newton method. *Biometrika*, 61(3), 439-447.
- [39] Yang, H., Acker, J. P., Cabuhat, M., Letcher, B., Larratt, L., & McGann, L. E. (2005). Association of post-thaw viable CD34+ cells and CFU-GM with time to hematopoietic engraftment. *Bone Marrow Transplantation*, 35(9), 881.
- [40] Zeileis, A., Cribari-Neto, F., Grün, B., & Kosmidis, I. (2010). Beta regression in R. *Journal of Statistical Software*, 34(2), 1-24.
- [41] Zhang, P., & Qiu, Z. (2014). Regression analysis of proportional data using simplex distribution. *Scientia Sinica Mathematica*, 44(1), 89-104.
- [42] Zhang, P., Qiu, Z., & Shi, C. (2016). simplexreg: An r package for regression analysis of proportional data using the simplex distribution. *Journal of Statistical Software*, 71(11), 1-21.
- [43] Zorn, C. J. (2001). Generalized estimating equation models for correlated data: A review with applications. *American Journal of Political Science*, 470-490.

Abstract

In regression problems, Logistic and Probit models are commonly used for proportional response variable in unit interval $(0, 1)$. But when the distribution of proportional response is highly skewed, these two models are not suitable, because their underlying distribution *i.e.* Logistic and Normal are symmetric, hence a contradiction is made. In this thesis, "Simplex regression" is introduced as an alternative solution. Firstly, Simplex distribution based on mixture inverse Gaussian distribution is introduced. Secondly, Simplex simulation studies are made to assert the capabilities of the distribution. As Simplex regression can be defined on the framework of "Dispersion models", we explain it first, thereafter, similar to generalized linear models methods, the regression parameters are estimated in two ways: heterogeneous and homogeneous variance. Modeling the longitudinal data by Simplex regression is the last subject in this note. All models are also evaluated by simulation and real examples.

Keywords: Dispersion models, Proportional data, Random variable generation, Generalized linear models, Simplex distribution, Simplex regression models, Longitudinal data, Generalized estimating equations.



Faculty Of Mathematical Sciences

MSc Thesis in Statistics

Regression Analysis of Proportion Data Using the Simplex Distribution

By: Sima Maleki

Supervisor:
Davood Shahsavani

January 2019