



دانشگاه علامه طباطبائی
دانشکده اقتصاد
گروه آمار، ریاضی و کامپیوتر
پایان نامه برای دریافت درجه کارشناسی ارشد آمار ریاضی

عنوان

تعیین برآورد پارامترهای معادله‌های برآوردساز تعمیم یافته در حضور داده‌های گم شده

پژوهشگر

مهدی عندلیب خواه

استاد راهنما

دکتر فرزاد اسکندری

استاد مشاور

دکتر حمیدرضا نواب پور

تیر ماه ۱۳۹۱

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

کلیه‌ی حقوق مادی و معنوی اعم از چاپ و تکثیر، نسخه‌برداری، ترجمه، اقتباس و ... از این پایان‌نامه
برای دانشگاه علامه طباطبائی محفوظ است. نقل مطالب با ذکر منبع مانعی ندارد.

تأیید پایان‌نامه‌ی کارشناسی ارشد توسط دانشجو

عنوان پایان‌نامه: تعیین برآورد پارامترهای معادله‌های برآوردساز تعمیم‌یافته در حضور داده‌های گم‌شده

نام دانشجو: مهدی عندلیب خواه

شماره‌ی دانشجویی: ۸۸۱۲۵۱۲۸۲۰۶

استاد راهنما: دکتر فرزاد اسکندری

این جانب مهدی عندلیب خواه دانشجوی کارشناسی ارشد رشته‌ی آمار ریاضی دانشکده‌ی اقتصاد دانشگاه علامه طباطبائی گواهی می‌نمایم پژوهش‌های ارائه شده در پایان‌نامه با عنوان مذکور توسط شخص این جانب انجام شده است و درستی مطالب نگارش یافته مورد تأیید می‌باشد. همچنین گواهی می‌نمایم مطالب مندرج در پایان‌نامه تاکنون برای دریافت هیچ نوع مدرک یا امتیازی توسط این جانب یا فرد دیگری در هیچ کجا ارائه نشده است و در نگارش متن پایان‌نامه شیوه‌ی نگارش مصوب دانشکده‌ی اقتصاد را به‌طور کامل رعایت نموده‌ام. چنان‌چه در هر زمان خلاف آن‌چه گواهی نموده‌ام مشاهده گردد خود را از آثار حقیقی و حقوقی ناشی از دریافت مدرک کارشناسی ارشد محروم می‌دانم و هیچ‌گونه ادعایی نخواهم داشت.

امضا دانشجو:

تاریخ:

تقدیم بہ

زیساترین ہمراہ راہ اندیت

سپاسگزاری

پیش‌تر از همه از خانواده‌ی عزیزم که همواره بی دریغ یاور و پشتیبانم بوده‌اند، از پدرم، مادرم که فداکارانه خودشان را وقف آینده‌ی فرزندان‌شان نموده‌اند، بی نهایت سپاسگزارم .

از دوستان محترم آقایان اکبر حیدریگی، مهرداد مددی، محمدفیاض و مهدی فیاض‌بخش به پاس توصیه‌های دلسوزانه ایشان نسبت به اینجانب سپاسگزارم.

از اساتید گرانقدر دست اندرکار پایان‌نامه که مهربانانه و صمیمانه دوشادوش، با صرف دقت، زمان و زحمات فراوان کار را آن‌چنان که باید سامان بخشیدند، دکتر فرزاد اسکندری راهنمای محترم پایان‌نامه، دکتر حمیدرضا نواب‌پور مشاور ارجمند، دکتر محمدرضا صالحی راد نماینده‌ی تحصیلات تکمیلی و دکتر خلیل شفیع‌ی داور جلسه‌ی دفاع نهایت تشکر را دارم.

امیدوارم بتوانم از عهده ادای حق این عزیزان برآیم.

فهرست مطالب

ب	فهرست مطالب
ث	فهرست جدول‌ها
ج	فهرست شکل‌ها
۱	۱ کلیات
۲	۱-۱ مقدمه
۳	۲-۱ مسئله‌ی مورد بررسی
۵	۳-۱ تاریخچه‌ی مسأله‌ی مورد بررسی
۵	۴-۱ اطلاع فیشر
۶	۵-۱ خانواده توزیع‌های نمایی
۷	۶-۱ روش امتیازدهی فیشر
۸	۷-۱ الگوریتم نیوتون-رافسون
۹	۸-۱ برازش مدل‌های خطی
۹	۱-۸-۱ مدل خطی کلاسیک
۱۱	۲-۸-۱ مدل حاشیه‌ای
۱۱	۳-۸-۱ مدل خطی با اثرهای تصادفی
۱۲	۴-۸-۱ مدل‌های خطی تعمیم یافته
۱۴	۵-۸-۱ تحلیل مانده‌ها
۱۵	۶-۸-۱ برآورد پارامتر پراکنندگی
۱۵	۹-۱ انواع مکانیسم‌های گم‌شدگی

۱۷	۱-۹-۱	مکانیسم گم شدگی MAR
۱۷	۲-۹-۱	مکانیسم گم شدگی MCAR
۱۷	۳-۹-۱	مکانیسم گم شدگی NMAR
۱۷	۱۰-۱	چشم‌انداز فصل‌های آینده
۱۸		۲	معادله‌های براوردساز تعمیم‌یافته
۱۹	۱-۲	مقدمه
۱۹	۲-۲	روش معادله‌های براوردساز تعمیم‌یافته
۲۰	۳-۲	برآورد ماتریس واریانس-کوواریانس
۲۲	۴-۲	انواع معادله‌های براوردساز تعمیم‌یافته
۲۲	۱-۴-۲	مدل جمعیت متوسط (PA)
۲۳	۲-۴-۲	مدل موضوع ویژه (SS)
۲۵	۵-۲	انواع ساختارهای ماتریس همبستگی در روش <i>GEEs</i>
۲۶	۱-۵-۲	ساختار همبستگی ثابت
۲۷	۲-۵-۲	ساختار همبستگی مستقل
۲۷	۳-۵-۲	ساختار همبستگی تبادل‌پذیر
۲۸	۴-۵-۲	ساختار همبستگی اتورگرسیو مرتبه‌ی اول
۲۸	۵-۵-۲	ساختار همبستگی مانا
۲۹	۶-۵-۲	ساختار همبستگی بی‌ساختار
۳۰	۷-۵-۲	نکته‌هایی در زمینه انتخاب بهترین ساختار همبستگی کاری
۳۰		۶-۲	روش معادله‌های براوردساز تعمیم‌یافته مرتبه‌ی اول <i>GEE1</i> و مرتبه‌ی دوم <i>GEE2</i>
۳۳	۱-۶-۲	مقایسه‌ی روش‌های <i>GEE1</i> و <i>GEE2</i>
۳۴	۷-۲	خلاصه‌ی فصل
۳۵		۳	داده‌های گم‌شده و معادله‌های براوردساز تعمیم‌یافته
۳۶	۱-۳	مقدمه
۳۶	۲-۳	روش معادله‌های براوردساز موزون
۴۰	۳-۳	روش تقریب توزیع چندمتغیره‌ی نرمال
۴۳	۴-۳	روش‌های جانهای در رویکرد <i>GEE</i>

۴۳ روش جانهی ساده در رویکرد <i>GEE</i>	۱-۴-۳
۴۵ روش جانهی چندگانه در رویکرد <i>GEE</i>	۲-۴-۳
۴۹ آگوریتم <i>AU</i>	۵-۳
۵۲ خلاصهی فصل	۶-۳
۵۳		۴ کاربرد
۵۴ مقدمه	۱-۴
۵۴ یک کاربرد	۲-۴
۵۶ ساختار داده‌ها و فرض‌های مورد نظر	۳-۴
۵۶ شبیه‌سازی	۴-۴
۶۶ نتیجه‌گیری	۵-۴
۶۷		کتاب‌نامه
۷۱		واژه‌نامه‌ی فارسی به انگلیسی
۷۶		پیوست ساختار کتابخانه‌های استفاده‌شده
۷۶ کتابخانه <i>geepack</i>	۶-۴
۷۷ کتابخانه <i>FastImputation</i>	۷-۴
۷۸ کتابخانه <i>mi</i>	۸-۴

فهرست جدول‌ها

- ۱-۴ برآورد پارامتر رگرسیونی β برای داده‌های آزمایش بالینی $CD4$ ۵۵
- ۲-۴ آماره‌های توصیفی متغیر تصادفی Y ۵۶
- ۳-۴ برآورد پارامترهای رگرسیونی خودگردان به شیوه‌های مختلف برای داده‌های
شبیه‌سازی- $AR(1)$ ساختار همبستگی اتورگرسیو مرتبه‌ی اول، UN ساختار
همبستگی بی‌ساختار، $Exch$ ساختار همبستگی قابل‌تغییر، Ind ساختار همبستگی
مستقل ۶۰
- ۴-۴ خطای استاندارد برآورد پارامترهای رگرسیونی به شیوه‌های مختلف برای
داده‌های شبیه‌سازی ۶۱

فهرست شکل‌ها

۵۷	نمودارهای همگرایی میانگین و انحراف معیار.	۱-۴
۵۸	نمودار بافت‌نگار جانهی	۲-۴
۵۹	نمودار پراکنش	۳-۴
		نمودار مقایسه‌ی مدل اصلی با مدل برازش داده‌شده به وسیله‌ی (a) روش	۴-۴
۶۲	مورد-کامل و (b) روش وزن‌دهی	
		نمودار مقایسه‌ی مدل اصلی با مدل برازش داده‌شده به وسیله‌ی (a) روش	۵-۴
۶۳	جانهی میانگین و (b) روش جانهی ساده	
		نمودار مقایسه‌ی مدل اصلی با مدل برازش داده‌شده به وسیله‌ی (a) روش	۶-۴
۶۴	جانهی چندگانه و (b) روش خودگردان بیزی	
		نمودار مقایسه‌ی مدل اصلی با مدل برازش داده‌شده به وسیله‌ی (a) روش	۷-۴
۶۵	تقریب نرمال و (b) روش الگوریتم AU	

چکیده

برای پاسخ به بسیاری از پرسش‌های علمی، نیاز به گردآوری داده‌ها است و گردآوری داده‌ها به تنهایی پاسخ‌گوی بسیاری از پرسش‌ها نیست. بنابراین امروزه تحلیل داده‌ها در تمام شاخه‌های علمی امری ضروری به شمار می‌آید. طی چند دهه‌ی اخیر تحلیل داده‌های رده‌بندی در حوزه‌های آمار رسمی و علوم پزشکی به طور چشم‌گیری متحول شده است. با توجه به کثرت وجود این داده‌ها، بدیهی است که تحلیل و مدل‌بندی آن‌ها می‌تواند تأثیر عمده‌ای در راستای تولید دانش آماری داشته باشد. اما در تحلیل این داده‌ها با دو مشکل عمده روبرو هستیم. داده‌های همبسته و داده‌های گم‌شده. این دو مشکل فرض آماری مشاهده‌های مستقل در روش‌های رگرسیونی سنتی را نقض می‌کند. با وجود این که اکثر مطالعه‌ها برای گردآوری تمام داده‌ها طراحی می‌شوند ولی بروز گم‌شدگی در داده‌ها نیز اجتناب ناپذیر خواهد بود. مکانیسم‌های گم‌شدگی یک مسئله‌ی حائز اهمیت است زیرا خصوصیات روش‌های برخورد با گم‌شدن داده‌ها به این مکانیسم‌ها مربوط می‌شود. با در نظر گرفتن یک متغیر نشانگر برای وضعیت گم‌شدگی و توزیعی برای آن، سه مکانیسم گم‌شدگی تصادفی، کاملاً تصادفی و غیرتصادفی برای داده‌ها تعریف می‌شوند. در چنین شرایطی روش معادله‌های براوردساز تعمیم یافته معرفی می‌شود. این رویکرد یکی از مناسب‌ترین روش‌ها را برای تحلیل فراهم می‌کند. در این معادله‌ها تنها فرض‌هایی که اختیار می‌شوند، فرض درباره‌ی امیدریاضی حاشیه‌ای مرتبه‌ی اول و دوم پاسخ‌ها است و هیچ فرضی درباره‌ی توزیع کامل آنها اختیار نمی‌شود. حال مسئله‌ی مورد نظر ما تعریف می‌شود. با چه روش‌هایی بر پایه *GEE* می‌توان به براورد پارامترها در حضور داده‌های گم‌شده پرداخت؟ از مهمترین روش‌ها در رویارویی با داده‌های گم‌شده روش جانهی است که این روش و انواع مختلف آن از جمله جانهی ساده، جانهی میانگینی و جانهی چندگانه را بیان می‌کنیم. روش دوم روش موسوم به وزندهی است که در آن، سهم هر مشاهده در معادله براوردساز برابر با معکوس احتمال مشاهده شده‌ی آن است. روش بعدی استفاده از تقریب توزیع چندمتغیره‌ی نرمال است که از براورد نرمال برای پارامترهای همبستگی بهره می‌برد. روش دیگر الگوریتم *AU* است که یک روش تکراری را برای حل معادله‌های براوردساز تعمیم‌یافته در حضور داده‌های گم‌شده ارائه می‌کند. در پایان تلاش خواهیم کرد با یک مطالعه شبیه‌سازی روش‌های گوناگون را مقایسه کنیم. در مجموع بر حسب آریبی و خطای استاندارد نتایج شبیه‌سازی نشان می‌دهند که روش الگوریتم *AU* از سایر روش‌ها در مواجهه با داده‌های گم‌شده بهتر عمل می‌کند. **واژگان کلیدی.** معادله‌های براوردساز تعمیم یافته، ماتریس همبستگی کاری، براوردگر استوار واریانس، داده‌های گم‌شده.

فصل ۱ کلیات

۱-۱ مقدمه

برای پاسخ به بسیاری از پرسش‌های علمی، نیاز به گردآوری داده‌ها است. این پرسش‌های علمی در شاخه‌های مختلف علوم به شکل‌های مختلف پیش می‌آیند. برای مثال در پزشکی تأثیر چند نوع داروی مختلف بر یک بیماری مورد ارزیابی قرار می‌گیرد و در اقتصاد، تأثیر عوامل مختلف بر تورم مورد کاوش قرار می‌گیرند. در چنین کاربردهایی، نیاز به گردآوری داده‌ها امری ضروری به شمار می‌آید. پس از گردآوری داده‌ها خلاصه‌سازی آنها که نیاز به جدول‌بندی دارد، صورت می‌گیرد. پس از جدول‌بندی نیز، استخراج اطلاعات توصیفی و در پایان تحلیل این داده‌ها صورت می‌پذیرد.

تحلیل داده‌های رده‌بندی در حوزه‌های آمار رسمی، علوم اجتماعی و علوم پزشکی طی چند دهه‌ی اخیر به طور چشم‌گیری متحول شده است. با توجه به کثرت وجود داده‌های رده‌بندی در مطالعات مقطعی و طولی، بدیهی است که تحلیل و مدل‌بندی آنها می‌تواند تأثیر عمده‌ای در راستای تولید دانش آماری داشته باشد. روش‌های مورد استفاده در تحلیل این داده‌ها از نیمه‌ی دوم قرن بیستم گسترش پیدا کرده‌اند. این روش‌ها برای تحلیل داده‌هایی که در مقیاس‌های رده‌بندی شده اندازه‌گیری می‌شوند مانند مقیاس‌های ترتیبی و اسمی اهمیت بسزایی دارند. اما در تحلیل این داده‌ها با دو مشکل عمده روبرو هستیم.

اول- داده‌های همبسته: داده‌های همبسته در مطالعات علوم اجتماعی بسیار معمول هستند. داده‌های طولی یا به صورت سلسله‌مراتبی ساماندهی شده، موقعیت‌های تحلیلی را ارائه می‌دهند که در آن داده‌ها همبسته‌اند. مثال کلاسیک در این زمینه، دانش‌آموزان دسته‌بندی شده درون کلاس‌ها در یک مدرسه هستند. این موضوع فرض آماری مشاهدات مستقل در روش‌های رگرسیون سنتی را نقض می‌کند.

دوم- داده‌های گم‌شده: به بیان ساده، مشکل داده‌های گم‌شده از تفاوت میان داده‌هایی که ما تصمیم به گردآوری داریم و داده‌هایی که گردآوری می‌کنیم ناشی می‌شود. به عنوان مثال در مطالعات طولی، این مشکل به علت عدم حضور فرد مورد آزمایش در قرار ملاقات و دست کشیدن از مشارکت در مطالعه یا خرابی تجهیزات گردآوری و در مطالعات پزشکی به علت عدم بهبودی و از دنیا رفتن فرد مورد آزمایش رخ خواهد داد. روش‌های برخورد با مسئله داده‌های گم‌شده در این گونه از مطالعات به طور پیوسته در حال پیشرفت می‌باشند. در این پایان‌نامه به بررسی یکی از مناسب‌ترین تکنیک‌ها برای مواجهه با این دو مشکل، یعنی روش معادله‌های برآوردساز تعمیم‌یافته پرداخته

می شود.

۲-۱ مسئله‌ی مورد بررسی

هدف اولیه از تحلیل در بیشتر مدل‌های آماری، بررسی تأثیر متغیرهای کمکی (مستقل) معین بر روی متغیر پاسخ (وابسته) است. با یک مشاهده برای هر فرد، مدل‌های خطی تعمیم‌یافته برای تحلیل می‌توانند به کار گرفته شوند، اما با مشاهدات مکرر همبستگی بین مقادیر نیز باید در تحلیل گنجانده شود. نتایج نظری برای برآوردیابی پارامترها در مدل‌های رگرسیونی، به ویژه برای توزیع نرمال و متغیرهای پیوسته در دسترس هستند ولی برای داده‌های رده‌بندی کمتر توسعه یافته‌اند. یکی از مشکلات در تحلیل داده‌های غیرنرمال، فقدان کلاس غنی از مدل‌ها مانند توزیع چندمتغیره‌ی نرمال برای توزیع توأم مشاهدات است. اعتبار این تحلیل‌ها بستگی به برقراری فرض‌های بیان شده در مدل‌های فوق دارد، اما این فرض‌های کلاسیک اغلب در شرایط تجربی به طور مثال در علوم پزشکی مشاهده نمی‌شوند زیرا در چنین موقعیت‌هایی بررسی‌های مکرر از یک پاسخ در زمان‌های مختلف صورت می‌گیرد.

به عنوان مثال در تحلیل داده‌های شمارشی، پاسخ‌های دودویی و ترتیبی به دلیل وجود وابستگی بین واحدهای آزمایشی و یا برقرار نبودن بسیاری از فرض‌ها همانند نرمال بودن توزیع مشاهدات، امکان استفاده از روشهای رگرسیون معمول مانند کمترین توان‌های دوم وجود ندارد. به دلیل گسترش و توسعه روزافزون مطالعه‌های طولی چشم‌پوشی از همبستگی‌های درون-واحدی در برازش مدل‌های مناسب می‌تواند پژوهشگر را به استنباط‌های نادرست هدایت نماید.

در چنین شرایطی رویکرد معادله‌های برآوردساز تعمیم‌یافته (GEE) اولین بار توسط لیانگ و زیگر (۱۹۸۶) ارائه گردید و سپس توسط محققان دیگری توسعه یافت. این رویکرد یکی از مناسب‌ترین روشها را برای تحلیل این‌گونه از داده‌ها یعنی پاسخ‌های شمارشی، دودویی و ترتیبی فراهم می‌کند. به طور کلی می‌توان گفت این رویکرد اصلاح مدل‌های خطی تعمیم‌یافته برای داده‌های همبسته است. معادله‌های برآوردساز تعمیم‌یافته می‌توانند برای حل دامنه وسیعی از مسائل آماری به کار روند و یک چارچوب متحد و هم‌شکل برای بسیاری از تکنیک‌های برآورد آشنا شامل ماکزیمم درست‌نمایی، کمترین توان‌های دوم، شبه‌درست‌نمایی و برآوردهای ناریب با کمترین واریانس را فراهم کنند. در این رویکرد تنها فرض‌هایی که اختیار می‌شوند، فرض درباره امیدریاضی حاشیه‌ای مرتبه اول و دوم پاسخ‌ها است و درباره توزیع کامل یا درست‌نمایی آن‌ها فرضی اختیار نمی‌شود. این رویکرد برای پارامترهای مدل رگرسیونی و واریانس آنها برآوردهای سازگار ارائه می‌کند. در

این رویکرد، تابعی از امیدریاضی حاشیه‌ای متغیر وابسته، به صورت تابعی خطی از متغیرهای کمکی مشخص می‌شود. همچنین فرض می‌شود که واریانس، تابعی معلوم از میانگین حاشیه‌ای است. زمانی که متغیر پاسخ نرمال چندمتغیره باشد، این روش معادل با ماکزیمم درست‌نمایی خواهد بود.

همبستگی بین مشاهدات در این رویکرد به عنوان پارامتر مزاحم عمل می‌کند. همچنین این همبستگی با در نظر گرفتن ماتریس‌های همبستگی کاری گوناگون که ساختارهای متفاوتی دارند، مدل‌سازی می‌شود. تعیین درست ساختار همبستگی در بهبود دقت و در نتیجه کارایی ضرایب رگرسیونی مؤثر خواهد بود. اصلی‌ترین مزیت *GEE* نیز برآورد استوار پارامترها است که حتی در صورت عدم تعیین درست ساختار همبستگی برآوردهای سازگار از پارامترهای رگرسیونی را به دست می‌دهد.

اگرچه اغلب مطالعه‌ها برای گردآوری داده‌های کامل بر روی تمام افراد شرکت‌کننده طراحی می‌شوند، داده‌های گم‌شده به طور معمول رخ می‌دهند و باید به شکل مناسب در تحلیل وارد گردند. خیلی اوقات محققان به سادگی از رکوردهای دارای مقادیر گم‌شده چشم‌پوشی کرده و از باقی داده‌ها برای انجام تحلیل استفاده می‌کنند. این روش تحلیل "مورد-کامل" نامیده می‌شود. انجام تحلیل مورد-کامل با نرم افزارهای موجود آسان است اما این تحلیل ناکارآمد خواهد بود زمانی که از متغیرهای پاسخ مشاهده شده متعلق به رکوردهای گم‌شده استفاده نمی‌کند. حال مسئله معادله‌های برآوردساز را در حضور داده‌های گم‌شده بررسی می‌کنیم. روشهای گوناگونی برای مواجهه با این مسئله وجود دارد:

رویکرد اول استفاده از معادله‌های برآوردساز موزون است. در این رویکرد، سهم هر مشاهده در معادله برآوردساز توسط معکوس احتمال مشاهده شدن وزن داده می‌شود.

رویکرد دوم استفاده از تقریب توزیع چندمتغیره‌ی نرمال است. این رویکرد از برآورد نرمال برای پارامترهای همبستگی بهره می‌برد، یعنی تابع برآوردسازی که برآورد پارامترهای همبستگی را نتیجه می‌دهد، از درست‌نمایی نرمال چند متغیره استفاده می‌کند.

رویکرد سوم استفاده از روش‌های جانهای است. در جانهای میانگینی، متغیر گم‌شده با میانگین نمونه‌ای جایگزین می‌شود. در جانهای چندگانه، متغیر گم‌شده با یک نمونه تصادفی از مقادیر برازش شده جایگزین می‌گردد.

۳-۱ تاریخچه‌ی مسأله‌ی مورد بررسی

رویکرد *GEE* اولین بار توسط لیانگ و زیگر (۱۹۸۶) مطرح گردید. بسط و اصلاح‌های گوناگونی از این روش صورت گرفته است. پرنتمیس (۱۹۸۸)، *GEE* را برای داده‌های دودویی همبسته بسط داد. میلر و همکاران (۱۹۹۳) رویکرد *GEE* را برای متغیرهای چندحالتی توسعه دادند. سپس رابینز و همکاران (۱۹۹۴) معادله‌های براوردساز موزون را مطرح کردند. ژائو و همکاران (۱۹۹۶) نشان دادند این معادله‌ها زمانی که داده‌ها گم‌شده هستند کاربردی و قابل اجراست. لایپسیتز و همکاران (۱۹۹۹) از معادله‌های براوردساز موزون برای برخورد با مسئله داده‌های گم‌شده استفاده کردند.

چی و پایک (۱۹۹۷) روش‌های جانهی را برای معادله‌های براوردساز تعمیم یافته در نظر گرفتند زمانی که داده‌ها کاملاً تصادفی گم‌شده هستند. پایک در همان سال جانهی را برای داده‌های غیر از *MCAR* مطرح کرد. چی و پایک در اواخر همان سال مدل معادله‌های براوردساز تعمیم یافته را برای پاسخ‌های دودویی زمانی که متغیرهای کمکی گم‌شده هستند ارائه دادند. لایپسیتز و همکاران (۲۰۰۰) روش‌هایی را براساس تقریب توزیع چندمتغیره‌ی نرمال پیشنهاد دادند. لایپسیتز و همکاران (۲۰۰۹) معادله‌های براوردساز توأم را برای داده‌های گم‌شده مطرح کردند. ژائو و همکاران (۲۰۰۹) آگوریتیم *AU* را پیشنهاد دادند.

۴-۱ اطلاع فیشر

فرض کنید برای پارامتر دلخواه θ ، متغیر تصادفی X دارای تابع چگالی $f(x; \theta)$ باشد به طوری که دامنه متغیر X در تابع چگالی به بردار پارامتر θ بستگی نداشته باشد. با فرض مشتق‌پذیری تابع

$$\text{چگالی و } f(x; \theta) > 0 \text{ برای کمیت } v = \frac{\frac{\partial f(x; \theta)}{\partial x}}{f(x; \theta)} \text{ به سادگی می‌توان نشان داد که}$$

$$E(v) = 0$$

و نیز

$$E(v^2) = -E(v') = \text{Var}(v)$$

که در هر دو مورد امیدریاضی نسبت به تابع چگالی $f(x; \theta)$ با فرض وجود آن، محاسبه می‌شود. هر اندازه $\text{Var}(v)$ بیشتر باشد، میزان تغییر نسبی چگالی X در θ بیشتر می‌باشد. از این رو $\text{Var}(v)$ را مقدار اطلاع درباره پارامتر θ در متغیر تصادفی X می‌نامند و آن را با $I_X(\theta)$ نشان می‌دهند.

این نوع اطلاع که در آمار اهمیت و همچنین کاربرد زیادی دارد را برای اولین بار رونالد فیشر انگلیسی معرفی نمود و از این رو آن را اطلاع فیشر می‌نامند.

در حالت بردار پارامتر $\theta = [\theta_1, \dots, \theta_k]$ ، ماتریس اطلاع فیشر دارای عناصر i_{jk} به شکل زیر است

$$i_{jj} = E \left(\frac{\partial \ln L}{\partial \theta_j} \right)^2$$

$$i_{jk} = E \left(\frac{\partial \ln L}{\partial \theta_j \partial \theta_k} \right)^2$$

۱-۵ خانواده توزیع‌های نمایی

متغیر تصادفی Y که تابع احتمال آن در صورت گسسته بودن، یا تابع چگالی احتمال آن در صورت پیوسته بودن بستگی به پارامتر θ دارد را در نظر بگیرید. توزیع به خانواده نمایی تعلق دارد هرگاه بتوان تابع چگالی آن را به صورت زیر نوشت

$$f_Y(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\} \quad (1-5-1)$$

که در آن θ پارامتر کانونی توزیع، b تابعی دلخواه، $a(\phi)$ تابعی از پارامتر پراکنندگی ϕ و $c(y, \phi)$ تابع ثابت نرمال‌ساز است. ثابت نرمال‌ساز تابعی مستقل از θ است و تضمین می‌کند که مجموع تابع احتمال در حالت گسسته یا انتگرال آن در حالت پیوسته برابر با یک است. در حقیقت خانواده نمایی راهکاری مناسب برای محققین در مدل‌سازی متغیرهای پاسخ پیوسته و گسسته فراهم می‌کند. بسیاری از توزیع‌های معروف از جمله نرمال، پواسون، برنولی، نمایی و گاما به خانواده نمایی تعلق دارند. (هاردین و هیلب ۲۰۰۲)

برآورد ماکزیمم درست‌نمایی پارامتر θ به کمک عبارت‌های زیر حاصل می‌شود

$$E \left(\frac{\partial \ell}{\partial \theta} \right) = 0 \quad (2-5-1)$$

و نیز

$$E \left(\frac{\partial^2 \ell}{\partial \theta^2} + \left(\frac{\partial \ell}{\partial \theta} \right)^2 \right) = 0 \quad (3-5-1)$$

که در آن $U = \frac{\partial \ell}{\partial \theta}$ ، تابع امتیاز فیشر نامیده می‌شود.

به کمک رابطه‌های (۲-۵-۱) و (۳-۵-۱) می‌توان برای هر مشاهده y_i از خانواده نمایی می‌توان عبارات زیر را برای مقادیر میانگین و واریانس به دست آورد

$$E(Y_i) = b'(\theta_i)$$

$$\text{Var}(Y_i) = b''(\theta_i)a(\phi)$$

اعضای معروف خانواده نمایی و عبارتهای $a(\phi), b(\theta)$ و θ مربوط به آن‌ها در جدول زیر آورده شده است. توجه کنید که در بعضی از توزیع‌ها مقدار پارامتر پراکندگی ثابت است.

جدول ۱-۵ عبارات $a(\phi), b(\theta)$ و θ برای توزیع‌های عضو خانواده نمایی

$a(\phi)$	$b(\theta)$	θ	توزیع
σ^2	$-\frac{\theta^2}{2}$	μ	نرمال
1	$-\log(-\theta)$	$-\mu^{-1}$	نمایی
α^{-1}	$-\log(-\theta)$	$-\mu^{-1}$	گاما
1	e^θ	$\log(\mu)$	پواسون
1	$\log(1+e^\theta)$	$\log(\mu/(1-\mu))$	برنولی
1	$n \log(1+e^\theta)$	$\log(\mu/(n-\mu))$	دوجمله‌ای
1	$-k \log(1-e^\theta)$	$\log(\mu/(k+\mu))$	دوجمله‌ای منفی

۱-۶ روش امتیازدهی فیشر

روش امتیازدهی فیشر برای حل عددی معادله‌های حاصل از تابع ماکسیم درست‌نمایی به کار می‌رود. فرض کنید X_1, X_2, \dots, X_n متغیرهای تصادفی مستقل و هم‌توزیع از تابع چگالی دوبار مشتق‌پذیر $f(x; \theta)$ باشند. برای محاسبه برآورد ML بردار پارامتر θ مراحل زیر را طی می‌کنیم. فرض می‌کنیم θ_0 نقطه‌ی شروع این روش باشد. در این حالت بسط سری تیلور تابع امتیاز $v(\theta)$ حول نقطه θ_0 به صورت زیر بیان می‌شود

$$v(\theta) \approx v(\theta_0) - \mathbf{J}(\theta_0)(\theta - \theta_0) \quad (1-6-1)$$

که در آن ماتریس اطلاع مشاهده‌شده در θ_0 است که به صورت زیر تعریف می‌شود

$$\mathbf{J}(\theta_0) = - \sum_{i=1}^n \left[\frac{\partial^2 l}{\partial \theta \partial \theta'} \right]$$

اکنون با قرار دادن $\theta = \theta^*$ و نیز به کمک رابطه (۱-۶-۱) می‌توان عبارت زیر را به دست آورد

$$\theta^* = \theta_0 + \mathbf{J}^{-1}(\theta_0)v(\theta_0)$$

بنابراین روش امتیازدهی فیشر به شکل زیر خواهد بود

$$\theta_{m+1} \approx \theta_m + \mathbf{J}^{-1}(\theta_m)v(\theta_m)$$

تحت شرایط نظم می‌توان نشان داد $\theta_m \xrightarrow{P} \theta^*$. در کاربرد $\mathbf{J}_X(\theta)$ توسط اطلاع فیشر که برابر با $\mathbf{I}_X(\theta) = E(\mathbf{J}(\theta))$ است، جایگزین می‌شود. در این صورت روش امتیازدهی فیشر به شکل زیر

بازنویسی می‌شود

$$\theta_{m+1} \approx \theta_m + I_X^{-1}(\theta_m)v(\theta_m)$$

۷-۱ آگوریتیم نیوتون-رافسون

این آگوریتیم، تابع لگاریتم درست‌نمایی را در همسایگی حدس اولیه انتخاب شده، تقریب می‌زند که این تقریب توسط یک تابع انجام می‌گیرد و محل ماکسیمم این چندجمله‌ای تقریبی، به عنوان حدس دوم برای برآورد ML در نظر گرفته می‌شود. سپس تابع لگاریتم درست‌نمایی در همسایگی نقطه دوم توسط تابع سهمی شکل دیگری تقریب زده می‌شود و حدس سوم به کمک محل ماکسیمم آن تابع به دست می‌آید و این روند ادامه می‌یابد.

برای محاسبه برآورد ضرایب رگرسیونی به روش نیوتون-رافسون مشتق تابع لگاریتم درست‌نمایی نسبت به بردار پارامترهای رگرسیونی β به کمک قانون زنجیره‌ای در مشتق‌گیری به دست می‌آید که یک معادله به صورت زیر است

$$\frac{\partial \ell}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial \ell_i}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j}$$

و وقتی جملات مشتق جداگانه محاسبه شوند، داریم

$$\frac{\partial \ell}{\partial \beta_j} = \sum_{i=1}^n \left(\frac{y_i - b'(\theta_i)}{a(\phi)} \right) \left(\frac{1}{v(\mu_i)} \right) \left(\frac{\partial \mu}{\partial \eta} \right)_i x_{ij} \quad (1-7-1)$$

که یک معادله برآورد برای محاسبه β است که به همراه تقریب خطی بسط سری تیلور زمانی که مشتق لگاریتم تابع درست‌نمایی وجود دارد، استفاده می‌شود.

دیدگاه عمومی در این زمینه برای برآورد پارامترهای رگرسیونی β ، حل معادله برآورد زیر است

$$\ell'(\beta) = 0$$

که به کمک بسط سری تیلور به صورت زیر تقریب زده می‌شود

$$0 = \ell'(\beta^{(0)}) + (\beta - \beta^{(0)})\ell''(\beta^{(0)}) + \frac{1}{2}(\beta - \beta^{(0)})'\ell'''(\beta^{(0)})(\beta - \beta^{(0)}) + \dots$$

$$\text{که در آن } \ell' = \frac{\partial \ell}{\partial \beta} \text{ و } \ell'' = \frac{\partial^2 \ell}{\partial \beta \partial \beta'}$$

برای انجام محاسبات لازم است مقدار $\beta^{(0)}$ به عنوان نقطه شروع آگوریتیم در نظر گرفته شود. در این مرحله به واسطه تقریب خطی تنها دو عبارت اول در محاسبات استفاده می‌شوند و معادله

برآورد بردار پارامترهای رگرسیونی به شکل ساده‌تر زیر در خواهد آمد

$$0 \approx \ell'(\beta^{(0)}) + (\beta - \beta^{(0)})\ell''(\beta^{(0)})$$