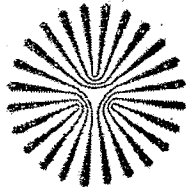


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه پیام نور

دانشکده علوم پایه

گروه آمار

عنوان:

## مدل های رگرسیون ترتیبی و کاربرد آن

نگارش:

مهناز محمدی

استاد راهنما:

دکتر مسعود یار محمدی

استاد مشاور:

دکتر مجتبی گنجعلی

پایان نامه

برای دریافت درجه کارشناسی ارشد

در رشته آمار ریاضی

۱۳۸۸/۷/۱

۱۰۱۵۲۴

کتابخانه و اطلاعیه مرکز علمی پژوهشی  
توسعه و ارتقاء

زمستان ۱۳۸۷

۱۱۶۰۰۲



جمهوری اسلامی ایران  
وزارت علوم، تحقیقات و فناوری



دانشگاه پیام نور  
دانشگاه پیام نور استان تهران

تاریخ

شماره

پیوست

## (( تصویب نامه ))

پایان نامه تحت عنوان :

"مدلهای رگرسیون ترتیبی و کاربرد آن"

تاریخ دفاع ۲۸/۱۰/۸۷ ساعت: ۱۴/۳۰-۱۳/۳۰

نمره: <sup>۱۸/۷۵</sup> هجده و هفتاد و پنج درجه: عالی

امضاء

اعضای هیات داوران مرتبه علمی

۱- استاد راهنما: دکتر یار محمدی

۲- استادمشاور: دکتر گنجعلی

۳- استاد داور داخلی: دکتر پرویز نصیری

۴- استاد داور خارجی: دکتر عادل محمدپور

۴- نماینده محترم گروه: دکتر پرویز نصیری

تهران، خیابان انقلاب،

آیة الله استاد نجات اللهی،

پیش خیابان سپند،

کد پستی ۲۳۳

تلفن: ۸۸۰۱۰۹۰

رنگار: ۸۸۹۰۳۱۵۸

سیستم الکترونیکی:

info@Tehran.pnu.ac

نشانی الکترونیکی:

http://www.Tehran.pnu.ac

تقدیم به:

پدر و مادر عزیزم و همسرفداکارم

با ژرف ترین سپاس ها:

- از لطف بی پایان الهی که هدایت گردون را شوق و انگیزه آموختن و رشد می دهد.

- از قلب پاک مادرم و پدر مهربانم که دعای خیرشان بخش عظیمی از موهبت های الهی است.

- از همسر مهربانم که همواره با صبر و بردباری خویش، صمیمانه پشتیبانم بوده است.

- از اساتید فرزانه ام، دکتر مسعود یار محمدی و دکتر مجتبی گنجعلی که مرا رشد و فراگیری دادند.

- از آقایان دکتر نصیری و دکتر محمدپور که به عنوان داور در جلسه دفاعیه اینجانب حضور یافتند.

- و همه آموزگارانم و دوستانم که مرا در بهتر شدن یاری داده اند.

## چکیده

در مطالعه وابستگی یک متغیر پاسخ به مجموعه ای از متغیرهای مستقل، انتخاب یک مدل به مقیاس اندازه گیری متغیر پاسخ وابسته است. در این پایان نامه، ابتدا با مدل های خطی تعمیم یافته آشنا می شویم. سپس روشهای آماری برای مدل بندی داده های پاسخ ترتیبی، نظیر مدل بخت های متناسب، را شرح می دهیم و به تفسیر و برآورد پارامترهای مدل و آزمون های نیکویی برازش می پردازیم. برای تحلیل اندازه های تکراری زمانی که متغیر پاسخ در هر زمان یک متغیر ترتیبی است، از روش معادلات برآوردگر تعمیم یافته استفاده می کنیم. در پایان با استفاده از نرم افزار SAS این روش را با یک مثال عددی تشریح می کنیم.

**واژه های کلیدی:** پاسخ ترتیبی، مدل بخت های متناسب، نیکویی برازش، اندازه های تکراری، معادلات برآوردگر تعمیم یافته.

## فهرست مطالب

پیشگفتار	۵
فصل اول مدل های خطی تعمیم یافته	۱
۱-۱ مقدمه	۱
۲-۱ مؤلفه های مدل های خطی تعمیم یافته	۲
۳-۱ ساختار رسمی مدل های خطی تعمیم یافته	۳
۴-۱ گشتاورها و معادلات درستنمایی	۵
۱-۴-۱ توابع میانگین و واریانس برای مؤلفه تصادفی	۶
۲-۴-۱ مؤلفه سیستماتیک و تابع ربط	۷
۳-۴-۱ معادلات درستنمایی مدل های خطی تعمیم یافته	۸
۴-۴-۱ ماتریس کواریانس مجانبی برآوردگرهای پارامترهای مدل	۱۰
۵-۱ استنباط مدل های خطی تعمیم یافته	۱۱
۱-۵-۱ کیش و نیکویی برازش	۱۱
۲-۵-۱ مانده ها	۱۳
فصل دوم مدل های رگرسیون با پاسخ ترتیبی	۱۴
۱-۲ مقدمه	۱۴
۲-۲ مدل رگرسیون لوژستیک	۱۵
۳-۲ مدل های رگرسیون ترتیبی	۱۶
۱-۳-۲ مدل های پیوسته گروه بندی شده	۱۷
۲-۳-۲ مدل نسبت دنباله ای	۲۱
۳-۳-۲ مدل بخت های متناسب جزئی	۲۳

۲۶	مدل لوژستیک چند حالتی	۴-۳-۲
۲۷	مدل لوژستیک قالبی	۵-۳-۲
۲۹	فصل سوم برآورد پارامتر و نیکویی برازش	
۲۹	مقدمه	۱-۳
۳۰	برآورد پارامتر	۲-۳
۳۴	آزمون نیکویی برازش	۳-۳
۳۵	روش لیستیز	۱-۳-۳
۳۶	آماره امتیاز	
۴۱	مانده ها	
۴۴	روش پولکستیز	۲-۳-۳
۴۵	آماره های خبی دو پی یرسون و کیش	
۴۶	آماره های خبی دو پی یرسون و کیش تعدیل یافته	
۴۸	مثال عددی	۴-۳
۵۲	فصل چهارم تحلیل داده های ترتیبی و همبسته	
۵۲	مقدمه	۱-۴
۵۳	مدل GEE لوژستیک برای داده های دودویی تکراری	۲-۴
۵۷	مثال کاربردی ۱	۳-۴
۶۲	GEE برای مدل بخت های متناسب	۴-۴
۶۶	برآورد $\alpha$	۱-۴-۴
۶۷	مثال کاربردی ۲	۵-۴
۷۳	بحث و نتیجه گیری	۶-۴
۷۴	پیوست الف. داده ها	
۷۶	پیوست ب. برنامه های رایانه ای	
۸۴	واژه نامه فارسی به انگلیسی	
۸۶	منابع	



## فهرست جداول

جدول ۱-۱	رابطه های متعارف برای مدل های خطی تعمیم یافته	۴
جدول ۱-۲	درجه بندی داروها روی آزمایش داروی مسکن	۲۳
جدول ۲-۲	نتایج برازاندن مدل بخت های متناسب به داده های آزمایش داروی مسکن	۲۴
جدول ۱-۳	رده بندی داده های مقادیر مشاهده شده برای آماره های خنثی دو پی یرسون و کیش	۴۵
جدول ۲-۳	آزمون بخت های متناسب	۴۸
جدول ۳-۳	تحلیل برآوردهای ماکسیمم درستنمایی	۴۹
جدول ۴-۳	برآوردهای نسبت بخت ها	۴۹
جدول ۵-۳	طبقه بندی الگوی متغیر کمکی	۵۰
جدول ۶-۳	فراوانی های مقادیر مشاهده شده و مورد انتظار در هر الگو	۵۰
جدول ۷-۳	مقادیر آماره آزمون	۵۰
جدول ۱-۴	تحلیل GEE برآورد پارامترهای مدل ۱	۵۸
جدول ۲-۴	آماره والد برای تحلیل مدل ۱	۵۸
جدول ۳-۴	برآوردها تحت ساختارهای همبستگی مختلف در مدل ۲	۵۹
جدول ۴-۴	ماتریس همبستگی اعمال شده	۶۱
جدول ۵-۴	معیار ارزیابی نیکویی برازش	۶۱
جدول ۶-۴	برآورد پارامترها در مدل ۱	۶۸
جدول ۷-۴	تحلیل GEE برآورد پارامترهای مدل ۲	۶۹
جدول ۸-۴	آماره والد برای تحلیل GEE	۶۹
جدول ۹-۴	نتایج برآورد مقایسه ای	۷۰
جدول الف-۱	داده های مثال افسردگی روحی (بخش ۳-۴)	۷۴
جدول الف-۲	داده های مثال بیماری پوستی (بخش ۴-۵)	۷۵

## پیشگفتار

تحلیل داده های گسسته جایگاه مهمی را در نوشتارهای آماری به خود اختصاص داده اند. در سالهای اخیر، روش مدل بندی آماری بیشتر از سایر روشها مورد توجه تحلیل گران قرار گرفته است. به کمک مدل بندی آماری می توان تأثیر متغیرهای کمکی مختلف یا عوامل خطر را بر یک یا چند متغیر پاسخ مورد بررسی قرار داد. با توجه به نوع متغیر پاسخ مدل های مختلفی را می توان برای تحلیل داده ها به کار برد.

زمانی که متغیر پاسخ رسته ای است آن متغیر نمی تواند توزیع نرمال داشته باشد، اما ممکن است از خانواده نمایی پیروی کند. مدل های رگرسیون با پاسخ غیر نرمال در خانواده نمایی اولین بار توسط نلدر و ودربرن (۱۹۷۲) مورد بررسی قرار گرفت، که آنها را مدل های خطی تعمیم یافته نامیدند و سپس مک کولا و نلدر (۱۹۸۹) این مدل ها را بیشتر مورد بررسی قرار دادند.

زمانی که پاسخ دو سطحی است، از مدل رگرسیون لوژیستیک استفاده می شود؛ اما مدل بندی داده های پاسخ ترتیبی معمولاً پیچیده تر از پاسخ های دو سطحی است. با مروری بر مقالات کاربردی منتشر شده در چند دهه گذشته، دیده می شود که برای تحلیل داده های ترتیبی و ارتباط آنها با متغیرهای مستقل یا عوامل خطر بیشتر از آزمون های آماری استفاده شده است تا مدل بندی آماری. در برخی موارد نیز این پاسخ های ترتیبی به پاسخ های دیگری، مثلاً دوحالتی، تغییر یافته و سپس به کمک آنها مورد بررسی قرار می گرفتند. به دلیل کارایی پایین این روشها مدل های جایگزین برای مدل بندی داده های ترتیبی مطرح شدند. در ابتدا والکر و دانکن (۱۹۶۷) مدل لوجیت تجمعی را پیشنهاد کردند و سپس مک کولا (۱۹۸۰) به معرفی آن پرداخت، و سپس مدل های دیگر رگرسیون ترتیبی توسط فینبرگ (۱۹۸۰)، انگل (۱۹۸۸) و پیترسن و هرل (۱۹۹۰) معرفی شدند.

در بسیاری از مطالعات متغیر پاسخ برای هر واحد نمونه گیری در زمان های متفاوت مشاهده می شود، که به آنها اندازه های تکراری گفته می شود. برای تحلیل داده های چند متغیره و همبسته از معادلات برآوردگر تعمیم یافته استفاده می شود، که برای اولین بار توسط لیانگ و زیگر (۱۹۸۶) برای

مدل های خطی تعمیم یافته ارائه شدند. این روش به طور مستقیم برای تحلیل پاسخ های ترتیبی همبسته قابل استفاده نبود، اما لیستیز و دیگران (۱۹۹۴) این روش را برای مدل بندی داده های ترتیبی تکراری تعمیم دادند. در این روشها، برای برآورد پارامترهای رگرسیونی از روش شبه درستنمایی به جای ماکسیمم درستنمایی استفاده می شود.

در این رساله هدف معرفی مدل های رگرسیون ترتیبی و کاربرد آنها، انجام آزمون های نیکویی برازش و برآورد پارامتر با استفاده از روش GEE برای مدل ها با پاسخ ترتیبی تکراری است. در فصل اول ابتدا به معرفی مدل های خطی تعمیم یافته و برآورد پارامترهای آن ها پرداخته و سپس برای آزمون پارامترها آماره کیش و مانده ها را معرفی می کنیم. در فصل دوم رگرسیون لوژستیک برای داده های پاسخ دودویی را معرفی کرده و سپس مدل های رگرسیون با پاسخ ترتیبی، نظیر مدل بخت های متناسب، نسبت دنباله ای و غیره را ارائه می کنیم. در فصل سوم در مورد برآورد پارامترها و آزمونهای نیکویی برازش بحث می کنیم، که برای آزمون نیکویی برازش به معرفی روشهای مطرح شده توسط لیستیز و همکاران (۱۹۹۶) و پولکستینز و رابینسون (۲۰۰۴) می پردازیم و سپس با یک مثال یکی از این روشها را به کار می بریم. در فصل چهارم، ابتدا معادلات برآوردگر تعمیم یافته را برای پاسخ های دودویی به همراه یک مثال تشریح و سپس از این روش برای تحلیل پاسخ های ترتیبی تکراری استفاده کرده و با یک مثال کاربردی بحث را به پایان می رسانیم.

## فصل اول

### مدل های خطی تعمیم یافته

#### ۱-۱ مقدمه

زمانی که با مدل های رگرسیون خطی و غیرخطی سروکار داریم توزیع نرمال نقش محوری را ایفا می کند. در حقیقت در روشهای استنباطی مربوط به الگوهای رگرسیون خطی و غیر خطی فرض بر این است که متغیر پاسخ  $Y$  از توزیع نرمال تبعیت می کند. وضعیت های عملی زیادی وجود دارند که این فرض حتی به طور تقریبی برقرار نیست.

مدل های خطی تعمیم یافته ( $GLM$ <sup>۱</sup>) معرفی شده توسط نلدر و ودربرن<sup>۲</sup> (۱۹۷۲) برای پردازش الگوهای رگرسیون به داده های پاسخ یک متغیری ارائه شده اند که متغیر پاسخ عضوی از توزیع بسیار جامعی که خانواده نمایی نامیده می شود، است. یک مدل خطی تعمیم یافته به وسیله سه مؤلفه مشخص می شود که عبارتند از: (۱) مؤلفه تصادفی که متغیر پاسخ  $Y$  و توزیع احتمال آن را نشان می دهد، (۲) مؤلفه سیستماتیک که متغیرهای تبیینی استفاده شده در یک تابع پیشگوی خطی را نشان می دهد، و (۳) تابع ربط که ارتباط  $E(Y)$  به مؤلفه سیستماتیک را توصیف می کند.

---

<sup>۱</sup> Generalized Linear Models

<sup>۲</sup> Nelder and Wedderburn

## ۲-۱ مؤلفه های مدل های خطی تعمیم یافته

مؤلفه تصادفی یک GLM عبارت است از یک متغیر پاسخ  $Y$  با مشاهدات مستقل  $(y_1, y_2, \dots, y_n)$  که از یک توزیع از خانواده نمایی پیروی می کند. این خانواده دارای تابع چگالی احتمال یا تابع جرم به شکل

$$f(y_i; \theta_i) = a(\theta_i) b(y_i) \exp[y_i Q(\theta_i)] \quad (1-1)$$

می باشد. چند توزیع مهم مانند توزیع های پواسون و دو جمله ای جزو موارد خاص و ویژه این خانواده هستند. مقدار پارامتر  $\theta_i$  برای  $i=1, 2, \dots, n$ ، به مقادیر متغیرهای تبیینی وابستگی دارد. جمله  $Q(\theta)$  پارامترکانونی نام دارد. در بخش ۱-۴ فرمول کلی تری که شامل یک پارامتر پراکندگی است ارائه می شود.

مؤلفه سیستماتیک یک GLM یک بردار  $(\eta_1, \dots, \eta_n)$  را به متغیرهای تبیینی در یک مدل خطی وابسته می کند. فرض کنید مقدار پیشگوی  $j$  ( $j=1, 2, \dots, p$ ) برای آزمودنی  $i$  باشد. در این صورت

$$\eta_i = \sum_j \beta_j x_{ij}, \quad i=1, 2, \dots, n$$

این ترکیب خطی از متغیرهای توضیحی، پیشگوی خطی نام دارد. معمولاً  $x_{i1}=1$  برای همه  $i$ ها، برای ضریب عرض از مبدا (اغلب با  $\alpha$  نشان داده می شود) در مدل استفاده می شود.

سومین مؤلفه GLM تابع ربط است، که مؤلفه های تصادفی و سیستماتیک را به هم متصل می کند. فرض کنید  $\mu_i = E(Y_i)$ ،  $i=1, 2, \dots, n$ .

این مدل  $\mu_i$  را توسط  $\eta_i = g(\mu_i)$  به  $\eta_i$  ربط می دهد، که در آن تابع ربط  $g$  یک تابع یکنوای مشتق پذیر می باشد، که  $E(Y_i)$  را به متغیرهای تبیینی از طریق رابطه زیر پیوند می دهد:

$$g(\mu_i) = \alpha + \sum_j \beta_j x_{ij}, \quad i=1, 2, \dots, n \quad (2-1)$$

تابع ربط  $g(\mu) = \mu$ ، که در آن  $\mu_i = \eta_i$ ، ربط همانی نامیده می شود. این یک تابع ربط برای رگرسیون معمولی با  $Y$  ای است که به طور نرمال توزیع شده است. تابع ربطی که میانگین را به

پارامتر کانونی تبدیل می کند، ربط متعارف نام دارد. در ربط متعارف  $g(\mu_i) = Q(\theta_i)$  و  $Q(\theta_i) = \sum_j \beta_j x_{ij}$

به طور خلاصه، یک GLM یک مدل خطی است برای میانگین تبدیل یافته از یک متغیر پاسخ که توزیعی از خانواده نمایی دارد. در اینجا خصیصه هایی که مدل های خطی تعمیم یافته را تعریف می کند، بیان می کنیم.

### ۳-۱ ساختار رسمی مدل های خطی تعمیم یافته

زمینه این توزیع خانواده نمایی است که با توزیع احتمال داده شده در معادله (۱-۱) تعریف می شود. این ساختار به شرح زیر است:

۱. مشاهدات پاسخ مستقل  $y_1, y_2, \dots, y_n$  را به ترتیب با میانگین های  $\mu_1, \mu_2, \dots, \mu_n$  داریم.

۲. مشاهده  $y_i$  دارای توزیعی است که عضوی از خانواده نمایی است.

۳. بخش منظمی از مدل پیشگوهای  $x_1, x_2, \dots, x_p$  را شامل می شود.

۴. مدل حول پیشگوی خطی  $\eta = \beta x = \alpha + \sum_{j=1}^p \beta_j x_j$  ساخته می شود. وجود این پیشگوی خطی اصطلاح الگوهای خطی تعمیم یافته را پیشنهاد می کند.

۵. مدل با استفاده از یک تابع ربط پیدا می شود:

$$\eta_i = g(\mu_i), \quad i = 1, 2, \dots, n \quad (3-1)$$

اصطلاح ربط از این حقیقت ناشی می شود که این رابطه، تابع ربط بین میانگین و پیشگوی خطی است. توجه می کنید که میانگین پاسخ مورد انتظار عبارت است از:

$$E(y_i) = g^{-1}(\eta_i) = g^{-1}(\beta x)$$

درحقیقت، در رگرسیون خطی چند گانه الگوی

$$\mu_i = \eta_i = \beta x_i, \quad i = 1, 2, \dots, n$$

حالت خاصی را پیشنهاد می کند که در آن  $g(\mu_i) = \mu_i$  است. بنابراین، تابع ربطی که مورد استفاده قرار می گیرد یک ربط همانی است.

۶. تابع ربط یک تابع یکنوای مشتق پذیر است.

۷. واریانس  $\sigma_i^2$  تابعی از میانگین  $\mu_i$  است. ( $i = 1, \dots, n$ )

برای جزئیات بیشتر در مورد ساختار GLM می توان به مک کولا و نلدر<sup>۱</sup> (۱۹۸۹) مراجعه کرد. جدول (۱-۱) توابع ربط متعارفی را برای متداول ترین انتخاب های توزیع هایی که با مدل خطی تعمیم یافته انتخاب می شوند، نشان می دهد.

توابع ربط دیگری وجود دارند که از آنها می توان با یک مدل خطی تعمیم یافته استفاده کرد. چند مورد از این توابع عبارتند از:

۱. ربط پروبیت

$$\eta_i = \Phi^{-1}[E(y_i)]$$

که در آن  $\Phi$  تابع توزیع تجمعی نرمال استاندارد را نشان می دهد.

۲. ربط لگ-لگ مکمل

$$\eta_i = \log[-\log(1 - \mu_i)]$$

جدول ۱-۱ ربط های متعارف برای مدل های خطی تعمیم یافته

توزیع	ربط متعارف	
نرمال	$\eta_i = \mu_i$	(ربط همانی)
دوجمله ای	$\eta_i = \log\left(\frac{p_i}{1 - p_i}\right)$	(ربط لوژستیک)
پواسن	$\eta_i = \log(\mu_i)$	(ربط لگاریتمی)
نمایی	$\eta_i = \frac{1}{\mu_i}$	(ربط وارون)

<sup>۱</sup> McCullagh and Nelder

۳. ربط خانواده توانی

$$\eta_i = \begin{cases} \mu_i^\lambda & \lambda \neq 0 \\ \log(\mu_i) & \lambda = 0 \end{cases}$$

از دیدگاه بنیادی مدل خطی تعمیم یافته شامل دو جزء می باشد، که عبارتند از: توزیع پاسخ (یا توزیع خطا) و تابع ربط. انتخاب تابع ربط را به نوعی مشابه انتخاب تبدیلی از پاسخ می توان تلقی کرد، در عین حال درک این مطلب که تابع ربط تبدیل میانگین جامعه است نه داده ها، با اهمیت است. به طور عینی همان طور که استفاده نکردن از یک تابع تبدیل درست می تواند مشکلاتی را در برازش الگوی خطی به بار آورد، انتخاب های نادرست تابع ربط نیز می تواند به مسائل معنی داری در مورد یک مدل خطی تعمیم یافته منجر شود.

#### ۴-۱ گشتاورها و معادلات درستنمایی

در این بخش توجه خود را به جزئیاتی نظیر معادلات درستنمایی و روشهایی برای برازاندن آنها معطوف می کنیم. مؤلفه تصادفی GLM فرض می کند که  $n$  مشاهده  $(y_1, y_2, \dots, y_n)$  روی  $Y$  مستقل و با تابع چگالی احتمال برای  $y_i$  به شکل زیر

$$f(y_i; \theta_i, \phi) = \exp\{[y_i \theta_i - b(\theta_i)] / a(\phi) + c(y_i, \phi)\}. \quad (4-1)$$

هستند، که خانواده نمایی نام دارد و  $\phi$  پارامتر پراکندگی نامیده می شود (جورگنسن<sup>۱</sup> ۱۹۸۷) و پارامتر  $\theta_i$  یک پارامتر مکانی است.

زمانی که  $\phi$  معلوم است، معادله (۴-۱) به معادله ساده (۱-۱) برای خانواده نمایی مکانی تبدیل می شود. که  $Q(\theta)$  و  $a(\theta)$  و  $b(y)$  در رابطه (۱-۱) به ترتیب با عبارتهای  $\theta / a(\phi)$  و  $\exp[-b(\theta)/a(\phi)]$  و  $\exp[c(y, \theta)]$  در (۴-۱) معادل هستند.

فرمول کلی تر (۴-۱) برای خانواده های یک پارامتری نظیر دو جمله ای و پواسون، به کار نمی رود. معمولاً  $a(\phi)$  به شکل  $a(\phi) = \phi / \omega_i$  است که  $\omega_i$  یک وزن معلوم است. برای مثال، زمانی که

<sup>۱</sup> Jørgensen



$y_i$  میانگین  $n_i$  مشاهده مستقل است، نظیر یک نسبت برای  $n_i$  آزمایش برنولی، آنگاه  $\omega_i = n_i$  می باشد.

### ۱-۴-۱ توابع میانگین و واریانس برای مؤلفه تصادفی

عبارتهای کلی برای  $E(Y_i)$  و  $\text{var}(Y_i)$  از رابطه (۴-۱) به دست می آیند. فرض کنید  $L_i = \log f(y_i; \theta_i, \phi)$ ، لگاریتم درستنمایی  $y_i$  باشد، به طوری که تابع لگاریتم درستنمایی برابر است با  $L = \sum_i L_i$ ، در نتیجه:

$$L_i = [y_i \theta_i - b(\theta_i)] / a(\phi) + c(y_i, \phi). \quad (۱-۵)$$

بنابراین:

$$\partial L_i / \partial \theta_i = [y_i - b'(\theta_i)] / a(\phi), \quad \partial^2 L_i / \partial \theta_i^2 = -b''(\theta_i) / a(\phi),$$

که در آن  $b'(\theta_i)$  و  $b''(\theta_i)$  به ترتیب عبارتند از مشتقات مرتبه اول و دوم  $b(\cdot)$  نسبت به  $\theta_i$ . با توجه به اینکه خانواده نمایی در شرایط نظم صدق می کند نتایج درستنمایی کلی عبارتند از:

$$E\left(\frac{\partial L}{\partial \theta}\right) = 0 \quad \text{و} \quad -E\left(\frac{\partial^2 L}{\partial \theta^2}\right) = E\left(\frac{\partial L}{\partial \theta}\right)^2$$

با استفاده از رابطه سمت چپ داریم:

$$E[Y_i - b'(\theta_i)] / a(\phi) = 0$$

یا

$$\mu_i = E(Y_i) = b'(\theta_i). \quad (۶-۱)$$

با استفاده از رابطه سمت راست داریم:

$$b''(\theta_i) / a(\phi) = [E(Y_i - b'(\theta_i)) / a(\phi)]^2 + \text{var}(Y_i) / [a(\phi)]^2$$

در نتیجه

$$\text{var}(Y_i) = b''(\theta_i) / a(\phi). \quad (۷-۱)$$

به طور خلاصه، تابع  $b(\cdot)$  در (۴-۱) گشتاورهای  $Y_i$  را تعیین می کند.

### ۲-۴-۱ مؤلفه سیستماتیک و تابع ربط

فرض کنید  $(x_{i1}, \dots, x_{ip})$  مقادیر متغیرهای توضیحی برای مشاهده  $i$ ام باشند. مؤلفه سیستماتیک یک GLM پارامترهای  $\{\eta_i\}$  را به متغیرهایی نسبت می دهد که از پیشگوی خطی زیر استفاده می کنند:

$$\eta_i = \sum_j \beta_j x_{ij}, \quad i = 1, 2, \dots, n$$

که شکل ماتریسی آن به صورت زیر می باشد:

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$$

که در آن  $\boldsymbol{\eta}_i = (\eta_1, \dots, \eta_n)'$  است و  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$  پارامترهای مدل می باشند، و  $X$  ماتریس  $n \times p$  از مقادیر متغیرهای توضیحی برای  $n$  مشاهده بوده و در مدل های خطی کلی ماتریس طرح نامیده می شود.

GLM توسط یک تابع ربط  $g(\cdot)$ ،  $\eta_i$  را به  $\mu_i = E(Y_i)$  ربط می دهد. بنابراین  $\mu_i$  توسط

$$\eta_i = g(\mu_i) = \sum_j \beta_j x_{ij}, \quad i = 1, 2, \dots, n$$

به متغیرهای تبیینی ربط داده می شود. تابع ربط  $g$  که  $g(\mu_i) = \theta_i$  در (۴-۱) ربط متعارف می باشد. برای آن، رابطه مستقیم

$$\theta_i = \sum_j \beta_j x_{ij}$$

بین پارامترهای مکانی و پیشگوی خطی رخ می دهد.

چون  $\mu_i = b'(\theta_i)$ ، پارامتر مکانی تابعی از میانگین است،  $\theta_i = (b')^{-1}(\mu_i)$ ، که در آن  $(b')^{-1}(\cdot)$  تابع معکوس  $b'$  است. بنابراین، ربط متعارف عکس  $b'$  است. در مورد پواسون به طور مثال داریم:  $b(\theta_i) = \exp(\theta_i)$ ، در نتیجه  $b'(\theta_i) = \exp(\theta_i) = \mu_i$ . بنابراین  $(b')^{-1}(\cdot)$  عکس تابع نمایی است که تابع لگاریتم می باشد (به طور مثال  $\theta_i = \log \mu_i$ ). ربط متعارف ربط لگاریتمی است.

## ۳-۴-۱ معادلات درستنمایی مدل های خطی تعمیم یافته

روش ماکسیمم درستنمایی مبنای نظری برآورد پارامتر در GLM است. برای  $n$  مشاهده مستقل، با استفاده از (۴-۱) تابع لگاریتم درستنمایی عبارت است از:

$$L(\beta) = \sum_i L_i = \sum_i \log f(y_i; \theta_i, \phi) = \sum_i \frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + \sum_i c(y_i, \phi). \quad (۸-۱)$$

عبارت  $L(\beta)$  وابستگی  $\theta$  را به پارامترهای مدل  $\beta$  منعکس می کند.

معادلات درستنمایی عبارتند از:

$$\frac{\partial L(\beta)}{\partial \beta_j} = \sum_i \frac{\partial L_i}{\partial \beta_j} = 0 \quad \forall j$$

با مشتق گیری از لگاریتم درستنمایی (۸-۱)، از قاعده زنجیری زیر استفاده می کنیم:

$$\frac{\partial L_i}{\partial \beta_j} = \frac{\partial L_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j} \quad (۹-۱)$$

چون  $\frac{\partial L_i}{\partial \theta_i} = [y_i - b'(\theta_i)] / a(\phi)$  و با توجه به روابط (۶-۱) و (۷-۱)

$$\mu_i = E(Y_i) = b'(\theta_i) \quad \text{و} \quad \text{var}(Y_i) = b''(\theta_i) / a(\phi)$$

لذا

$$\frac{\partial L_i}{\partial \theta_i} = (y_i - \mu_i) / a(\phi) \quad \text{و} \quad \frac{\partial \mu_i}{\partial \theta_i} = b''(\theta_i) = \text{var}(Y_i) / a(\phi)$$

همچنین، چون  $\eta_i = \sum_j \beta_j x_{ij}$  پس  $\frac{\partial \eta_i}{\partial \beta_j} = x_{ij}$ .

و چون  $\eta_i = g(\mu_i)$  آنگاه  $\frac{\partial \mu_i}{\partial \eta_i}$  به تابع ربط برای مدل بستگی دارد. به طور خلاصه

$$\frac{\partial L_i}{\partial \beta_j} = \frac{y_i - \mu_i}{a(\phi)} \frac{a(\phi)}{\text{var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} x_{ij} = \frac{(y_i - \mu_i) x_{ij}}{\text{var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i} \quad (۱۰-۱)$$

معادلات درستنمایی عبارتند از:

$$\sum_{i=1}^N \frac{(y_i - \mu_i)x_{ij}}{\text{var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_j} = 0, \quad j = 1, \dots, p. \quad (11-1)$$

اگرچه  $\beta$  در این معادلات ظاهر نمی شود، اما چون  $\mu_i = g^{-1}(\sum_j \beta_j x_{ij})$ ،  $\beta$  به طور ضمنی از طریق  $\mu_i$  در معادلات درستنمایی موجود است. توابع ربط مختلف مجموعه معادلات متفاوتی را نتیجه می دهند.

معادلات درستنمایی (۱۱-۱) فقط از طریق  $\mu_i$  و  $\text{var}(Y_i)$  به توزیع  $Y_i$  بستگی دارند. واریانس خودش از طریق یک شکل تابع خاص

$$\text{var}(Y_i) = v(\mu_i)$$

به میانگین بستگی دارد که تابع  $v$  به صورت  $v(\mu_i) = \mu_i$  برای توزیع پواسون؛  $v(\mu_i) = \mu_i(1 - \mu_i)$  برای توزیع برنولی؛ و  $v(\mu_i) = \sigma^2$  برای توزیع نرمال می باشد. زمانی که  $Y_i$  توزیعی از خانواده نمایی داریم، ارتباط بین میانگین و واریانس، نوع توزیع  $Y_i$  را مشخص می کند (جورگنسن ۱۹۸۷). برای مثال اگر  $Y_i$  دارای توزیعی از خانواده نمایی باشد و  $v(\mu_i) = \mu_i$ ، می توان گفت که  $Y_i$  دارای توزیع پواسون است.

معادلات درستنمایی (۱۱-۱) نسبت به  $\beta$  غیر خطی می باشند، به همین دلیل از روشهای تکراری برای پیدا کردن برآوردهای ماکسیم درستنمایی استفاده می شود. اگرستی<sup>۱</sup> (۲۰۰۲) دو روش نیوتن-رافسون و امتیاز بندی فیشر را برای به دست آوردن برآوردهای ماکسیم درستنمایی به کار برده است.

<sup>۱</sup> Agresti