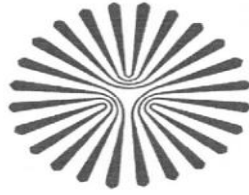


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



**دانشگاه پیام نور**  
**دانشکده علوم پایه**

**پایان نامه**  
**برای دریافت مدرک کارشناسی ارشد**  
**رشته آمار ریاضی**  
**گروه آمار**

**عنوان پایان نامه:**  
**استوار سازی معادلات بر آوردگر تعمیم یافته**  
**سعید معدنی**

**استاد راهنما:**  
**دکتر مسعود یار محمدی**

**استاد مشاور:**  
**دکتر پرویز نصیری**

**شهریور ۹۰**



جمهوری اسلامی ایران  
وزارت علوم، تحقیقات و فناوری

مجمع علوم پایه کشاورزی



دانشگاه پیام نور  
دانشگاه پیام نور استان تهران  
العلم بل لیک العز و العابد والنسر

شماره .....

تاریخ .....

پیوست .....

## صور تجلسه دفاع از پایان نامه کارشناسی ارشد

جلسه دفاع از پایان نامه کارشناسی ارشد آقای سعید معدنی

دانشجوی رشته آمار ریاضی به شماره دانشجویی: ۸۶۷۱۰۵۲۶۶

تحت عنوان:

**"استوار سازی معادلات برآوردگر تعمیم یافته"**

جلسه دفاع با حضور داوران نامبرده ذیل در روز دوشنبه مورخ: ۹۰/۰۶/۰۷ ساعت: ۹-۱۰ در محل

مجمع علوم پایه و کشاورزی برگزار شد. و پس از بررسی پایان نامه مذکور با نمره به عدد ۱۸.۵

به حروف ..... و با درجه ارزشیابی ..... مورد قبول واقع شد.  نشد

ردیف	نام و نام خانوادگی	هیات داوران	مرتبه دانشگاهی	دانشگاه / موسسه	امضاء
۱	دکتر مسعود یاز محمدی	استاد راهنما	رئیس	پیام نور	
۲	دکتر پرویز نصیری	استاد مشاور	د	س	
۳	دکتر صادق رضایی	استاد داور	دیار	امیرکبیر	
۴	دکتر سیما نصیری	نماینده علمی گروه	استاد	پیام نور	
۵	دکتر سیما نصیری	نماینده تحصیلات تکمیلی	استاد	پیام نور	

تهران، خیابان استاد نجات الهی

خیابان شهید فلاح پور، پلاک ۲۷

تلفن: ۸۸۸۰۰۲۵۲

دورنگار: ۸۸۳۱۹۴۷۵

WWW.TPNU.AC.IR

science.agri@tpnu.ac.ir

اینجانب سعید معدنی دانشجوی ورودی سال ۱۳۸۶ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه در پایان نامه خود از فکر، ایده و نوشته دیگری بهره گرفته‌ام با نقل قول مستقیم یا غیر مستقیم منبع و ماخذ آن را نیز در جای مناسب ذکر کرده‌ام. بدیهی است مسئولیت تمامی مطالبی که نقل قول دیگران نباشد بر عهده خویش میدانم و جوابگوی آن خواهم بود.

سعید معدنی

اینجانب سعید معدنی دانشجوی ورودی سال ۱۳۸۶ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه بر اساس مطالب پایان نامه خود اقدام به انتشار مقاله، کتاب، و... نمایم ضمن مطلع نمودن استاد راهنما، بانظر ایشان نسبت به نشر مقاله، کتاب، و... به صورت مشترک و با ذکر نام استاد راهنما مبادرت نمایم.

سعید معدنی

کلیه حقوق مادی مترتب از نتایج مطالعات، آزمایشات و نوآوری ناشی از تحقیق این پایان نامه متعلق به دانشگاه پیام نور می‌باشد

شهریور ۱۳۹۰

به پاس گرمای امید بخش وجودشان

این مجموعه را به همسر و دخترم تقدیم می کنم.

## تقدیر و تشکر

در آغاز بر خود واجب میدانم که از زحمات بی دریغ استاد بزرگوار جناب آقای دکتر مسعود یار محمدی که بارانمایی های راهگشای خود مرا در تهیه و تدوین این اثر همراهی کردند سپاسگزاری کنم.

همچنین از استاد بزرگوار جناب آقای دکتر پرویز نصیری که در مسیر تدوین پایان نامه از مشاوره ارزنده ایشان بهره گرفته ام تشکر می کنم.

## چکیده

معادلات برآوردگر تعمیم یافته ارایه شده توسط لیانگ و زیگر (۱۹۸۶)، روشی برای برآورد پارامترها در مدل‌های خطی تعمیم یافته است، هنگامی که همبستگی نامشخصی در میان مشاهدات موجود باشد. این روش در برابر داده‌های غیر معمول و نقاط دورافتاده تحت تاثیر قرار دارد. روشهای تشخیصی و روشهای استوارسازی روشهایی هستند که به ترتیب جهت شناسایی نقاط دورافتاده و نیز کاهش اثرات این نقاط به کار می‌روند.

در این تحقیق دو روش شناسایی داده‌های دورافتاده و استوارسازی را مورد بررسی قرار داده و آنها را در مورد معادلات برآوردگر تعمیم یافته به کار می‌بریم. روشهای تشخیصی به کاررفته تعمیم روشهایی می‌باشند که توسط کوک (۱۹۷۷) و بلسلی و همکاران (۱۹۸۰) در رابطه با رگرسیون خطی مطرح شده‌اند. روش استوارسازی به کار رفته که در اینجا آنرا معادلات برآوردگر تعمیم یافته استوار می‌نامیم تعمیمی از روش معادلات برآوردگر تعمیم یافته می‌باشد که در برابر داده‌هایی با قدرت نفوذ بالا استوار است. همچنین با ارائه یک مثال کاربردی روشهای تشخیصی و روش‌های استوارسازی را بررسی کرده و با استفاده از روش شبیه‌سازی روشهای استوار را با روش معادلات برآوردگر تعمیم یافته مقایسه خواهیم کرد.

**واژه‌های کلیدی:** معادلات برآوردگر تعمیم یافته - معادلات برآوردگر تعمیم یافته استوار

## فهرست مطالب

مقدمه	۱
<b>فصل اول: معادلات برآوردگر تعمیم یافته</b>	<b>۵</b>
۱-۱ مقدمه	۵
۲-۱ مدل‌های خطی تعمیم یافته	۶
۱-۲-۱ توزیع متغیر پاسخ	۷
۲-۲-۱ پیشگوی خطی	۸
۳-۲-۱ تابع پیوند (تابع ربط)	۸
۳-۱ معادلات درستمایی برای الگوهای خطی تعمیم یافته	۹
۴-۱ روش حداقل مربعات تعمیم یافته و موزون	۱۰
۵-۱ روش شبه درستمایی	۱۳
۶-۱ معادلات برآوردگر تعمیم یافته	۱۳
۱-۶-۱ ساختار داده‌ها و نمادگذاری	۱۴
۲-۶-۱ تشریح اجزاء معادلات برآوردگر	۱۵
۳-۶-۱ ساختارهای مختلف ماتریس همبستگی کاری	۱۷
۴-۶-۱ برآورد پارامتر پراکندگی $\phi$	۲۰
۵-۶-۱ روش حل معادلات برآوردگر تعمیم یافته	۲۱
۶-۶-۱ مراحل برآورد $\beta$ در روش معادلات برآوردگر تعمیم یافته	۲۲

## فصل دوم: بررسی تأثیر و میزان نفوذ داده‌های دورافتاده روش‌های «تشخیصی» و روش‌های

<b>«استوارسازی»</b>	<b>۲۵</b>
۱-۲ مقدمه	۲۵
۲-۲ مدل‌های خطی یک متغیره	۲۶
۱-۲-۲ حذف مجموعه‌ای از مشاهدات	۳۰
۲-۲-۲ رگرسیون استوار	۳۰
۳-۲ مدل‌های خطی تعمیم یافته	۳۵
۱-۳-۲ رگرسیون استوار در مدل‌های خطی تعمیم یافته	۳۹
۴-۲ مدل با پاسخ همبسته	۴۳
۱-۴-۲ رگرسیون استوار در مدل‌های همبسته	۴۴

## فصل سوم: روش‌های «تشخیصی» برای شناسایی نقاط دورافتاده و نقاط بانفوذ در معادلات برآوردگر

<b>تعمیم یافته</b>	<b>۴۹</b>
۱-۳ مقدمه	۴۹
۲-۳ برآورد پارامترها در معادلات برآوردگر تعمیم یافته، به روش کمترین مربعات دوباره وزنی شده تکراری	۵۰
۳-۳ اندازه‌گیری میزان نفوذ بر مبنای روش حذفی	۵۳
۱-۳-۳ ارزیابی میزان نفوذ در حالت حذفی روی $\hat{\beta}$	۵۴
۲-۳-۳ تعیین میزان نفوذ در حالت حذفی روی مقادیر برازش داده شده	۵۹



۳-۴ یک مثال کاربردی..... ۶۱

۳-۴-۱ تعیین مشخصات مدل ..... ۶۲

۳-۴-۲ برنامه مورد استفاده در محیط SAS و نتایج آن..... ۶۳

۳-۴-۳ استفاده از روشهای تشخیص برای یافتن نقاط دورافتاده و نقاط پرنفوذ..... ۶۶

## ۷۷ فصل چهارم: معادلات برآوردگر تعمیم یافته‌ی استوار.....

۴-۱ مقدمه ..... ۷۷

۴-۲ معادلات برآوردگر تعمیم یافته‌ی استوار..... ۷۸

۴-۲-۱ برآورد پارامترهای رگرسیون و برآورد واریانس..... ۸۰

۴-۲-۲ طبقه بندی برآوردها..... ۸۰

۴-۲-۳ برآورد پارامترهای مزاحم..... ۸۲

۴-۲-۴ روش مالوس در مقایسه با روش شوئیپ..... ۸۳

۴-۲-۵ معادلات برآوردگر تعمیم یافته ی استوار برای پاسخهای دودویی همبسته..... ۸۴

۴-۳ مقایسه معادلات برآوردگر تعمیم یافته با معادلات برآوردگر تعمیم یافته ی استوار به روش شبیه سازی..... ۸۵

۴-۳-۱ مقایسه مقدار اریبی روشهای معادلات برآوردگر تعمیم یافته استوار با روش GEE در برآورد  $\beta_1$ ..... ۸۸

۴-۳-۲ مقایسه میانگین مربعات خطا در برآورد  $\beta_1$  در معادلات برآوردگر تعمیم یافته و روشهای استوار آن..... ۹۰

۴-۴ مرور مجدد داده‌های پزشکی..... ۹۳

۴-۴-۱ معادلات برآوردگر تعمیم یافته ی استوار در کلاس «کاهش وزن مشاهده ی مالوس و شوئیپ»..... ۹۴

۴-۴-۲ مقایسه نتایج بدست آمده در روش های استوار مالوس و شوئیپ با روش GEE..... ۹۷

## ۹۸ نتیجه گیری.....

معادلات برآوردگر تعمیم یافته<sup>۱</sup> (GEE) که توسط لیانگ<sup>۲</sup> و زیگر<sup>۳</sup> (۱۹۸۶) ارائه شده است امروزه به طور وسیعی در بیولوژی و پزشکی به کار رفته و کاربرد آن در رشته‌های علمی دیگر مانند کشاورزی و صنایعی که در تحلیل داده‌های آنها همبستگی بین مشاهدات وجود دارد دیده می‌شود. به عنوان مثال در کشاورزی ممکن است داده‌های جمع‌آوری شده روی یک کرت از زمین دارای همبستگی باشند و یا در تحقیقات صنعتی مشاهدات جمع‌آوری شده مربوط به یک کوره مجموعه‌ای از داده‌های همبسته را بوجود آورند.

هرگاه از هر عضو نمونه در چند موقعیت مکانی یا زمانی مشاهداتی جمع‌آوری شود، این داده‌ها را اندازه‌گیری‌های مکرر<sup>۴</sup> می‌نامیم. داده‌های مربوط به هر عضو مجموعه‌ای از مشاهدات همبسته را تشکیل می‌دهند که آنها را خوشه نامیده و داده‌های مکرر را که در آن موقعیت تکرار مشاهدات، نقاط زمانی هستند، داده‌های طولی و مطالعات از این نوع را مطالعات طولی می‌نامند. در مطالعات طولی اغلب بین متغیرهای درون هر خوشه همبستگی وجود داشته و به همین دلیل برای تحلیل این نوع از داده‌ها باید تمهیدات خاصی را به کار برد که روش معادلات برآوردگر تعمیم یافته یکی از آنها است.

به علت کاربرد زیاد روش (GEE) در برآورد پارامترهای مدل‌های خطی تعمیم یافته بررسی تأثیر مشاهدات جمع‌آوری شده روی برآورد پارامترها مهم بوده و اینکه مشاهدات دور افتاده کدامند و تأثیر آنها روی برآورد پارامترها چیست حائز اهمیت است.

---

<sup>۱</sup> .generalized estimating equations

<sup>۲</sup> .liang

<sup>۳</sup> .zeger

<sup>۴</sup> Repeated measurements

به منظور بررسی و تشخیص اینکه در اثر حذف مشاهدات یا خوشه ها چه تغییراتی در برازش مدل اتفاق می افتد، روش های تشخیصی رگرسیون را بکاربرده ایم. استفاده از روش های تشخیصی برای اثرات داده ها، خصوصاً در مدل های چند متغیره کارایی بالایی نداشته بعلاوه شناسایی نقاط دور افتاده و تشخیص آن در رگرسیون دودویی که در آنها متغیر پاسخ به صورت دو حالتی ۰ یا ۱ مطرح می شود از تشخیص هویت آنها در داده های پیوسته دشوارتر است.

از آنجا که هدف از روش های استوار سازی شناسایی و کاهش اثر داده های دورافتاده و پرنفوذ است، در این تحقیق ابتدا روش های شناسایی نقاط دورافتاده و پرنفوذ را مطرح کرده و اثر داده های غیرمتجانس را روی پارامترهای مدل های خطی تعمیم یافته که به روش GEE تحلیل می شوند مورد بررسی قرار می دهیم. سپس روش های استوار سازی را به منظور کاهش اثر نقاط دورافتاده و پرنفوذ مورد توجه قرار داده و آنرا معادلات برآوردگر تعمیم یافته استوار و به اختصار<sup>۱</sup> ( REGEE ) می نامیم. این روشها بر اساس وزن دهی به مشاهدات برمبنای میزان تأثیر آنها روی برآورد پارامترهای مدل طراحی شده اند. روش های برآورد پارامتر در ( REGEE ) توسعه یافته ی روش های استوار سازی مدل های خطی تعمیم یافته اند که توسط پرگیبون<sup>۲</sup> (۱۹۸۲) و کارول و پدرسن<sup>۳</sup> (۱۹۹۳) در استوار سازی رگرسیون لوژستیک و پرایسرو کاکیش<sup>۴</sup> (۱۹۹۵) در استوار سازی مدل های با پاسخ همبسته مطرح شده اند ، در این روش ها کاهش وزن مشاهدات با استفاده از الگوریتم معینی انجام شده و برای هر مشاهده وزنی بین ۰ و ۱ در نظر گرفته می شود که این وزن متناسب با اثر آن مشاهده روی برآورد پارامترهای مدل است.

<sup>۱</sup> Resistant generalized estimating equations

<sup>۲</sup> Pergibon

<sup>۳</sup> Carroll and Pederson

<sup>۴</sup> Preisser and Qaqish

در فصل اول راجع به مدل‌های خطی تعمیم یافته و روش معادلات برآوردگر تعمیم یافته بحث می‌شود. روش‌های تشخیصی و اصول استوارسازی در فصل دوم مورد بررسی قرار می‌گیرد. در فصل سوم روش‌های تشخیصی را برای معادلات برآوردگر تعمیم یافته مورد بحث قرار داده و مثالی کاربردی از داده‌های پزشکی را تحلیل خواهیم کرد، روش‌های استوارسازی را برای معادلات برآوردگر تعمیم یافته در فصل چهارم مورد بررسی قرار داده و با استفاده از شبیه‌سازی و یک مثال کاربردی روش‌های استوار معادلات برآوردگر تعمیم یافته را با روش معمول آن مقایسه خواهیم کرد.

# فصل اول

## معادلات بر آوردگر تعمیم یافته

## فصل اول: معادلات برآوردگر تعمیم یافته

۱-۱ مقدمه

با توجه به کاربرد فراوان الگوهای خطی تعمیم یافته وضعیت های بسیاری وجود دارد که پاسخ های تکراری از یک واحد، مشاهداتی همبسته را به وجود می آورد. مثالی از این مورد را می توان در مطالعات طولی که در آن پاسخها روی یک آزمودنی در طول زمان به تکرار اندازه گیری می شوند در نظر گرفت. بدیهی است که مشاهدات روی یک آزمودنی همبسته اند و در نظر گرفتن همبستگی میان مشاهدات مشکلاتی را بوجود می آورد. در واقع وجود همبستگی استفاده از روشهای حداکثر درستنمایی معمول را مشکل کرده و لازم است از روش های شبه درستنمایی و یا تعمیمی از آن که روش معادلات برآوردگر تعمیم یافته است استفاده شود. در این فصل ابتدا مدل های خطی تعمیم یافته ، روش شبه درستنمایی و در پایان روش معادلات برآورد گر تعمیم یافته را معرفی می کنیم .

## ۱-۲ مدل‌های خطی تعمیم یافته

در الگوهای رگرسیون خطی و غیر خطی توزیع نرمال نقش محوری داشته و در برخی از روش‌های استنباطی مربوط به رگرسیون خطی و غیر خطی فرض آن است که متغیر پاسخ  $Y$  از توزیع نرمال پیروی می‌کند. اما مثال‌های کاربردی زیادی وجود دارد که این فرض حتی به طور تقریبی هم برقرار نیست. مثلاً زمانی که متغیر پاسخ یک متغیر گسسته و شمارشی است این فرض نقض می‌گردد. به عنوان مثال در مدل‌هایی با پاسخ‌های دوتایی نیز که فقط دو حالت موفقیت (۱) و شکست (۰) مورد نظر است فرض نرمال بودن نقض شده است. مواردی وجود دارد که متغیر پاسخ پیوسته بوده ولی فرض نرمال بودن برقرار نیست به عنوان مثال می‌توان به توزیع فشار در اجزاء مکانیکی و یا زمان زوال اجزای الکترونیکی اشاره کرد.

الگوهای خطی تعمیم یافته<sup>۱</sup> (GLM) برای برازش الگوهای رگرسیون به داده‌هایی هستند که متغیر پاسخ در آنها از توزیع خانواده نمایی تبعیت می‌کنند. خانواده توزیع‌های نرمال، دو جمله‌ای، پواسون، هندسی، دو جمله‌ای منفی، نمایی، گاما، از این نوع هستند. شکل کلی مدل خطی تعمیم یافته به صورت زیر است:

$$g(\mu) = g[E(y)] = X'\beta \quad (1-1)$$

$X$  بردار متغیرهای کمکی و  $\beta$  بردار پارامترها یا ضرایب رگرسیون است.  $g$  تابع پیوند یا ربط می‌باشد می‌توان گفت هر الگوی خطی تعمیم یافته سه جزء اصلی دارد:

۱. توزیع متغیر پاسخ.

۲. پیشگوی خطی که متغیرهای برگشت یا کمکی را شامل می‌شود.

---

<sup>۱</sup>. Generalized Linear model.

۳. تابع پیوند  $g$  که پیشگوی خطی را به میانگین متغیر پاسخ مربوط می کند.

### ۱-۲-۱ توزیع متغیر پاسخ

همانطور که گفته شد متغیر پاسخ  $y$  با میانگین  $E(y) = \mu$  در مدل‌های خطی تعمیم یافته دارای

تابع چگالی از خانواده نمایی می باشد. این تابع چگالی به صورت کلی زیر مطرح می شود.

$$f(y, \theta, \phi) = \exp\left\{\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)\right\} \quad (۲-۱)$$

اگر تابع لگاریتم درست‌نمایی را  $L(\theta)$  بنامیم ، داریم :

$$L(\theta) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + c(y_i, \phi) \right\}$$

و از اصول استنباط آماری خواهیم داشت :

$$E\left(\frac{\partial L(\theta)}{\partial \theta}\right) = 0 \quad \Rightarrow \quad \mu = E(y) = \frac{db(\theta)}{d\theta} = b'(\theta)$$

$$\text{var}(y) = \frac{d^2 b(\theta)}{d\theta^2} a(\phi) = b''(\theta) a(\phi) = \frac{d\mu}{d\theta} a(\phi)$$

خانواده تابع چگالی نمایی روش یکسانی را برای برآورد پارامترها در مدل‌های خطی تعمیم یافته

فراهم می کند . مثلاً توزیع نرمال که برجسته ترین عضو خانواده نمایی است دارای تابع چگالی به

صورت زیر است :

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

$$= \exp\left\{\left[\frac{y\mu - \frac{\mu^2}{2}}{\sigma^2}\right] - \frac{1}{2}\left[\frac{y^2}{\sigma^2} + L_n(2\pi\sigma^2)\right]\right\}$$

با توجه به صورت کلی توزیع نمایی در (۲-۱) داریم:

$$\theta = \mu \quad b(\theta) = \frac{\mu^2}{2} \quad a(\phi) = \phi = \sigma^2$$



$$c(y, \phi) = -\frac{1}{\phi} \left[ \frac{y^2}{\sigma^2} + L_n(\sqrt{\pi} \sigma^2) \right]$$

و یا در مورد توزیع دو جمله ای داریم:

$$f(y, n, p) = \binom{n}{y} p^y (1-p)^{n-y}$$

$$= \exp \left\{ L_n \binom{n}{y} + y L_n p + (n-y) L_n (1-p) \right\}$$

$$\theta = L_n \left( \frac{p}{1-p} \right) \quad b(\theta) = n L_n (1 + e^\theta) \quad \phi = 1 \quad c(\cdot) = L_n \binom{n}{y} \quad a(\theta) = 1$$

### ۲-۲-۱ پیشگوی خطی

پیشگوی خطی یا همان مؤلفه سیستماتیک متغیرهایی مانند  $X$  را نشان می دهد که نقش متغیرهای تبیینی را در رگرسیون به عهده دارند. فرض کنیم  $\mu = E(y)$  میانگین احتمال باشد مقدار  $\mu$  در GLM با تغییر متغیرهای تبیینی تغییر می کند. به طور کلی مؤلفه سیستماتیک به صورت زیر مطرح می شود.

$$X'\beta = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n \quad (3-1)$$

این ترکیب خطی، متغیرهای پیشگو نامیده می شود.

### ۳-۲-۱ تابع پیوند (تابع ربط)

یک تابع پیوند چگونگی وابستگی  $E(y) = \mu$  را به  $X'\beta$  شرح می دهد و معمولاً به صورت

$g(\mu)$  نوشته می شود که یک تابع غیر خطی بوده و به صورت زیر مطرح می شود.

$$g(\mu) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n = X'\beta \quad (4-1)$$

ساده ترین تابع پیوند دارای شکل  $g(\mu) = \mu$  است که پیوند همانی نامیده می شود. این تابع میانگین را به طور مستقیم مدل سازی کرده و مدلی شبیه مدل رگرسیون معمولی پدید می آورد که معمولاً برای پاسخ های پیوسته مورد استفاده قرار می گیرد .

$$\mu = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n = X'\beta \quad (5-1)$$

از جمله توابع پیوند دیگر  $g(\mu) = \log(\mu)$  است که مدل را بر حسب لگاریتم میانگین بیان می کند. انتخاب توابع پیوند بستگی به ماهیت متغیر پیشامد کنترل دارد. انواع توابع ربط (پیوند) برای توزیع های مختلف در جدول (۱-۱) آمده است .

جدول ۱-۱

توزیع	پیوند متعارف
نرمال	( پیوند همانی ) $\eta_i = \mu_i$
دو جمله ای	( پیوند لجوجیت ) $\eta_i = L_n \left( \frac{p_i}{1-p_i} \right)$
پواسون	( پیوند لگاریتمی ) $\eta_i = L_n(\mu_i)$
نمایی	( پیوند وارون ) $\eta_i = \left( \frac{1}{\mu_i} \right)$
گاما	( پیوند وارون ) $\eta_i = \left( \frac{1}{\mu_i} \right)$

### ۳-۱ معادلات درستنمایی برای الگوهای خطی تعمیم یافته

روش برآورد پارامترها در GLM، روش حداکثر درستنمایی می باشد ولی تغییرات واقعی حداکثر درستنمایی به الگوریتمی می رسد که برپایه روش تکراری کمترین مربعات وزنی قرار دارد . در

صورتی که از پیوند متعارف  $\eta_i = g[\mu_i] = X_i' \beta$  استفاده کنیم، تابع لگاریتم درستنمایی عبارت است از:

$$L = \log L(y, \beta) = \sum_{i=1}^n \left\{ \frac{(y_i \theta_i - b(\theta_i))}{a(\phi)} + c(y_i, \phi) \right\}$$

$$\frac{\partial L}{\partial \beta} = \frac{\partial L}{\partial \theta_i} \times \frac{\partial \theta_i}{\partial \beta} = \sum_{i=1}^n \frac{1}{a(\phi)} \left[ y_i - \frac{db(\theta_i)}{d\theta_i} \right] X_i = \sum_{i=1}^n \frac{1}{a(\phi)} (y_i - \mu_i) X_i$$

$$= \sum_{i=1}^n \frac{1}{a(\phi)} (y_i - \mu_i) X_i \quad (6-1)$$

می توانیم برآوردهای حداکثر درستنمایی پارامترها را با حل دستگاه معادلات زیر برای  $\beta$  پیدا کنیم از آنجا که در بیشتر موارد  $a(\phi)$  یک عدد ثابت است، این معادلات به صورت زیر در می آیند:

$$\sum_{i=1}^n (y_i - \mu_i) X_i = 0 \quad (7-1)$$

این دستگاه دارای  $p = k + 1$  معادله بوده و هر معادله برای یک پارامتر الگو می باشد. شکل ماتریسی این معادله به صورت زیر است :

$$X'(y - \mu) = 0 \quad (8-1)$$

که در آن  $\mu' = [\mu_1, \mu_2, \dots, \mu_n]$  می باشد. این معادلات را معادلات امتیاز حداکثر درستنمایی می نامند.

#### ۴-۱ روش حداقل مربعات تعمیم یافته و موزون

در مدل رگرسیون خطی  $y = X\beta + \varepsilon$  مفروضاتی مانند  $E(\varepsilon) = 0$  و  $V(\varepsilon) = \sigma^2 I$  در نظر گرفته می شود گاهی این مفروضات قابل اجرا نیستند. به عنوان مثال اگر  $V(\varepsilon) = \sigma^2 V$  که  $V$  یک ماتریس  $n \times n$  شناخته شده باشد اصلاحاتی باید روی روش حداقل مربعات معمولی در نظر گرفته

شود. می دانیم که اگر  $V$  یک ماتریس قطری باشد اما اعضای قطر مساوی نباشد در اینصورت متغیرهای  $y$  ناهمبسته اند اما واریانسهای مساوی ندارند. اگر بعضی از اعضای غیر قطر اصلی صفر نباشد مشاهدات همبسته بوده و برآورد حداقل مربعات معمولی  $\hat{\beta} = (X'X)^{-1}X'y$  بهینه نیست. چون  $\sigma^2 V$  ماتریس کوواریانس خطاها می باشد،  $V$  بایستی ناویژه و معین مثبت باشد بنابراین ماتریس  $n \times n$  ناویژه متقارن  $k$  موجود است به طوری که  $k'k = kk = V$ . ماتریس  $k$  ریشه دوم  $V$  نامیده می شود. متغیرهای زیر را در نظر می گیریم:

$$Z = k^{-1}y \quad B = k^{-1}X \quad g = k^{-1}\varepsilon$$

به عبارت دیگر مدل رگرسیون  $y = X\beta + \varepsilon$  به صورت زیر خواهد شد

$$k^{-1}y = k^{-1}X\beta + k^{-1}\varepsilon$$

$$Z = B\beta + g \quad (9-1)$$

خطاهای تبدیل یافته در مدل جدید دارای امید ریاضی صفر بوده و ماتریس کوواریانس  $g$  هم به صورت زیر خواهد بود.

$$\begin{aligned} V(g) &= E\left\{[g - E(g)][g - E(g)]'\right\} = E(gg') \\ &= E(k^{-1}\varepsilon\varepsilon'k^{-1}) = k^{-1}E(\varepsilon\varepsilon')k^{-1} = \sigma^2 k^{-1}V k^{-1} = \sigma^2 k^{-1}kkk^{-1} = \sigma^2 I \end{aligned}$$

بنابراین اعضای  $g$  دارای میانگین صفر و واریانس ثابت بوده و ناهمبسته می باشند. چون خطاهای  $g$  در مدل  $Z = B\beta + g$  و در مفروضات معمول صدق می کند می توانیم روش حداقل مربعات معمول را برای آن به کار ببریم. تابع حداقل مربعات و در نتیجه معادلات نرمال به صورت زیر خواهد بود.

$$S(\beta) = g'g = \varepsilon'V^{-1}\varepsilon = (y - X\beta)'V^{-1}(y - X\beta) \quad (10-1)$$