

# تحلیل یک مسئله با استفاده از دیدگاه‌های متفاوت

مجتبی گنجعلی<sup>†,\*</sup> و لیلا لطیفیان<sup>‡</sup>

<sup>†</sup> دانشگاه شهید بهشتی

<sup>‡</sup> بانک مرکزی جمهوری اسلامی ایران

چکیده. در این مقاله یک مسئله‌ی کاربردی که در آن، پاسخ مورد نظر، تعداد موفقیت‌ها در یک آزمایش به‌خصوص است مطرح می‌شود و برای مدل‌بندی از دیدگاه‌های متفاوت مورد بررسی قرار می‌گیرد. بررسی تأثیر مقدار پاسخ دورافتاده بر نتایج یک تحلیل رگرسیونی از اهمیت فراوان برخوردار است. به این دلیل با استفاده از روش‌های تشخیص، داده‌های دورافتاده شناسایی می‌شوند. نشان داده شده است که استفاده از روش تبدیل آرک‌سینوس ممکن است تشخیص دورافتاده بودن برخی داده‌ها را منحرف کند. در صورتی که حذف داده‌های دورافتاده در مدل‌بندی امکان‌پذیر نباشد، استفاده از روش‌های استوار پیشنهاد شده است. روش ماکسیمم درست‌نمایی که از آن در مدل‌های خطی تعمیم‌یافته و روش تبدیل، برای برآورد پارامترها استفاده می‌شود، و روش شبه‌درست‌نمایی، استوار نیستند. روش دیگری که به نتایج استوار منجر می‌شود و به روش شبه‌درست‌نمایی استوار موسوم است، در این مقاله بازبینی و روشی که از نظر محاسباتی برای استفاده آسان‌تر است، برای محاسبه‌ی مقادیر معناداری در آن ارائه شده است. شیوه‌های مختلف مدل‌بندی نیز در مثال کاربردی مقایسه شده‌اند. © ۱۳۸۳ پژوهشکده‌ی آمار. همه‌ی حقوق محفوظ است.

واژگان کلیدی. روش تبدیل؛ رگرسیون لوژیستیک؛ مدل‌های خطی تعمیم‌یافته؛ مانده‌ی پی‌یرسونی؛ شبه‌درست‌نمایی استوار.

## ۱ مقدمه

تحلیل مسائل کاربردی آمار به تأمل و اندیشه‌ی عمیق نیاز دارد. گاهی در دید نخست، تحلیل یک مسئله‌ی آماری، ساده به نظر می‌رسد و روش‌های موجود یا بسیار سنتی، به‌نظر به رهیافتی اساسی می‌انجامند. اما

\* نویسنده‌ی عهده‌دار مکاتبات.

پس از تحلیل دقیق و در نظر گرفتن روش‌های مختلف، از جمله روش‌های سنتی و مدرن است که می‌توان در مورد داده‌ها و آنچه آن‌ها به ما می‌آموزند اظهار نظر کرد و گاهی به اندیشه‌های جدید دست یافت. هر چه در تحلیل توصیفی یک مجموعه داده‌ی آماری بیش‌تر بیندیشیم و سعی کنیم تا این مجموعه داده‌ها را با روش‌های آماری متعدد تحلیل کنیم، نه تنها استنباط بهتری می‌توانیم به دست آوریم، بلکه می‌توانیم اطلاع کسب کنیم که در چه موقعیت‌هایی چه روشی سودمندتر است؛ هر چند ممکن است هیچ روشی کامل نباشد. به همین دلیل است که شیوه‌های متعدد آماری در تحلیل‌های آماری رگرسیونی مورد استفاده قرار می‌گیرند.

به عنوان یک نمونه‌ی بارز کاربردی، حالتی را در نظر بگیرید که پاسخ مورد نظر، تعداد موفقیت‌ها در یک آزمایش به خصوص است و در صددیم تا تأثیر برخی متغیرهای تبیینی بر میانگین موفقیت‌ها را بررسی کنیم. یکی از روش‌های بسیار سنتی برای تحلیل چنین داده‌هایی، استفاده از تبدیلی از نسبت موفقیت‌ها در نمونه، به گونه‌ای است که متغیر تبدیل‌یافته متغیری پیوسته و دارای توزیع نرمال با واریانس تثبیت‌شده باشد [۱۰]، [۲]. یکی دیگر از روش‌های تحلیل، استفاده از مدل‌های خطی تعمیم‌یافته است که از سوی نلدر و ودربرن [۱۵] معرفی شده است. البته برای پاسخ دوجمله‌ای، مدل رگرسیونی لوژیستیک که حالت خاصی از مدل‌های خطی تعمیم‌یافته است، قبلاً توسط فینی [۸]، [۹] و برکسون [۳] بیان شده است. به هر حال، براوردگرهای ماکسیمم درست‌نمایی در مدل رگرسیونی لوژیستیک نسبت به داده‌های دورافتاده نالاستوارند و اخیراً رده‌ی براوردگرهای استوار، مورد تأکید قرار گرفته است. رهیافت شبه‌درست‌نمایی نیز توسط ودربرن [۱۷] معرفی شده است، که از معادلات براوردگر معرفی‌شده در آن نیز می‌توان برای براورد پارامترها استفاده کرد. به هر حال، براوردهای به دست آمده از این روش نیز نسبت به داده‌های دورافتاده نالاستوارند.

باکس [۴] واژه‌ی «استواری» را برای نخستین بار ارائه کرد، هوبر [۱۲] نظریه‌ی استواری را مورد بررسی قرار داد، همپل [۱۱] آمار استوار را بر اساس تابع تأثیر بیان کرد و کانتونی و رونجی [۵] استنباط استوار را برای مدل‌های خطی تعمیم‌یافته، بر اساس کبیض (deviance) استوار که تعمیم‌یافته‌ی شبه‌درست‌نمایی است، به دست آوردند. در این مقاله از این روش، روش تبدیل آرک‌سینوس و رگرسیون لوژیستیک برای تحلیل داده‌هایی با پاسخ دوجمله‌ای استفاده می‌کنیم.

برای ایجاد انگیزه، در بخش بعد، مثال کاربردی را تشریح می‌کنیم. در بخش ۳ روش‌های مختلف تحلیل را می‌آوریم. در بخش ۴ نتایج به دست آمده از روش‌های مختلف را مقایسه و در بخش ۵ نتیجه‌گیری می‌کنیم.

## ۲ مسئله‌ی کار بردی

در این بخش، داده‌هایی که از فلیس [۱۶] گرفته شده است، مورد ارزیابی قرار می‌گیرد. این داده‌ها به وسیله‌ی ویلیامز [۱۸] نیز مورد بررسی قرار گرفته است. مکولا و نلدر [۱۴] این داده‌ها را به کار گرفته‌اند تا روش‌هایی برای ارزیابی انحرافات زیاد از مدل، که به دلیل حضور مشاهده‌ی دورافتاده ایجاد می‌شود، مورد بررسی قرار گیرد. این داده‌ها نسبت گوجه‌فرنگی‌های آفت‌زده را در آزمایشی با سه بلوک و هشت مقدار مصرف مختلف حشره‌کش نشان می‌دهد. این داده‌ها در جدول ۱ نشان داده شده است. در این جدول، نسبت گوجه‌فرنگی‌های آفت‌زده و لگاریتم مقدار مصرف حشره‌کش، در همه‌ی بلوک‌ها داده شده است. همچنان‌که در جدول ۱ دیده می‌شود، در بیش‌تر موارد با افزایش مقدار مصرف حشره‌کش، نسبت گوجه‌فرنگی‌های آفت‌زده در بلوک‌ها کاهش می‌یابد. در بلوک ۲ برای مقدار مصرف ۲/۱۲ (که از این به بعد آن را مشاهده‌ی ۱۴ می‌نامیم) تناقضی غیر عادی رخ می‌دهد. نسبت گوجه‌فرنگی‌های آفت‌زده برای این مقدار مصرف، ۵/۴۰ است که مقدار آن به نظر بالاست. این مشاهده را با این نگاه کاوشگرانه، می‌توان یک مشاهده‌ی دورافتاده تلقی کرد. بررسی آماری این که آیا این مشاهده واقعاً دورافتاده است یا خیر، و تأثیری که این مشاهده می‌تواند در نتایج یک مدل آماری داشته باشد، در بخش‌های بعد مورد بررسی قرار می‌گیرد.

جدول ۱. داده‌های مربوط به گوجه‌فرنگی‌های آفت‌زده

بلوک ۳			بلوک ۲			بلوک ۱			لگاریتم مقدار مصرف
$\frac{Y}{n}$	$n$	$Y$	$\frac{Y}{n}$	$n$	$Y$	$\frac{Y}{n}$	$n$	$Y$	
۰/۲۹۴	۳۴	۱۰	۰/۴۴۷	۳۸	۱۷	۰/۲۸۶	۳۵	۱۰	۱/۵۲۰
۰/۲۶۳	۳۸	۱۰	۰/۲۵۰	۴۰	۱۰	۰/۳۸۱	۴۲	۱۶	۱/۶۴۰
۰/۱۳۹	۳۶	۵	۰/۲۴۲	۳۳	۸	۰/۱۶۰	۵۰	۸	۱/۷۶۰
۰/۰۸۶	۳۵	۳	۰/۲۰۵	۳۹	۸	۰/۱۴۳	۴۲	۶	۱/۸۸۰
۰/۰۴۱	۴۹	۲	۰/۱۰۶	۴۷	۵	۰/۲۵۷	۳۵	۹	۲/۰۰۰
۰/۰۲۵	۴۰	۱	۰/۴۰۵	۴۲	۱۷	۰/۲۱۴	۴۲	۹	۲/۱۲۰
۰/۱۳۶	۲۲	۳	۰/۱۷۱	۳۵	۶	۰/۰۳۱	۳۲	۱	۲/۲۴۰
۰/۰۶۵	۳۱	۲	۰/۱۱۴	۳۵	۴	۰/۰۷۱	۲۸	۲	۲/۳۶۰

## ۳ تحلیل‌های گوناگون

در این بخش، روش‌های تحلیل تبدیل آرک‌سینوس، رگرسیون لوژیستیک، و شبه‌درست‌نمایی استوار را مورد بررسی قرار می‌دهیم.

### ۳/۱ روش تبدیل آرک سینوس

یکی از روش‌هایی که از آن می‌توان برای تحلیل داده‌های جدول ۱ استفاده کرد، این است که نخست تبدیلی مناسب برای پاسخ یافت تا آن را به متغیری تبدیل کند که دارای توزیع نرمال است و سپس با استفاده از رگرسیون معمولی، تأثیر متغیرهای تبیینی بر پاسخ تبدیل شده را یافت. این که آیا در این روش، پس از تحلیل، با بررسی مانده‌ها می‌توان داده‌های دورافتاده را تشخیص داد یا خیر، سؤالی است که در پی یافتن جوابی برای آن هستیم.

یکی از تبدیلی‌هایی که از آن می‌توان برای نسبت‌ها استفاده کرد تبدیل آرک سینوس، معرفی شده توسط آنسکامب [۲] است که سعی در تثبیت واریانس دارد و به صورت زیر تعریف می‌شود:

$$(۱) \quad \tilde{Y}_i = \sin^{-1} \sqrt{\frac{Y_i + \frac{3}{4}}{n_i + \frac{3}{4}}}.$$

$\tilde{Y}_i$  دارای واریانس مجانبی  $\frac{1}{4n_i+3}$  است. از آن‌جا که واریانس مجانبی به اندازه‌ی نمونه وابسته است، می‌توان برآورد پارامترها را برای بردار  $p$ -بعدی از متغیرهای تبیینی  $(x_i)$ ، به روش کم‌ترین توان‌های دوم وزنی (WLS) با مینیمم کردن  $\sum_i w_i (\tilde{Y}_i - \eta_i)^2$  که در آن  $w_i = 4n_i + 3$  و  $\eta_i = \beta^T x_i$  پیشگوی خطی است، به دست آورد [۷]. راه هم‌ارز دیگر آن است که از متغیر

$$(۲) \quad Z_i = (4n_i + 3)^{\frac{1}{2}} \tilde{Y}_i$$

که انتظار می‌رود به‌طور مجانبی برای هر  $i$  دارای واریانس تثبیت شده‌ی ۱ باشد، استفاده کنیم. در این صورت می‌توان از روش LS یا ماکسیمم درست‌نمایی، برآورد پارامترها را یافت. شکل ۱ نمودار پاسخ تبدیل شده‌ی  $Z_i$  را در مقابل لگاریتم مقدار مصرف برای همه‌ی بلوک‌ها نشان می‌دهد. همان‌طور که شکل

۱ نشان می‌دهد، با ابزار چشمی، هر چند می‌توان روند کلی کاهش پاسخ تبدیل شده‌ی  $Z_i$  با افزایش لگاریتم مقدار مصرف را مشاهده کرد، ولی خیلی سخت می‌توان در مورد تفاوت تأثیر بلوک‌های مختلف بر پاسخ تبدیل شده اظهار نظر کرد. مشاهده‌ی ۱۴ ممکن است با تأثیر زیاد، نتایج تحلیل رگرسیونی را تحت الشعاع خود قرار دهد. برای داده‌های بخش ۲ مدل زیر را در نظر می‌گیریم:

$$(۳) \quad Z_i = \beta_0 + \beta_1 \log(dose_i) + \beta_2 block_{2i} + \beta_3 block_{1i} + \varepsilon_i,$$

که در آن،  $Z_i$ ،  $i = 1, 2, \dots, 24$  پاسخ تبدیل شده برای  $i$ امین آزمایش،  $\beta_0, \beta_1, \beta_2$  و  $\beta_3$  پارامترهای مدل، و  $block_{ji}$  برای  $j = 1, 2$  متغیرهای نشانگر می‌باشند که اگر آزمایش  $i$ ام در بلوک  $j$ ام انجام گیرد مقدار ۱ و در غیر این صورت مقدار ۰ می‌گیرند.  $\varepsilon_i$  خطای اندازه‌گیری برای  $i$ امین آزمایش است که فرض

می‌شود دارای توزیعی نرمال با میانگین  $\mu$  و واریانس تثبیت‌شده  $\sigma^2$  است.  
مانند هر مدل رگرسیونی خطی دیگر، برای بررسی نیکویی برازش مدل، مانده‌های پی‌یرسونی عبارت‌اند از:

$$(4) \quad r_i = \frac{z_i - \hat{\mu}}{\hat{\sigma}},$$

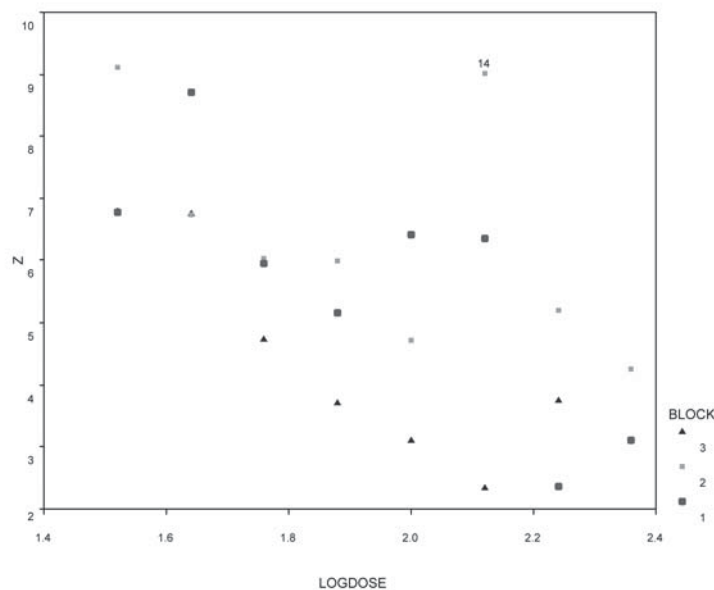
که در آن  $\hat{\mu}_i = E(Z_i | \mathbf{x}_i)$ . مانده‌ی پی‌یرسونی استاندارد شده عبارت است از:

$$(5) \quad \tilde{r}_i = \frac{r_i}{\sqrt{1 - \hat{h}_i}},$$

که در آن  $\hat{h}_i$  مؤلفه‌ی براورد شده‌ی ماتریس کلاهدار (معروف به تأثیر) می‌باشد. مانده‌ی کیش [۱] عبارت است از:

$$(6) \quad \sqrt{d_i} \operatorname{sgn}(z_i - \hat{\mu}_i),$$

که در آن  $\operatorname{sgn}(\cdot)$  تابع علامت است و  $d_i = (z_i - \hat{\mu}_i)^2$ .



شکل ۱. نمودار پاسخ تبدیل‌شده‌ی  $Z_i$  در مقابل لگاریتم مقدار مصرف برای همه‌ی بلوک‌ها

در رگرسیون خطی، نمودارهای (۱) مانده‌ها در مقابل مقادیر برازش‌یافته، (۲) نمودار Q-Q، (۳) جذر قدرمطلق مانده‌ی استاندارد شده در مقابل مقادیر برازش‌یافته، و (۴) نمودار فاصله‌ی کوک ( $h_i$ ) در مقابل شماره‌ی مشاهده، می‌توانند در تشخیص داده‌های دورافتاده و مؤثر مورد استفاده قرار گیرند [۷].

## ۳/۲ رگرسیون لوژستیک

در این روش فرض می‌شود متغیر پاسخ (تعداد موفقیت‌ها) دارای توزیع دوجمله‌ای است و چون این توزیع از خانواده‌ی نمایی است می‌توان از مدل‌های خطی تعمیم‌یافته (GLM) [۱۴] برای یافتن تأثیر متغیرهای تبیینی بر میانگین تعداد موفقیت‌ها (تعداد آفت زدن‌های گوجه‌فرنگی،  $\hat{Y}_i$ ، در مثال ما) کمک گرفت. در واقع اگر برای  $\mu_i = \mu(\mathbf{x}_i) = E(Y_i|\mathbf{x}_i)$  فرض کنیم

$$(۷) \quad \ln \frac{\mu(\mathbf{x}_i)}{n_i - \mu(\mathbf{x}_i)} = \beta^T \mathbf{x}_i,$$

که در آن  $\mathbf{x}$  بردار  $p$ -بعدی از متغیرهای تبیینی است و  $\beta$  بردار پارامترهای نامعلوم است، از یک GLM با پاسخ دوجمله‌ای و تابع ربط لوجیت استفاده کرده‌ایم. برآورد پارامترهای این مدل به راحتی می‌تواند با استفاده از اکثر نرم‌افزارهای آماری همچون SPSS، S-Plus یا R به دست آید. تابع درست‌نمایی برای پاسخ دوجمله‌ای با  $N$  آزمودنی و احتمال موفقیت  $\pi_i = \pi(\mathbf{x}_i)$  برای آزمودنی  $i$ ام، عبارت است از:

$$L(\beta|\mathbf{Y}, \mathbf{X}) = \prod_{i=1}^N \left[ \binom{n_i}{y_i} \pi^{y_i}(\mathbf{x}_i) \{1 - \pi(\mathbf{x}_i)\}^{n_i - y_i} \right],$$

که در آن  $\mathbf{Y} = (Y_1, \dots, Y_N)^T$  و  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  ماتریس طرح است و از آن‌جا

$$(۸) \quad \begin{aligned} \ell(\beta|\mathbf{Y}, \mathbf{X}) &= \ln\{L(\beta|\mathbf{Y}, \mathbf{X})\} \\ &= \sum_{i=1}^N \left[ \ln \binom{n_i}{y_i} + y_i \ln \pi(\mathbf{x}_i) + (n_i - y_i) \ln \{1 - \pi(\mathbf{x}_i)\} \right]. \end{aligned}$$

این تابع لگ درست‌نمایی به روش WLS تکراری [۱۴، ص ۴۰] ماکسیمم می‌شود و برآورد پارامترها به دست می‌آید.

مکولا و نلدر [۱۴] برای بررسی نیکویی برازش، مدل کپیش را به صورت زیر تعریف کرده‌اند:

$$(۹) \quad D(\mathbf{y}; \hat{\mu}) = -2\{\ell(\hat{\mu}; \mathbf{y}) - \ell(\mathbf{y}; \mathbf{y})\},$$

که در آن،  $\ell(\hat{\mu}; \mathbf{y})$  ماکسیمم لگ درست‌نمایی برای مدل را نشان می‌دهد و  $\ell(\mathbf{y}; \mathbf{y})$  ماکسیمم لگ درست‌نمایی برای کلی‌ترین مدل است که برای هر یک از مشاهدات، پارامتر مجزایی دارد و برآورد  $\mu$  برای آن،  $\mathbf{y}$  است. این مدل، مدل اشباع‌شده نامیده می‌شود.  $D(\mathbf{y}; \hat{\mu})$  دارای توزیع تقریبی خی دو با درجه‌ی آزادی  $N - p$  است، که  $N$  تعداد مشاهدات و  $p$  تعداد پارامترهای مدل است. اگر  $Y_i$  دارای توزیع دوجمله‌ای با پارامترهای  $n_i$  و  $\pi_i$  باشد، کیش عبارت است از:

$$D = \sum_{i=1}^N \left\{ y_i \ln \frac{y_i}{n_i \hat{\pi}_i} + (n_i - y_i) \ln \frac{n_i - y_i}{n_i - n_i \hat{\pi}_i} \right\}.$$

برای مقایسه‌ی دو مدل آشیانی به‌ترتیب با  $p$  و  $q$  پارامتر ( $q < p < N$ ) از اختلاف بین کیش آن‌ها می‌توان استفاده کرد. تغییر در کیش به‌صورت زیر تعریف می‌شود:

$$(۱۰) \quad \Delta D = 2 \{ \ell(\mathbf{b}_1; \mathbf{y}) - \ell(\mathbf{b}_0; \mathbf{y}) \},$$

که در آن  $\mathbf{b}_1$  و  $\mathbf{b}_0$  به‌ترتیب برآورد بردار پارامترها تحت مدل‌های با  $q$  و  $p$  پارامتر هستند.  $\Delta D$  دارای توزیع تقریبی خی دو با  $p - q$  درجه‌ی آزادی است. مانده‌ی پی‌یرسونی  $i$ امین آزمودنی برای نیکویی برازش مدل، عبارت است از:

$$(۱۱) \quad r_i = \frac{y_i - \hat{\mu}_i}{\{\widehat{\text{var}}(y_i)\}^{\frac{1}{4}}},$$

و مانده‌ی کیش عبارت است از:

$$(۱۲) \quad \sqrt{d_i} \text{sgn}(y_i - \hat{\mu}_i),$$

که در آن،

$$d_i = y_i \ln \frac{y_i}{n_i \hat{\pi}_i} + (n_i - y_i) \ln \frac{n_i - y_i}{n_i - n_i \hat{\pi}_i}.$$

$i$ امین مانده‌ی پی‌یرسونی استاندارد شده عبارت است از:

$$(۱۳) \quad \tilde{r}_i = \frac{r_i}{\sqrt{1 - \hat{h}_i}}.$$

برای داده‌های ارائه شده در بخش ۲، مدل رگرسیون لوژیستیک زیر را در نظر خواهیم گرفت:

$$(۱۴) \quad \ln \frac{\mu(\mathbf{x}_i)}{n_i - \mu(\mathbf{x}_i)} = \beta_0 + \beta_1 \log(\text{dose}_i) + \beta_2 \text{block}_{2i} + \beta_3 \text{block}_{3i}.$$

### ۳/۳ روش شبه‌درست‌نمایی استوار

برای بررسی استواری برآوردها از معیار تابع تأثیر استفاده می‌شود. در روش اِم‌برآوردها، به جای مینیم کردن مجموع توان دوم مانده‌ها، مجموع  $\sum_{i=1}^n \rho(r_i)$  برای یافتن برآوردها مینیم می‌شود، که در آن  $\rho$  تابعی متقارن است. تابع تأثیر یک اِم‌برآوردها، متناسب با تابع امتیازش است، که تابع امتیاز، مشتق  $\sum_{i=1}^n \rho(r_i)$  نسبت به ضرایب رگرسیونی است (همیل [۱۱] را ببینید). بنا براین، برآوردهای ماکسیم درست‌نمایی در GLM که نوعی اِم‌برآوردها هستند، استوار نیستند؛ یعنی، انحراف‌های بزرگ متغیر پاسخ از میانگین‌اش یا نقاط مؤثر روی محور رگرسیونی، می‌توانند تأثیر بزرگی روی برآوردها داشته باشند. در نتیجه، با استفاده از یک تابع امتیاز کراندار می‌توان برآوردها را به دست آورد. کانتونی و رونجی [۵] با توجه به برآوردهای شبه‌درست‌نمایی، برآوردهای شبه‌درست‌نمایی استوار را به دست آورده‌اند. در GLM، برآوردهای شبه‌درست‌نمایی پارامترهای مدل، از حل معادلات برآوردها زیر به دست می‌آیند:

$$(۱۵) \quad \sum_{i=1}^n \frac{y_i - \mu_i}{V(Y_i)} \mu'_i = 0,$$

که در آن،  $\mu_i = E(Y_i | \mathbf{x}_i)$ ،  $\mu'_i = \frac{\partial \mu_i}{\partial \beta}$  و  $V(Y_i) = \text{var}(Y_i)$ . توجه کنید که مانند هر GLM دیگر، برای تابع ربط یکنوا  $h(\cdot)$  داریم  $h(\mu_i) = \beta^T \mathbf{x}_i$ . در این روش تابع ربط، ارتباط واریانس پاسخ با میانگین آن و استقلال پاسخ‌ها را دانسته فرض می‌کنند؛ اما فرضی در مورد توزیع پاسخ اعمال نمی‌شود. در روشی که توسط کانتونی و رونجی [۵] پیشنهاد شده است، برآوردها از حل معادلات برآوردها زیر به دست می‌آیند:

$$(۱۶) \quad \sum_{i=1}^n \Psi(y_i, \mu_i) = 0,$$

که  $\Psi(y_i, \mu_i) = \nu(y_i, \mu_i) w(\mathbf{x}_i) \mu'_i - a(\beta)$  و

$$a(\beta) = \frac{1}{n} \sum_{j=1}^n E\{\nu(Y_j, \mu_j)\} w(\mathbf{x}_j) \mu'_j$$

جمله‌ای تصحیح‌کننده است که ما را از سازگاری فیشر مطمئن می‌کند. یا به‌طور معادل، عبارت زیر نسبت به  $\beta$  مینیم می‌شود:

$$Q_M(\mathbf{y}, \boldsymbol{\mu}) = \sum_{i=1}^n Q_M(y_i, \mu_i),$$

که:

$$Q_M(y_i, \mu_i) = \int_{\tilde{s}}^{\mu_i} \nu(y_i, t) w(\mathbf{x}_i) dt - \frac{1}{n} \sum_{j=1}^n \int_{\tilde{t}}^{\mu_j} E\{\nu(Y_j, t)\} w(\mathbf{x}_i) dt,$$



و  $\tilde{t}$  و  $\tilde{s}$  به ترتیب به گونه‌ای انتخاب می‌شوند که  $\nu(y_i, \tilde{s}) = 0$  و  $E\{\nu(Y_i, \tilde{t})\} = 0$ . در این معادلات، اثر داده‌ی دورافتاده‌ی پاسخ با انتخاب تابع کرانداری برای  $\nu(y_i, \mu_i)$  کنترل می‌شود. همچنین اثر مشاهدات مؤثر در متغیرهای تبیینی با انتخاب مناسب  $w(\mathbf{x}_i)$  کنترل می‌شود [۵].

برای داده‌های مثال بخش ۲، با وجود یک مشاهده‌ی دورافتاده در متغیر پاسخ و نداشتن نقطه‌ی مؤثر در متغیرهای تبیینی، در این معادله فرض می‌کنیم:

$$w(x_i) = 1 \quad \text{و} \quad \nu(y_i, \mu_i) = \frac{\Psi_c(r_i)}{\sqrt{V(\mu_i)}}$$

که در آن،  $r_i = \frac{y_i - \mu_i}{\sqrt{V(\mu_i)}}$  مانده‌های پی‌یرسونی و  $\Psi_c(r_i)$  تابع هوبر است که برابر است با:

$$\Psi_c(r_i) = \begin{cases} r_i, & |r_i| \leq c \\ c \times \text{sgn}(r_i), & |r_i| > c \end{cases}.$$

بنا بر این، برآوردها از حل معادلات برآوردها زیر به دست می‌آیند:

$$(۱۷) \quad \sum_{i=1}^n \left\{ \frac{\Psi_c(r_i)}{\sqrt{V(\mu_i)}} \mu'_i - a(\beta) \right\} = 0,$$

که در آن،

$$a(\beta) = \frac{1}{n} \sum_{j=1}^n E \left\{ \frac{\Psi_c(r_j)}{\sqrt{V(\mu_j)}} \right\} \mu'_j$$

برای مدل‌های دوجمله‌ای و پواسون دارای فرمی بسته است و نیازی به انتگرال‌گیری عددی ندارد [۵]. ثابت تنظیم‌کننده‌ی هوبر،  $c$ ، برای حفظ کارایی تقریبی، برابر با  $\frac{1}{4}$  در نظر گرفته می‌شود. با این شرایط، برآوردها به دست آمده، برآوردها شبه‌درست‌نمایی هوبر نامیده می‌شود. برآوردها با حل معادله‌ی (۱۷) با استفاده از روش نیوتون-رافسون یا امتیازبندی فیشر به دست می‌آیند.

برای ارزیابی نیکویی برازش، مدل شبه‌کیش به صورت زیر تعریف می‌شود:

$$(۱۸) \quad D_{QM}(\mathbf{y}, \boldsymbol{\mu}) = -2Q_M(\mathbf{y}, \boldsymbol{\mu}) = -2 \sum_{i=1}^n Q_M(y_i, \mu_i).$$

برای انتخاب بین دو مدل آشیانی، از اختلاف شبه‌کیش دو مدل استفاده می‌شود، که عبارت است از:

$$(۱۹) \quad \Lambda_{QM} = \{D_{QM}(\mathbf{y}, \boldsymbol{\mu}^\circ) - D_{QM}(\mathbf{y}, \hat{\boldsymbol{\mu}})\},$$

که در آن،  $\mu^*$  و  $\hat{\mu}$  به ترتیب با محاسبه‌ی برآورد پیشگوی خطی تحت فرض صفر (دارای  $p - q$  پارامتر) و برآورد پیشگوی خطی تحت فرض مقابل (دارای  $p$  پارامتر) به دست می‌آیند. برای به دست آوردن مقدار احتمال، لازم است که توزیع تقریبی آماره‌ی آزمون محاسبه شود. کانتونی و رونجی [۵] نشان داده‌اند که توزیع تقریبی  $\Lambda_{QM}$  مانند توزیع  $\sum_{i=1}^q d_i N_i^*$  است، که  $N_i$  ها متغیرهای نرمال استاندارد مستقل،  $d_i$  ها ویژه‌مقدارهای مثبت ماتریس  $Q(\Psi, F)C(\Psi, F)$ ، و  $q$  بعد پارامتر تحت فرض صفر است، که

$$Q(\Psi, F) = E\{\Psi(Y, \mu)\Psi^T(Y, \mu)\},$$

$$C(\Psi, F) = M^{-1}(\Psi, F) - \widetilde{M}^+(\Psi, F),$$

که در آن،

$$M(\Psi, F) = -E\left\{\frac{\partial}{\partial \beta}\Psi(Y, \mu)\right\},$$

و اگر ماتریس  $M(\Psi, F)$  به صورت زیر افراز شود:

$$M(\Psi, F) = \begin{pmatrix} M_{11}(\Psi, F) & M_{12}(\Psi, F) \\ M_{21}(\Psi, F) & M_{22}(\Psi, F) \end{pmatrix},$$

که در آن،  $M_{11}(\Psi, F)$  از مرتبه‌ی  $(p - q) \times (p - q)$ ،  $M_{12}(\Psi, F)$  از مرتبه‌ی  $(p - q) \times q$ ،  $M_{21}(\Psi, F)$  از مرتبه‌ی  $q \times (p - q)$ ، و  $M_{22}(\Psi, F)$  از مرتبه‌ی  $q \times q$  است، آن‌گاه  $\widetilde{M}^+(\Psi, F)$  ماتریسی است که  $\widetilde{M}_{11}^+(\Psi, F) = M_{11}^{-1}(\Psi, F)$  و  $\widetilde{M}_{12}^+(\Psi, F) = \widetilde{M}_{21}^+(\Psi, F) = \widetilde{M}_{22}^+(\Psi, F) = \mathbf{0}$  است. کانتونی و رونجی [۵]، از روش ایمهاف [۱۳] برای محاسبه‌ی مقادیر بحرانی توزیع  $\Lambda_{QM}$  استفاده کرده‌اند، که نیاز به اجرای الگوریتم‌های عددی طولانی دارد. ما برای محاسبه‌ی مقدار احتمال، از روش گشتاوری با تقریب ساتریتویت ([۶] را برای یک مثال داده‌شده ملاحظه کنید) که محاسبات ساده‌تری را در بر دارد استفاده می‌کنیم و توزیع تجمعی ترکیب خطی از متغیرهای تصادفی با توزیع خی‌دو را با توزیع خی‌دویی که به درجه‌ی آزادی آن تقسیم شده است و این درجه‌ی آزادی برابر است با

$$\nu = \frac{\left(\sum_{i=1}^q d_i N_i^*\right)^2}{\sum_{i=1}^q d_i^2 N_i^{*2}}, \quad (20)$$

تقریب می‌زنیم.

## ۴ نتایج تحلیل‌ها

جدول ۲ نتایج تحلیل با استفاده از روش‌های گوناگون را نشان می‌دهد.

جدول ۲. نتایج تحلیل با روش‌های مختلف

پارامتر	برآورد ML در مدل لوژیستیک	برآورد ML در روش تبدیل	برآورد شبه‌درست‌نمایی استوار
$\beta_0$	۱,۴۸۰ (۰,۶۵۶)	۱۳,۱۰۷ (۱,۹۵۹)	۱,۹۳۹ (۰,۷۰۰)
$\beta_1$	-۱,۸۱۷ (۰,۳۴۳)	-۴,۵۴۵ (۰,۹۸۱)	-۲,۰۴۹ (۰,۳۷۰)
$\beta_2$	۰,۸۴۳ (۰,۲۲۶)	۲,۰۹۲ (۰,۶۶۰)	۰,۶۸۵ (۰,۲۴۰)
$\beta_3$	۰,۵۴۲ (۰,۲۳۲)	۱,۳۱۳ (۰,۶۶۰)	۰,۴۵۰ (۰,۲۴۰)
مانده‌ی کبیش درجه‌ی آزادی کبیش	۳۹,۹۷۶ ۲۰	۳۴,۸۸۶ ۲۰	۳۲,۷۴۰ ۲۰

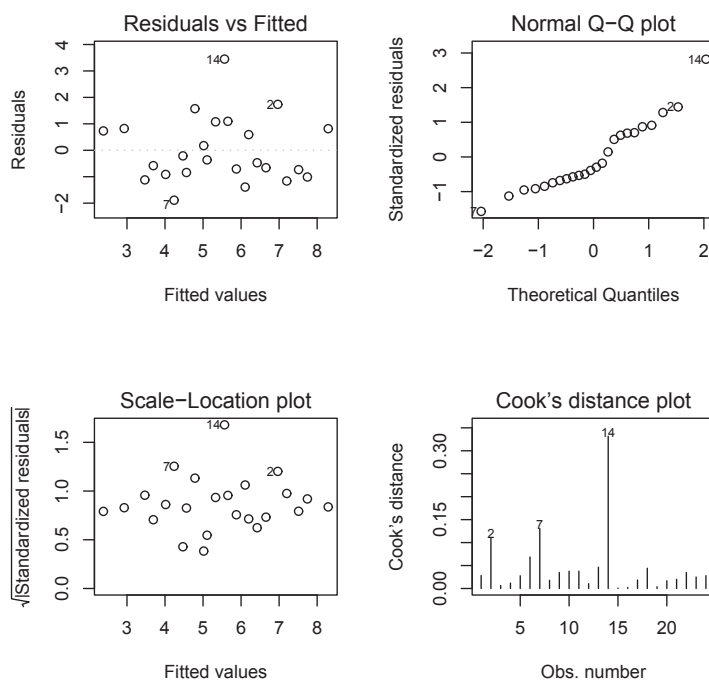
همان‌طور که جدول ۲ نشان می‌دهد، مانده‌ی کبیش در همه‌ی مدل‌ها نسبت به درجه‌ی آزادی آن زیاد است، ولی در مدل شبه‌درست‌نمایی استوار، این مقدار، کم‌ترین است. کبیش زیاد به دو دلیل ممکن است رخ دهد: (۱) بیش‌پراکندگی و (۲) وجود نقطه‌ی دورافتاده [۱۴]. در روش تبدیل،  $\hat{\sigma} = ۱,۳۲۱$  و  $R^2 = ۰,۶۱۴$  به دست می‌آید. این برآورد زیاد  $\sigma$  ممکن است به دلیل وجود داده‌ی دورافتاده باشد. شکل‌های ۲ و ۳ چهار نمودار بحث‌شده در زیربخش ۳/۱ را به‌ترتیب برای مدل با استفاده از روش تبدیل آرک‌سینوس و مدل رگرسیون لوژیستیک نشان می‌دهند.

شکل ۲ نشان می‌دهد که مشاهده‌ی ۱۴، هم دارای مانده‌ای زیاد است و هم نسبت به سایر مشاهدات، دارای تأثیری زیاد است. به هر صورت، از آن‌جا که مانده بر انحراف معیار بیش‌برآورد شده تقسیم می‌شود تا مانده‌ی پی‌یرسونی به دست آید، ممکن است این روش، توانایی تشخیص داده‌ی دورافتاده را نداشته باشد و داده‌ی دورافتاده‌ای به غلط داده‌ای معمولی تلقی شود. شکل ۳ نشان می‌دهد که مشاهده‌ی ۱۴ یک داده‌ی دورافتاده است. البته در این حالت، رگرسیون لوژیستیک و بررسی‌های تشخیصی آن، تأکید بیش‌تری بر دورافتاده بودن این مشاهده دارد.

حال که با استفاده از رگرسیون لوژیستیک از دورافتاده بودن مشاهده‌ی ۱۴ مطمئن شده‌ایم، دو راه برای

مواجهه با این مشاهده وجود دارد: یا این مشاهده باید از تحلیل خارج شود و بدون آن، تحلیل دوباره انجام شود، که این با مشورت با محقق اصلی باید انجام گیرد، یا این که از روش‌های استواری که نسبت به داده‌ی دورافتاده حساس نیستند استفاده شود. اجازه دهید تا در حالت اول، تغییر در نتایج را جویا شویم. پس از حذف مشاهده‌ی ۱۴، در مدل با استفاده از روش تبدیل،  $R^2 = 0.727$  به دست می‌آید و مانده‌ی کپش به مقدار  $21.0^\circ$  تقلیل می‌یابد.  $\hat{\sigma} = 17.51$  به دست می‌آید و برازش پس از بررسی نمودارهای تشخیصی، برازش نیکویی را نتیجه می‌دهد. در این روش، همه‌ی متغیرهای تبیینی، معنادارند.

پس از حذف مشاهده‌ی ۱۴، در رگرسیون لوژیستیک، مانده‌ی کپش به مقدار قابل ملاحظه‌ای کاهش یافته، برابر با  $25.289$  می‌شود. در این روش نیز همه‌ی متغیرها معنادارند و نمودارهای تشخیصی، برازش مناسبی را نشان می‌دهند. بنا بر این، در صورت حذف داده‌ی دورافتاده، هر یک از دو روش تبدیل و رگرسیون لوژیستیک می‌تواند مورد استفاده قرار گیرد. پس از حذف مشاهده‌ی ۱۴، در روش استوار نیز مانده‌ی کپش کاهش می‌یابد و به  $24.63$  می‌رسد و همه‌ی متغیرهای تبیینی معنادار می‌شوند. البته در این روش، متغیر مربوط به بلوک ۱ به طور ضعیف معنادار است (پی مقدار  $= 0.08$ ).

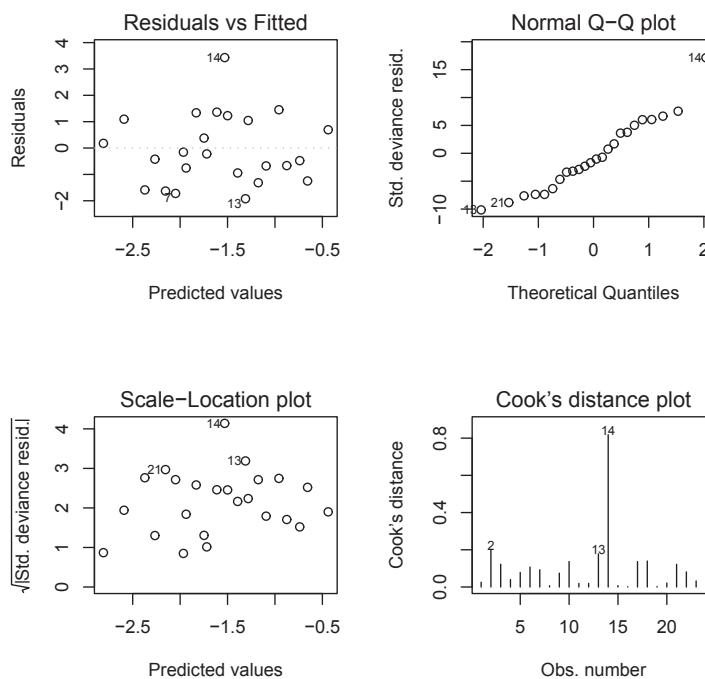


شکل ۲. نمودارهای تشخیصی برای روش تبدیل آرک سینوس

در حالتی که داده‌ی دورافتاده در تحلیل باقی بماند، از آن‌جا که در روش استوار، به این داده وزن کمتری داده می‌شود، نتایج تحلیل، کم‌تر تحت تأثیر این داده قرار می‌گیرد. وزن داده‌شده به آزمودنی  $i$ ام،  $\frac{\Psi_c(r_i)}{r_i}$  است که برای مشاهده‌ی ۱۴، برابر با ۰/۲۶ و برای بقیه‌ی مشاهدات بین ۰/۷ تا ۱ است.

## ۵ نتیجه‌گیری

روش‌های تشخیصی دورافتاده بودن مشاهدات در تبدیل آرک‌سینوس، کم‌تر از روش‌های تشخیصی رگرسیون لوژیستیک، توانایی نشان دادند. این ممکن است به دلیل آن باشد که در نظر گرفتن واریانس مجانبی برای برخی از مشاهدات در بعضی موارد امکان‌پذیر نیست. در هر صورت، وجود مشاهده‌ی دورافتاده باعث می‌شود که روش تبدیل آرک‌سینوس، واریانس را بیش‌برآورد کند و این امر نشان‌دهنده‌ی این است که این روش در تشخیص دورافتاده بودن برخی مشاهده‌ها توانایی کم‌تری نشان می‌دهد. اگر پس از تشخیص دورافتاده



شکل ۳. نمودارهای تشخیصی برای روش رگرسیون لوژیستیک

بودن داده‌ای، بتوان این داده را از تحلیل خارج کرد، می‌توان با برازش دوباره و بررسی‌های تشخیصی، از نیکویی برازش اطمینان حاصل کرد. در غیر این صورت، استفاده از روش‌های استوار که نتایج آن‌ها زیاد تحت تأثیر داده‌ی دورافتاده قرار نمی‌گیرد، ترجیح داده می‌شود.

## مرجع‌ها

- [1] Agresti, A. (2002). *Categorical Data Analysis*, 2nd ed. Wiley, New York.
- [2] Anscombe, F.J. (1948). The transformation of Poisson, binomial and negative binomial data. *Biometrika* **35**, 246-254.
- [3] Berkson, J. (1957). Tables for the maximum likelihood estimate of the logistic function. *Biometrics* **13**, 28-34.
- [4] Box, G.E.P. (1953). Non-normality and tests on variances. *Biometrika* **40**, 318-335.
- [5] Cantoni, E.; Ronchetti, E. (2001). Robust inference for generalized linear models. *J. Amer. Statist. Assoc.* **96**, 1022-1030.
- [6] Casella, G.; Berger, R.L. (1990). *Statistical Inference*. Brooks/Cole, Pacific Grove, CA.
- [7] Draper, N.R.; Smith, H. (1981). *Applied Regression Analysis*, 2nd ed. Wiley, New York.
- [8] Finney, D.J. (1952). *Probit Analysis*, 2nd ed. Cambridge University Press, UK.
- [9] Finney, D.J. (1964). *Statistical Method in Biological Assay*, 2nd ed. Griffin, London.
- [10] Fisher, R.A.; Yates, F. (1938). *Statistical Tables for Biological, Agricultural and Medical Research*. Oliver and Boyd, Edinburgh.
- [11] Hampel, F.R. (1974). The influence curve and its role in robust estimation. *J. Amer. Statist. Assoc.* **69**, 383-393.
- [12] Huber, P.J. (1964). Robust estimation of a location parameter. *Ann. Math. Statist.* **35**, 73-101.
- [13] Imhof, J.P. (1961). Computing the distribution of quadratic forms in normal variables *Biometrika* **48**, 419-426.
- [14] McCullagh, P.; Nelder, J.A. (1989). *Generalized Linear Models*, 2nd ed. Chapman and Hall, London.
- [15] Nelder, J.A.; Wedderburn, R.W.M. (1972). Generalized linear models. *J. Roy. Statist. Soc. A* **135**, 370-384.

- [16] Phellps, K. (1982). Use of the complementary log-log function to describe dose response relationships in insecticide evaluation field trails. In *Proceedings of the International Conference on Generalized Linear Models*, R. Gilchrist, ed. Springer, New York.
- [17] Wedderburn, R.W.M. (1974). Quasi-likelihood functions, generalized linear models, and the Gauss-Newton method. *Biometrika* **61**, 439-447.
- [18] Williams, D.A. (1987). Generalized linear model diagnostics using the deviance and single case deletions. *Appl. Statist.* **36**, 181-191.

دریافت: ۱۰ مرداد ۱۳۸۳

آخرین اصلاح: ۱۱ آذر ۱۳۸۳

انتشار: ۲۱ دی ۱۳۸۳

مجتبی گنجعلی

گروه آمار، دانشکده‌ی علوم ریاضی،

دانشگاه شهید بهشتی،

بولوار دانشگاه، اوین،

تهران، ایران.

پایم‌نگار: *m-ganjali@cc.sbu.ac.ir*

# Analysis of a Problem Using Various Visions

M. Ganjali<sup>†</sup> and L. Latifian<sup>‡</sup>

<sup>†</sup> Shahid Beheshti University

<sup>‡</sup> Central Bank of the Islamic Republic of Iran

**Abstract.** In this paper an applied problem, where the response of interest is the number of success in a specific experiment, is considered and by various visions is studied. The effects of outlier values of response on results of a regression analysis are so important to be studied. For this reason, using diagnostic methods, outlier response values are recognized. It is shown that use of arc-sine transformation may be misleading in recognizing response outliers. If deleting of outliers is not possible, use of robust modeling approach is suggested. Method of maximum likelihood for estimating parameters in generalized linear model, transformation method and also pseudo-likelihood method are not robust. A method, which is called robust pseudo-likelihood and leads to robust results, is reviewed and a simpler method of computing  $P$ -values for model selection is presented. Various approaches for modeling are also compared in the applied example.

© 2004 Statistical Research Center. All rights reserved.

**Keywords.** transformation method; logistic regression; generalized linear models; Pearson residuals; robust pseudo-likelihood.