

## مقایسه مدل‌های آماری حاشیه ای و آمیخته در تحلیل داده های پزشکی

میر سعید یکانی نژاد<sup>۱</sup>، مهدی یاسری<sup>۲\*</sup>، کرامت نوری جلیانی<sup>۲</sup>، آرش اکابری<sup>۳</sup>

<sup>۱</sup> دانشجوی دوره دکتری تخصصی آمار زیستی، گروه اپیدمیولوژی و آمار زیستی، دانشکده بهداشت، دانشگاه علوم پزشکی تهران، تهران، ایران

<sup>۲</sup> دانشیار آمار زیستی، گروه اپیدمیولوژی و آمار زیستی، دانشکده بهداشت، دانشگاه علوم پزشکی تهران، تهران، ایران

<sup>۳</sup> کارشناس ارشد آمار زیستی، دانشگاه علوم پزشکی خراسان شمالی، بجنورد، ایران

\* نویسنده مسئول: تهران، دانشگاه علوم پزشکی تهران، دانشکده بهداشت، گروه اپیدمیولوژی و آمار زیستی

پست الکترونیک: m.yasari@yahoo.com

### چکیده

**زمینه و هدف:** امروزه داده های زیادی وجود دارند که در آنها فرض استقلال داده ها که پیش فرض اصلی بسیاری از مدل‌های آماری است برقرار نیست. داده های حاصل از نمونه گیری خوشه ای، مطالعات طولی با اندازه گیری های مکرر و یا داده های زوجی مانند داده های دو چشم و همچنین مطالعات همسان سازی شده نمونه هایی از این داده ها هستند

**مواد و روش کار:** دو مدل آماری با در نظر گرفتن همبستگی بین مشاهدات، مدل‌های آمیخته و مدل‌های حاشیه ای هستند که در این مقاله مورد مقایسه قرار گرفته اند. این مدل‌ها از نظر روشی که همبستگی بین داده ها را در نظر می گیرند و تفسیر ضرایب رگرسیونی با یکدیگر تفاوت دارند.

**یافته ها:** در مدل‌های غیر خطی ضرایب رگرسیونی مدل حاشیه ای تغییرات در سطح جامعه را نشان می دهد در حالیکه در مدل‌های آمیخته این ضرایب تغییرات را در یک فرد یا خوشه خاص نشان می دهند. در مدل‌های خطی تفسیر ضرایب رگرسیونی در هر دو مدل یکسان است.

**نتیجه گیری:** به عنوان یک نتیجه گیری کلی در مدل‌های غیر خطی که تفسیر نتایج متفاوت است بر اساس کاربرد مورد نیاز بهتر است از یکی از این دو مدل استفاده شود. در واقع برای سیاستگذاران بهداشتی که دید جامعه نگر دارند استفاده از مدل حاشیه ای توصیه می شود در حالیکه برای پزشکان معالج که بیشتر به تغییرات فردی در یک بیمار علاقمند هستند استفاده از مدل آمیخته مناسبتر است.

**واژه های کلیدی:** مدل حاشیه ای، مدل آمیخته، مدل خطی، مدل غیر خطی

### مقدمه

برقرار نمی باشد لذا استفاده از این مدل‌ها گرچه ممکن است به آریبی برآوردگرها نیانجامد و یا مقدار این آریبی قابل اغماض باشد با این وجود خطای معیار برآورد ضرایب به شدت تحت تاثیر پیش فرض استقلال داده ها است [۱-۳]. به عبارت دیگر در صورتیکه داده ها از هم مستقل نباشند خطای معیار و در نتیجه فاصله اطمینان و نتیجه آزمون فرضها برای صفر بودن ضرایب رگرسیونی غیر قابل اعتماد خواهند بود. لذا در آنالیز این داده ها استفاده از روشهایی که بتواند این وابستگی را در نظر بگیرد ضروری است. روشهای مختلفی برای حل این مساله ابداع شده

در بسیاری از داده های پزشکی فرض استقلال داده ها برقرار نمی باشد برای مثال داده هایی که به روش نمونه گیری خوشه ای جمع آوری می شوند و یا مطالعاتی که بصورت طولی انجام می شوند و در آن اندازه گیری های مکرر وجود دارد، و یا داده هایی که به طور همزمان از یک یا چند عضو بدن جمع آوری می شود و همچنین هنگامیکه در یک مطالعه همسان سازی انجام می شود فرض استقلال داده ها را نمی توان پذیرفت زیرا داده های یک خوشه، یک فرد در طول زمان یا نمونه های همسان شده به هم وابسته هستند. پیش فرض مورد نیاز در برازش مدل‌های خطی تعمیم یافته فرض استقلال داده هاست

(۳)

$$E(y_i | b_i) = X_i \beta + z_i b_i$$

که  $b$  یک بردار از اثرات تصادفی است که از خوشه ای به خوشه ی دیگر در تغییر است. معرفی اثرات تصادفی منجر به ایجاد ساختار کوواریانس اثرات تصادفی می شود.

(۴)

$$\sum i = Z_i G Z_i' + \delta^2 I_{ni}$$

در دو مدل مذکور تفسیر ضرایب  $\beta$  یکسان است و در هر دوی آنها ضریب  $\beta$  نشان دهنده ی آن است که متغیرهای کمکی به چه میزان، میانگین متغیر پاسخ را در جامعه مورد مطالعه تغییر می دهند. [۲،۳] این تفسیر از  $\beta$  در معادله (۱) کاملاً آشکار است و مدل‌های حاشیه‌ای مستقیماً به مدل (۱) اشاره دارند از این رو می‌توان تفسیر  $\beta$  را به تغییرات در میانگین جامعه تعمیم داد. جالب است که در مورد مدل‌های آمیخته ی خطی نیز می‌توان تفسیر حاشیه ای از میانگین متغیر پاسخ داشت زیرا:

$$E(y_i) = E\{E(y_i | b_i)\}$$

(۵)

$$= E(X_i \beta + z_i b_i)$$

$$= X_i \beta + z_i E(b_i)$$

$$= X_i \beta$$

پارامتر  $\beta$  در معادلات (۱) و (۳) تفسیر یکسانی دارد. زیرا معادله ی (۳) یک مدل آمیخته خطی است. (به عبارتی سمت راست (۳) یک تابع خطی از  $b$  است.)

به بیان ساده‌تر در مدل اثرات آمیخته خطی، ضرایب  $\beta$  تفسیر حاشیه‌ای دارد زیرا میانگین تغییرات خطی متغیر پاسخ بر اساس متغیرهای کمکی در هر سطح برابر تغییرات میانگین متغیر پاسخ در جامعه بر اساس متغیرهای کمکی است.

#### بخش دوم: مدل‌های تعمیم یافته ی خطی

در این بخش به مقایسه ی مدل‌های حاشیه‌ای و مدل‌های تعمیم یافته آمیخته می‌پردازیم همانطور که می‌دانیم یکی از مولفه‌های مدل‌های تعمیم یافته، تابع ربط<sup>۵</sup> است.

است که مهمترین آنها مدل‌های آمیخته<sup>۱</sup> و مدل‌های حاشیه‌ای<sup>۲</sup> هستند. گرچه از ابداع این مدل‌ها مدت زمان زیادی می‌گذرد اما هنوز اتفاق نظر در مورد برتری یکی از این مدل‌ها بر دیگری وجود ندارد. [۱، ۴-۶] برخی از آمار شناسان برای هر یک از این دو مدل هویت مستقل قائل هستند و برتری هریک از این مدل‌ها را بر اساس کاربرد آنها تعریف می‌کنند. [۴، ۵، ۷] در این مقاله به بررسی تفاوت‌های مدل‌های آمیخته و مدل‌های حاشیه ای می‌پردازیم. این مقاله شامل دو بخش اصلی است که در این دو بخش تفاوت‌های میان مدل‌های آمیخته و مدل‌های حاشیه ای در دو کلاس بزرگ مدل‌های خطی<sup>۳</sup> و مدل‌های خطی تعمیم یافته<sup>۴</sup> مورد بررسی قرار می‌گیرد. در هر یک از این دو بخش ابتدا تفاوت در نحوه در نظر گرفتن همبستگی در این دو مدل بررسی می‌شود و سپس تفسیر ضرایب رگرسیونی در آنها توضیح داده می‌شود و در پایان با یک مثال عددی و یک مثال واقعی حاصل از طرح سلامت و بیماری در سال ۱۳۹۰ اختلاف در تفسیر ضرایب در این دو مدل نشان داده شده و تفاوت در کاربرد آنها بیان می‌شود.

#### بخش اول: مدل‌های خطی

در مدل‌های کلاسیک آماری میانگین متغیر پاسخ به صورت تابعی خطی از متغیرهای مستقل بیان می‌شود و به منظور در نظر گرفتن همبستگی مشاهدات از دو روش عمده استفاده می‌شود در مدل‌های حاشیه ای یک الگوی همبستگی تعریف می‌شود که به صورت یک ماتریس همبستگی است

(۱)

$$E(y_i) = X_i \beta$$

(۲)

$$\sum i = \text{cor}(y_i)$$

روش دوم معرفی اثرات تصادفی در مدل برای میانگین متغیر پاسخ است.

1 - Mixed model

2 - Marginal model

3 - Linear model

4 - Generalized linear models

5 - Link function

اکنون مدل‌های تعمیم یافته ی آمیخته ی خطی را برای میانگین شرطی متغیر پاسخ، به شرط برداری از اثرات تصادفی  $b_i$  در نظر می‌گیریم.

(۹)

$$g[E(y_i | b_i)] = X_i \beta^* + Z_i b_i$$

در اینجا پارامترهای ثابت را با  $\beta^*$  نمایش می‌دهیم تا از نظیر آنها در مدل‌های حاشیه ای متمایز باشند. در اینجا پارامترهای رگرسیونی تفسیر شرطی برای یک خوشه یا یک فرد دارند. یعنی  $\beta^*$  بر حسب یک واحد تغییر در متغیر کمکی، زمانی که  $b_i$  ثابت در نظر گرفته می‌شود بیان می‌شود (۸). معمول ترین روش برای ثابت در نظر گرفتن اثر تصادفی آن است که میانگین شرطی متغیر پاسخ را در نظر بگیریم بنابراین بر خلاف  $\beta$  در مدل‌های حاشیه ای،  $\beta^*$  بر اساس تغییر در میانگین متغیر پاسخ در هر خوشه و یا هر فرد تفسیر می‌شود و پارامترهای رگرسیونی  $\beta^*$  در مورد تغییرات متغیر پاسخ در جامعه اطلاعی را نمی‌دهند.

در صورتی که بخواهیم میانگین حاشیه ای متغیر پاسخ را بدست آوریم باید از توزیع اثرات تصادفی  $b_i$  میانگین بگیریم. که این مستلزم میانگین گرفتن از یک تابع غیر خطی از اثرات تصادفی،  $b_i$ ، است.

$$\mu_i = E(y_i)$$

$$= E\{E(y_i | b_i)\}$$

$$= E\{h(X_i \beta^* + z_i b_i)\}$$

برای میانگین گرفتن از یک تابع غیر خطی از اثرات تصادفی، میانگین وزنی بر اساس توزیع اثرات تصادفی می‌گیریم:

$$\mu = E(y_i) = E\{h(X_i \beta^* + z_i b_i)\} = \int_{-\infty}^{\infty} h(X_i \beta^* + z_i b_i) f(b_i) db_i$$

در معادله بالا، انتگرال گیری معادل میانگین گرفتن می‌باشد و مقادیر تابع چگالی احتمال  $f(b_i)$  وزنه‌های مورد استفاده هستند همانطور که از انتگرال فوق بر می‌آید

$g(\mu_i)$  که میانگین متغیر پاسخ را به تابعی خطی از متغیرهای کمکی مرتبط می‌سازد. در بحثی که در مقایسه ی دو مدل، در بخش اول انجام گرفت به عنوان یک حالت خاص از مدل‌های تعمیم یافته، تابع ربط یک تابع همانی بود.  $g(\mu_i) = \mu_i$  که در این حالت تفسیر پارامترهای رگرسیونی  $\beta$  در دو مدل آمیخته و مدل حاشیه ای یکسان است. در اینجا حالتی را مورد بررسی قرار خواهیم داد که تابع ربط غیر خطی است. مدل حاشیه‌ای تعمیم یافته ی خطی در این حالت بصورت زیر بیان می‌شود.

(۶)

$$g(\mu_i) = g\{E(y_i)\} = X_i \beta$$

که در آن  $g(\cdot)$  یک تابع ربط غیر خطی متناسب با متغیر پاسخ است. در مدل‌های حاشیه‌ای تفسیر پارامترهای رگرسیونی  $\beta$  به صورت تغییر در میانگین متغیر پاسخ در جامعه مورد مطالعه بر اساس تغییرات متغیرهای کمکی است. برای مثال، زمانی که پاسخ بصورت بردار  $y_i$  دو وضعیتی است و تابع ربط لجیت برای آنها استفاده می‌شود تفسیر پارامترهای رگرسیونی  $\beta$  به صورت میزان تغییر در لگاریتم نسبت شانس در جامعه است

(۷)

$$\log \text{it}(\mu_i) = X_i \beta$$

برای هر تابع ربط مفروض،  $g(\cdot)$ ، میانگین جامعه را می‌توان بر اساس معکوس تابع ربط بیان کرد.

(۸)

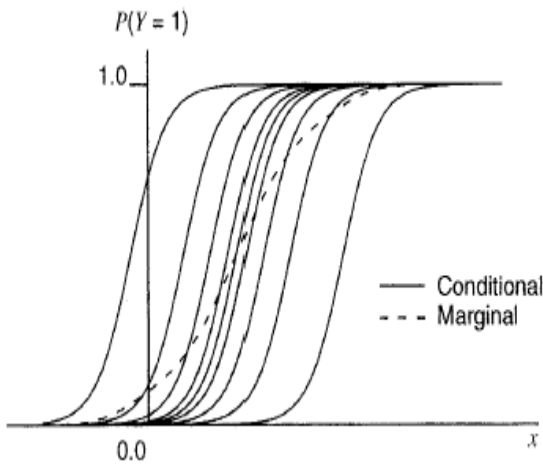
$$h\{g(\mu_i)\} = \mu_i = E(y_i) = h(X_i \beta)$$

به عنوان مثال، زمانی که  $y_i$  ها دو وضعیتی هستند و تابع ربط لجیت است مدلی که برای  $\mu_i$  در نظر گرفته می‌شود به صورت زیر است:

که در اینجا  $h(\cdot)$  تابع معکوس لجیت است. پارامترهای رگرسیونی  $\beta$ ، در هر دو فرم (۶) و (۸) نشان دهنده ی تغییر در میانگین متغیر پاسخ در جامعه به ازای تغییرات در متغیرهای کمکی هستند.

شکل ۱ علت کوچکتر بودن ضرایب در مدل‌های حاشیه‌ای نسبت به مدل با عرض از مبدا تصادفی به نمایش درآمده است

برای یک متغیر کمکی با نام  $X$  نمودارهای مختص فردی  $P(Y_{ij} = 1 | b_i)$  را برای افراد گوناگون وقتی تغییرپذیری قابل ملاحظه‌ای وجود دارد نمایش داده است. این تغییر پذیری به علت وجود یک تاثیر تصادفی با  $\sigma$  بزرگ است.



در هر مقدار ثابت از  $X$  تغییر پذیری زیادی، در میانگین‌های شرطی وجود دارد  $E(Y_{ij} | b_i)$  میانگین این میانگین‌های شرطی (با تاثیر تصادفی) همان میانگین حاصل از مدل حاشیه ای است یعنی  $E(Y_{ij})$  این میانگین گیری برای مقادیر مختلف  $X$  موجب ایجاد یک منحنی کلی برای مدل حاشیه ای می شود. این منحنی شیب کمتری نسبت به تمامی منحنی‌های مدل‌های با تاثیر تصادفی دارد. در حقیقت این منحنی دقیقاً از توزیع لجستیک تبعیت نمی‌کند. چنین خاصیتی برای بسیاری از توابع ربط نیز وجود دارد.

البته لازم به ذکر است که علاوه بر تابع ربط همانی در تابع ربط لگاریتم نیز هر دو ضریب یک تفسیر دارند. مدل با تاثیرات تصادفی عرض از مبدا را در نظر بگیرید.

$$\ln[E(y_i | b_i)] = X_i \beta^* + b_i$$

$$\mu_i = E(y_i) = E[E(y_i | b_i)] = E[e^{X_i \beta^* + b_i}] = e^{X_i \beta^*} \times E(e^{b_i})$$

$$\Rightarrow \ln[E(y_i)] = X_i \beta^* + \ln E(e^{b_i})$$

میانگین متغیر پاسخ برابر با  $h(X_i \beta)$  نمی باشد.

$$E(y_i) \neq h(X_i \beta) \quad (11)$$

برای مثال در مدل رگرسیون لجستیک با عرض از مبدا تصادفی داریم:

$$b_i \approx N(0, \delta_b^2)$$

$$\logit \{E(y_i | b_i)\} = X_i \beta^* + b_i$$

میانگین حاشیه ای از متغیر پاسخ برابر است با:

$$\mu_i = E(y_i)$$

$$= E\{E(y_i | b_i)\}$$

$$= E\left\{\frac{e^{(X_i \beta^* + b_i)}}{1 + e^{(X_i \beta^* + b_i)}}\right\}$$

$$= \int_{-\infty}^{\infty} \frac{e^{(X_i \beta^* + b_i)}}{1 + e^{(X_i \beta^* + b_i)}} \frac{1}{\sqrt{2\pi} \delta_b} e^{-\frac{b_i^2}{2\delta_b^2}} db_i$$

این انتگرال به راحتی قابل محاسبه نیست و واضح است که از فرم رگرسیون لجستیک نمی باشد.

$$(12)$$

$$\frac{e^{(X_i \beta^* + b_i)}}{1 + e^{(X_i \beta^* + b_i)}}$$

رابطه حاصله را می توان بصورت تقریبی بفرم زیر بیان کرد.

$$\logit E(y_i) \approx (1 + k \sigma_b^2)^{-\frac{1}{2} X_i \beta^*}$$

که در آن  $k = \frac{16\sqrt{3}}{15\pi}$  طرز رسیدن به این فرمول خیلی مهم نیست. آنچه که اهمیت دارد این است که چگونه ضریب در مدل با عرض از مبدا تصادفی بیشتر از مدل های حاشیه‌ای می گردد.

$$\beta = \frac{\beta^*}{\sqrt{1 + 0.346 \sigma_b^2}} \quad k^2 = 0.346$$

بنابراین هر گاه  $Var(b_i) = \sigma_b^2 > 0$  پارامترهای مدل حاشیه ای رگرسیون لجستیک  $\beta$  مقادیر کوچکتر از پارامترهای مدل رگرسیونی لجستیک با تاثیر تصادفی  $\beta^*$  می‌گردند. و این تفاوت با افزایش  $\sigma_b^2$  افزایش می یابد. در

اگر تفاضل ریسک را به عنوان یک مدل خطی از ریسک بیمار شدن در نظر بگیریم، (در اینجا تفاضل ریسک بیماری پس از درمان و قبل از درمان، به عنوان شاخصی از اثر بخشی درمان بکار می رود). این اختلاف ریسک ها در ستون چهارم جدول (۱) نشان داده شده است و مقادیر آن از  $-0.09$  تا  $-0.17$  متغیر است.

این مقادیر اثرات فردی درمان می باشند. حال به دو طریق می توان یک شاخص از اثر بخشی دارو ارائه داد در روش اول از اثرات فردی دارو برای هر سه فرد میانگین می گیریم .

$$\frac{-0.13 - 0.17 - 0.09}{3} = -0.13$$

در روش دیگر میانگین ریسک بیمار شدن را قبل از درمان ( $0.5$ ) با پس از درمان ( $0.37$ ) مقایسه می کنیم. این روش در واقع مقایسه ریسک بیمار شدن در جامعه است .

$$0.13 = (0.37 - 0.50)$$

بنابراین تفاضل میانگین های ریسک بیماری قبل از درمان و پس از درمان در جامعه برابر میانگین اختلاف ریسک ها در افراد است. بطور خلاصه «تفاضل میانگین ها» برابر «پمیانگین تفاضل» است. این مثال ساده‌ی عددی موید آن است که تفسیر ضرایب رگرسیونی در مدل‌های خطی آمیخته و مدل‌های حاشیه ای یکسان است.

حال، یک تابع غیر خطی از ریسک بیمار شدن را در نظر می گیریم . لگاریتم نسبت شانس را به عنوان یک تابع غیر خطی، که ریسک بیمار شدن را پس از درمان و قبل از درمان مقایسه می کند در نظر می گیریم. لگاریتم نسبت شانس برای افراد A, B, C در ستون پنجم جدول (۱) نمایش داده شده است.

و از آنجایی که  $e^{b_i}$  یک عدد ثابت است ( اگر  $b_i \sim N(0, \sigma_b^2)$  آنگاه  $[\ln[E(e^{b_i})]] = \sigma_b^2$ ). بنابراین برای دو تابع ربط همانی و لگاریتم ضرایب رگرسیونی یکسان هستند.

اما بطور کلی حتی با انتگرال گیری روی کلیه مقادیر  $b_i$  نیز بصورت مستقیم ضرایب مدل حاشیه ای قابل حصول نیست.

### بحث

حال با بیان یک مثال عددی از داده های فرضی و یک مثال کاربردی از داده های طرح سلامت و بیماری به بررسی بیشتر مدل‌های حاشیه ای و آمیخته و مقایسه بین آنها می پردازیم.

### مثال عددی

برای روشن تر شدن مطلب به بیان یک مثال ساده می پردازیم تا تفاوت‌های اساسی مدل‌های آمیخته و حاشیه ای را مشخص کنیم. یک مجموعه از داده های فرضی را در جدول (۱) در نظر بگیرید. این مجموعه از داده ها نشان دهنده ی ریسک واقعی بیمار شدن سه فرد در شروع مطالعه و پس از درمان با یک داروی جدید به منظور کاهش دادن خطر بیماری می باشد. این سه فرد از نظر

ریسک بیمار شدن در ابتدای مطالعه کاملا متفاوت هستند. که این تفاوت را می توان به صورت اثرات تصادفی  $b_i$  بیان کرد. یعنی سه فرد A, B, C دارای ریسک های شدید، متوسط، ضعیف از بیمار شدن هستند اگر فرض کنیم که همه جامعه نیز به نسبت مساوی دارای ریسک های فوق می باشند. آنگاه سطر آخر جدول (۱) میانگین های ریسک جامعه را قبل از درمان ، پس از درمان و اختلاف ریسک را نشان می دهند.

جدول ۱: داده های فرضی از ریسک بیمار شدن، قبل و پس از درمان، برای سه فرد با ریسک متفاوت برای بیمار شدن

شخص	قبل از درمان	پس از درمان	تفاضل	لگاریتم نسبت شانس
A	۰/۸۰	۰/۶۷	-۰/۱۳	-۰/۶۸
B	۰/۵۰	۰/۳۳	-۰/۱۷	-۰/۷۱
C	۰/۲۰	۰/۱۱	-۰/۹۰	-۰/۷۰
میانگین جامعه	۰/۵۰	۰/۳۷	۰/۱۳	

دهنده ی تغییر در شیوع بیماری در جامعه مورد مطالعه است زمانیکه تمام افراد جامعه تحت درمان باشند. یعنی با مصرف دارو در جامعه شانس بیمار شدن در جامعه ۴۰٪ کاهش می یابد. ( $1 - e^{-0.532} \approx 0.4$ ) این مقدار برآورد برای محققان و سیاست گذاران بهداشتی حائز اهمیت است. زیرا نشان دهنده ی اثر کلی درمان در سطح جامعه است.

### مثال کاربردی

در این بخش به مقایسه مدل‌های آماری رگرسیون لجستیک معمولی، حاشیه ای و آمیخته در تعیین عوامل موثر بر تغذیه با شیر مادر خواهیم پرداخت. داده‌های مورد استفاده در این مثال مربوط به مطالعه‌ی بررسی سلامت و بیماری است. طرح سلامت و بیماری در سال ۱۳۷۸ در سطح کشور جمهوری اسلامی ایران اجرا شده که در آن ۱۳۴۷۵ خانوار بر اساس یک نمونه گیری چند مرحله ایی خوشه ایی و به نسبت یک هزارم جمعیت کل خانوارهای کشور و بطور تصادفی انتخاب گردیده اند

از مجموع ۱۳۴۷۵ خانوار طرح سلامت ۱۷۵۶ خانوار کودک زیر دو سال دارند، که در این ۱۷۵۶ خانوار مجموعاً ۲۰۱۰ کودک کمتر از ۲ سال وجود دارند. که نمونه مورد مطالعه را تشکیل می دهند. جهت تعیین عوامل موثر بر تغذیه با شیر مادر در کودکان کمتر از دو سال، بر اساس اطلاعات پرسشنامه کودکان به دو گروه استفاده کننده از شیر مادر ( $Y=1$ ) و کودکان که با شیر مادر تغذیه نمی شوند ( $Y=0$ ) تقسیم شدند.

نتایج حاصل از سه مدل مورد بررسی در جدول ۲ ارائه شده است با نگاهی به این جدول در می یابیم که خطای معیار ضرائب متغیرها در مدل آمیخته و حاشیه‌ای نسبت به مدل معمولی که در آن همبستگی درون خوشه‌ای در نظر گرفته شده بیشتر برآورد شده است مقایسه  $P$ - مقادیرهای این مدل‌ها نشان داد که در مدل معمولی  $P$  مقدار کمتری نسبت به دو مدل آمیخته و حاشیه‌ای دارد. با مقایسه مدل حاشیه‌ای و آمیخته در می یابیم که با وجود اینکه در هر دو مدل خطای معیار ضرایب افزایش یافته، اما میزان افزایش خطای معیار در مدل حاشیه ای به اندازه مدل آمیخته نیست. از آنجایی که مدل‌های آمیخته و حاشیه‌ای روشهای کاملاً متفاوتی برای در نظر گرفتن

برای مثال لگاریتم نسبت شانس برای فرد  $A$  برابر است با:

$$\log \left\{ \frac{0.67/(1-0.67)}{0.8/(1-0.8)} \right\} = -0.68$$

مقادیر نسبت شانس برای افراد مورد مطالعه خیلی شبیه هستند. و از  $-0.68$  تا  $-0.71$  در حال تغییرند. که این مقادیر اثرات فردی دارو بر ریسک بیماری می باشند. بار دیگر به دو روش می توانیم یک شاخص کلی از اثر بخشی دارو ارائه بدهیم. در روش اول از اثرات فردی، که همان مقادیر لگاریتم نسبت شانس برای تک تک افراد هستند میانگین می گیریم:

$$\frac{-0.68 - 0.71 - 0.70}{3} = -0.697$$

این مقدار نشان دهنده ی آن است که درمان، خطر بیمار شدن را به نصف کاهش می دهد. (زیرا  $e^{-0.697} \approx 0.5$ ) در روشی دیگر به عنوان شاخصی از اثر بخشی درمان لگاریتم نسبت شانس قبل از درمان،  $\log\left(\frac{0.5}{0.5}\right) = 0$  را با لگاریتم نسبت شانس پس از درمان،  $\log\left(\frac{0.37}{0.63}\right) = -0.532$  مقایسه می کنیم که در نتیجه این مقایسه مقدار  $-0.532$  برای تفاضل لگاریتم نسبت شانس برای هر فرد حاصل می شود که از  $-0.697$  کوچکتر است. بنابراین در یک تابع غیر خطی از ریسک بیمار شدن « مقایسه غیر خطی از میانگین ها » برابر با « میانگین مقایسه های غیر خطی » نیست.

حال این سوال پیش می آید که کدام یک از دو مقدار بدست آمده برای اثر بخشی دارو واقعیت را نشان می دهند؟ جواب این سوال این است که هر دوی این مقادیر، مقادیری واقعی هستند اما به سوالهای متفاوتی پاسخ می دهند.

برآورد  $-0.697$  یک معیار از تغییرات مورد انتظار در ریسک بیماری برای فردی است که با دارو درمان شده است. یعنی خطر بیماری برای فردی که دارو را استفاده می کنند به میزان ۵۰٪ درصد کاهش ( $1 - e^{-0.532} \approx 0.5$ ) می یابد.

این برآورد برای پزشک معالج بیمار ارزش ویژه ای دارد زیرا او علاقه مند به دانستن اثر متوسط درمان بر روی هر فرد است. از طرف دیگر، مقدار برآورد  $-0.532$  نشان

همبستگی مشاهدات دارند خطای معیار حاصل از این مدلها متفاوت است. همانطور که در جدول ۲ مشاهده می‌شود متغیرهای

بنابراین تفسیر ضرایب رگرسیونی در مدل آمیخته وابسته به اثرات تصادفی است.

جدول ۲: مقایسه نتایج سه مدل رگرسیونی در بررسی عوامل موثر بر تغذیه با شیر مادر

رگرسیون لجستیک						رده	متغیر
آمیخته		حاشیه‌ای		معمولی			
P	(خطای معیار) ضریب	P	(خطای معیار) ضریب	P	(خطای معیار) ضریب		
<۰/۰۰۱	-۰/۱۰۴ (۰/۰۲۱)	<۰/۰۰۱	-۰/۰۹۸ (۰/۰۱۴)	<۰/۰۰۱	-۰/۱۰۶ (۰/۰۱۲)		سن کودک (ماه)
	پایه		پایه		پایه	مرد	جنس کودک
۰/۰۲۷	-۰/۴۴۷ (۰/۱۸۱)	۰/۰۰۵	-۰/۳۸۵ (۰/۱۷۴)	۰/۰۰۲	-۰/۴۲۸ (۰/۱۴۹)	زن	
	پایه		پایه		پایه	۲۵ و کمتر	سن اولین
۰/۰۵۹	-۰/۵۷۵ (۰/۳۴۸)	۰/۰۲۱	-۰/۵۲۳ (۰/۲۷۶)	۰/۰۴۰	-۰/۶۰۶ (۰/۲۳۵)	بالاتر از ۲۵	
	پایه		پایه		پایه	۲۴ و کمتر	فاصله دو تولد آخر
۰/۰۰۴	۰/۵۸۷ (۰/۲۴۷)	<۰/۰۰۱	۰/۵۰۵ (۰/۱۷۳)	<۰/۰۰۱	۰/۵۷۵ (۰/۱۴۸)	بیش از ۲۴	

جنس کودک، سن اولین حاملگی، و فاصله تولد در مدل حاشیه‌ای در مقایسه با مدل آمیخته اثر معنی‌داری بر روی تغذیه با شیر مادر دارند که این به دلیل است که خطای معیار در مدل آمیخته بیشتر است.

همانطور که دیدیم تفسیر ضرایب رگرسیونی در مدل‌های حاشیه‌ای و مدل‌های آمیخته کاملاً متفاوت می‌باشد و این دو مدل به دو جامعه هدف متفاوت اشاره دارند در مدل‌های حاشیه‌ای ضرایب رگرسیونی مقدار تغییرات در میانگین جامعه یا مقدار تبدیل یافته میانگین جامعه را به ازای تغییر در متغیر مستقل نشان می‌دهند. و یک مشخصه منحصر بفرد مدل‌های حاشیه‌ای در آن است که مدل‌های رگرسیونی برای میانگین متغیر پاسخ و الگوی همبستگی مشاهدات به صورت مجزا تعیین می‌شوند. تفکیک مدل رگرسیونی برای میانگین متغیر پاسخ و مدل همبستگی داده‌ها از یکدیگر موجب می‌شود که تفسیر ضرایب رگرسیونی در مدل حاشیه‌ای مستقل از الگوی در نظر گرفته شده برای وابستگی مشاهدات باشد و همانطور که از لفظ حاشیه‌ای برمی‌آید ضرایب رگرسیونی وابسته به هیچ اثری جز متغیرهای کمکی نیستند.

در مقایسه با مدل‌های حاشیه‌ای اساس مدل‌های آمیخته، تغییر پذیری ذاتی در خوشه‌هاست و همبستگی بین مشاهدات از این واقعیت ناشی می‌شود که ضرایب رگرسیونی برای مجموعه‌ای از مشاهدات یکسان است.

به عنوان مثال . در مدل حاشیه‌ای با هر ماه افزایش سن کودک لگاریتم شانس تغذیه با شیر مادر در جامعه به میزان ۰/۰۹۸ کاهش می‌یابد در حالیکه براساس مدل آمیخته با هر ماه افزایش سن کودک لگاریتم شانس تغذیه با شیر مادر در یک خوشه بطور متوسط ۰/۱۰۴ کاهش می‌یابد.

#### نتیجه گیری

زمانیکه تابع ربط به فرم همانی و یک تابع خطی باشد تفسیر پارامترهای رگرسیونی در مدل‌های آمیخته همانند مدل‌های حاشیه‌ای است. در حالیکه وقتی تابع ربط مورد استفاده به شکل غیر خطی است تفسیر پارامترهای رگرسیونی در دو مدل کاملاً متفاوت است. تفسیر نتایج برآوردهای حاصل از مدل آمیخته برای پزشکان معالج ارزش ویژه‌ای دارد زیرا آنها علاقه مند به دانستن اثر متوسط درمان بر روی هر فرد هستند. از طرف دیگر، تفسیر نتایج برآوردهای حاصل از مدل حاشیه‌ای تغییر در پارامتر مورد بررسی در جامعه مورد مطالعه است این مقدار برآورد برای محققان و سیاست‌گذاران بهداشت حائز اهمیت است. گرچه باید دقت کرد که در برخی از توابع ربط مانند همانی این دوضرب مفهوم یکسانی می‌یابند و هر دو کاربرد را دارند.

**References**

1. Lee Y, Nelder J, Conditional and Marginal Models: Another View, *Statist Sci*, 2004;27(2):20.
2. Fitzmaurice GM, Larid N, Ware J, *Applied longitudinal analysis*, New Jersey: wiley & sons; 2004.
3. Agresti A, *Categorical Data Analysis* Hoboken, New Jersey :wiley & sons; 2003.
4. Carriere I, Bouyer J, Choosing marginal or random-effects models for longitudinal binary responses: application to self-reported disability among older persons, *BMC Med Res Methodol*, 2002; Dec 5;2:15.
5. Hadgu A, Koch G, Westrom L, Analysis of ectopic pregnancy data using marginal and conditional models, *Stat Med* , 1997 ;Nov 15;16(21):2403-17.
6. Ten Have TR, Ratcliffe SJ, Reboussin BA, Miller ME, Deviations from the population-averaged versus cluster-specific relationship for clustered binary data, *Stat Methods Med Res*, 2004 Feb;13(1):3-16.
7. Lindsey JK, Lambert P, On the appropriateness of marginal models for repeated measurements in clinical trials, *Stat Med*, 1998 Feb 28;17(4):447-69.
8. Larsen K, Petersen JH, Budtz-Jorgensen E, Endahl L, Interpreting parameters in the logistic regression model with random effects, *Biometrics* , 2000 Sep;56(3):909-14.



## Comparison of marginal and mixed models in medical data analysis

yekaninejad MS<sup>1</sup>, Yaseri M<sup>1\*</sup>, Nourijelyani K<sup>2</sup>, Akaberi A<sup>3</sup>

<sup>1</sup> PH.D student of Biostatistics, Department of Epidemiology and Biostatistics, School of Public Health, Tehran University of Medical Sciences, Tehran, Iran

<sup>2</sup> Associated professor of Biostatistics, Department of Epidemiology and Biostatistics, School of Public Health, Tehran University of Medical Sciences, Tehran, Iran

<sup>3</sup> M.Sc of Biostatistic, North Khorasan University of Medical Sciences, Bojnurd, Iran

**\*Corresponding Author:**  
Department of Epidemiology  
and Biostatistics, School of  
Public Health, Tehran  
University of Medical Sciences,  
Tehran, Iran  
Email: m.yaseri@yahoo.com

---

### Abstract

**Background & Objectives:** In medical researchers, there are lots of correlated data which cannot be analyzed using the usual classical statistical methods because the assumption of independency between observations is not met. Data from cluster sampling, longitudinal studies, observations on paired organs and matched studies are examples of such data.

**Materials & Methods:** Two statistical methods that can be used to correctly handle these kinds of data are marginal and mixed models. These models are different in considering correlation between subjects and Interpretation of the regression coefficients. These models were compared in this paper.

**Results:** The regression coefficients in marginal models with non-identity link functions show the change in population whereas in mixed models they represent changes within a subject or a cluster.

**Conclusion:** In result, in nonlinear models, application of these two kinds of models depends on the areas of their usage. While the marginal models are more attractive to the Health policy makers who are considering the potential effects of a variable on the population as a whole, the mixed model will be of most interest to a physician in a physician/patients context.

**Key words:** Marginal model, Mixed model, Linear model, Non linear model

---



