

تحلیل پاسخ‌های شمارشی طولی برای تعداد ادعاهای خسارت با تعداد صفر زیاد در پرتفوی بیمه شخص ثالث کشور ایران

فرید صلواتی^۱

تاریخ دریافت: ۱۳۹۱/۱۰/۱۸

احسان بهرامی سامانی^۲

تاریخ پذیرش: ۱۳۹۳/۰۳/۱۲

چکیده

در بیمه شخص ثالث به دلیل وجود سیستم پاداش- جریمه و استفاده از سیستم پاداش آخر سال، بیمه‌گذار خسارت‌های کوچک خود را به شرکت بیمه گزارش نمی‌دهد. این کار باعث ایجاد صفرهای زیاد در تعداد ادعای خسارت بیمه‌گذار می‌شود. از سوی دیگر تحلیل تعداد ادعای خسارت و عوامل تشکیل‌دهنده خطر روی این پاسخ، برای شرکت‌های بیمه حایز اهمیت است. برای این منظور، برخی از مدل‌ها با پاسخ‌های شمارشی با استفاده از توزیع سری‌های توانی مانند مدل رگرسیون پواسون و مدل رگرسیون دوجمله‌ای منفی و توزیع سری‌های توانی آماسیده صفر مانند مدل رگرسیون پواسون آماسیده صفر و رگرسیون دوجمله‌ای منفی آماسیده صفر برای تحلیل داده‌های بیمه شخص ثالث با تعداد صفر زیاد استفاده می‌شود. در این مقاله می‌توان این مدل‌ها را برای داده‌های طولی بیمه شخص ثالث با تعداد صفر زیاد تعمیم داد. یک شیوه درستی مبنای برای به‌دست‌آوردن برآورد پارامترهای مدل استفاده شده است. در این روش از الگوریتم EM نیز در برآورد پارامترها برای مدل‌هایی با پاسخ آماسیده صفر استفاده شده است. در نهایت برای تشریح سودمندی مدل‌های پیشنهادشده، داده‌های واقعی طولی بیمه شخص ثالث، مورد تجزیه و تحلیل قرار گرفته است. **واژگان کلیدی:** الگوریتم EM ، بیمه شخص ثالث، تعداد ادعای خسارت، توزیع آماسیده-صفر، داده‌های طولی

۱. مقدمه

از موضوع‌های مهم در بیمه شخص ثالث، تحلیل و انتخاب مدل مناسب برای برآزش روی داده‌های تعداد ادعای خسارت است. تحلیل رگرسیون با پاسخ‌های شمارشی، اجازه شناسایی عامل‌های ریسک و پیشگویی فراوانی مورد انتظار ادعاها با توجه به ویژگی‌های قرارداد را می‌دهد. در بیمه منظور از پاسخ‌ها، تعداد ادعاهای خسارتی است که شخص بیمه‌گذار به شرکت‌های بیمه گزارش می‌دهد. شرکت‌های بیمه برای به‌دست‌آوردن حق‌بیمه از مدل‌های شمارشی شناخته‌شده مانند مدل پواسون^۱ برای تعداد ادعای خسارت استفاده می‌کنند که این مدل‌ها در بیمه اتومبیل به دلیل آنکه این نوع داده‌ها دارای صفرهای زیادی می‌باشند، کارایی پایین‌تری نسبت به مدل‌های آماسیده صفر^۲ دارد. بیشتر این صفرها به دلیل گزارش دروغین بیمه‌گذارها برای استفاده از سیستم تخفیف آخر سال است. مدل‌های آماسیده صفر (که در ادامه معرفی می‌شود) به دلیل اینکه یک پارامتر احتمالی را برای این صفرها در نظر می‌گیرد، توانایی برآورد احتمال این ادعاهای دروغین را بر اساس سابقه بیمه‌گذار - که همان متغیرهای تبیینی و تعداد ادعای خسارت سال‌های گذشته بیمه‌گذار است - دارد.

این تحقیق بر آن است که با مقایسه مدل‌های پواسون و دو جمله‌ای منفی^۳ با مدل‌های پواسون آماسیده صفر^۴ و دو جمله‌ای منفی آماسیده صفر^۵، بهترین مدل را بر اساس معیارهای مورد نظر به‌دست‌آورد. برای این منظور علاوه بر تعداد ادعای خسارت، متغیرهای تبیینی مانند نوع اتومبیل، سن اتومبیل و محل رانندگی نیز در نظر گرفته می‌شود، سپس با وارد کردن این متغیرها در مدل‌های مورد بررسی، تأثیر آنها را در مدل مورد ارزیابی قرار می‌دهد. در نهایت بهترین مدل از بین مدل‌های شمارشی رایج و مدل‌های

-
1. Poisson
 2. Zero-Inflated
 3. Negative Binomial
 4. Zero-Inflated Poisson
 5. Zero-Inflated Negative Binomial

آماسیده صفر را برای برازش به داده‌های طولی تعداد ادعای خسارت با صفر زیاد معرفی می‌نماید.

۲. مروری بر ادبیات تحقیق

۲-۱. بیان مسئله و اهمیت تحقیق

برای تحلیل تعداد ادعای خسارت، پاسخ Y_{it} ، $(i=1, \dots, k, t=1, \dots, T)$ به عنوان تعداد ادعای خسارت بیمه‌گذار i ام در زمان t ام در نظر گرفته می‌شود. به طوری که این پاسخ‌ها یک پاسخ طولی است که به شرکت بیمه گزارش داده شده است. همچنین یک سری از عوامل تشکیل‌دهنده خطر وجود دارند که روی این پاسخ اثرگذار می‌باشند، به عنوان مثال در بیمه شخص ثالث به صورت نمونه برای متغیرهای تبیینی، می‌توان اطلاعات مربوط به راننده، وسیله نقلیه بیمه‌شده و محل رانندگی را نام برد. همچنین در این حالت، به دلیل وجود عوامل تشکیل‌دهنده خطری که قابل مشاهده و اندازه‌گیری نمی‌باشند اما بر تعداد تصادف‌ها و همچنین تعداد ادعاهای خسارت تأثیر دارند، نیاز به در نظر گرفتن اثرهای تصادفی در مدل‌های پیشنهادی با پاسخ‌های طولی است.

تاکنون هیچ تحقیقی در زمینه برازش و مقایسه مدل‌های رگرسیون آماسیده صفر برای تعداد ادعای خسارت در طول چند سال صورت نگرفته است. در این مقاله با معرفی مدل‌های آماسیده صفر طولی برای تعداد ادعای خسارت طولی به تشخیص عوامل تشکیل‌دهنده خطر مؤثر روی این پاسخ‌ها می‌پردازیم، همچنین با مقایسه مدل‌های شمارشی در این زمینه، بهترین مدل ممکن روی داده‌های بیمه شخص ثالث معرفی می‌شود.

۲-۲. پیشینه تحقیق

برای تحلیل و مدل‌سازی روی تعداد ادعای خسارت، مدل‌های شمارشی همچون مدل‌های رگرسیون پواسون^۱ و رگرسیون دوجمله‌ای منفی^۲ توسط توماس و سامسون^۳

-
1. Poisson Regression
 2. Negative Binomial Regression
 3. Thomas and Samson, 1987

با پاسخ‌های مقطعی مورد استفاده قرار گرفت. پژوهش‌های مختلفی نیز در زمینه محاسبه حق بیمه با استفاده از توزیع‌های چندمتغیره طولی پواسون و دوجمله‌ای منفی توسط بوچر و همکاران^۱ صورت گرفته است. از آنجایی که در تعداد ادعای خسارت ممکن است با صفر زیاد روبه‌رو شویم، نیاز به استفاده از مدل‌های رگرسیون آماسیده صفر برای برازش و تحلیل روی این داده‌هاست.

مدل رگرسیون پواسون آماسیده صفر نیز برای برازش به داده‌های شمارشی با صفر زیاد توسط لمبرت^۲ مورد استفاده قرار گرفت. همچنین با در نظر گرفتن توزیع دوجمله‌ای منفی به جای پواسون در مدل‌های آماسیده صفر، هیلبرن^۳ مدل دوجمله‌ای منفی آماسیده صفر را معرفی کرد. هال^۴ اظهار کرد که این مدل‌ها قابل تعمیم به پاسخ‌های طولی نیز می‌باشند و تحت عنوان مدل رگرسیون آماسیده صفر با پاسخ‌های طولی مطرح شد. این دو مدل کاربردهای بسیاری را روی تحلیل تعداد ادعای خسارت با صفر زیاد ایفا می‌کنند. ییپ و یائو^۵ نیز به مقایسه مدل‌های آماسیده صفر و مدل‌های شمارشی رایج در حالت مقطعی پرداخته‌اند. همچنین برای تعیین حق بیمه‌ها در بیمه اتومبیل برای تعداد ادعای خسارت در طول چند سال بوچر و گیلن^۶ از توزیع‌های آماسیده صفر استفاده کرده‌اند.

۲-۳. چند مفهوم مهم در این تحقیق

در این بخش با مفهوم پاسخ‌های طولی و توزیع‌های آماسیده صفر آشنا شده و سپس به بررسی مدل‌های رگرسیون آماسیده صفر مقطعی و طولی مورد استفاده برای تعداد ادعای خسارت با تعداد صفر زیاد در داده‌های بیمه شخص ثالث می‌پردازیم. از مهم‌ترین این

-
1. Boucher et al., 2008
 2. Lambert, 1992
 3. Heilbron, 1994
 4. Hall, 2000
 5. Yip and Yau, 2005
 6. Boucher and Guille'n, 2009

مدل‌ها می‌توان به مدل رگرسیون پواسون آماسیده صفر، مدل رگرسیون دوجمله‌ای منفی آماسیده صفر در حالت‌های مقطعی از زمان و در طول چند سال اشاره نمود.

۱-۳-۲. پاسخ‌های طولی

پاسخ‌هایی هستند که در طول زمان و فضای مشخصی جمع‌آوری می‌شوند که هدف اولیه آن آشکارسازی تغییرات پاسخ در طول زمان و همچنین عامل‌های تأثیرگذار روی این تغییرات است. این پاسخ‌ها می‌تواند تعداد ادعای خسارت‌ها در طول چند سال برای هر بیمه‌گذار باشد. اگر تعداد ادعای خسارت برای یک سال مشخص در نظر گرفته‌شود این پاسخ‌ها به صورت مقطعی خواهند بود اما اگر تعداد ادعای خسارت گزارش داده‌شده برای هر بیمه‌گذار در طول چند سال مورد نظر باشد، تعداد ادعای خسارت به صورت طولی در نظر گرفته می‌شود.

۲-۳-۲. توزیع‌های سری توانی

X دارای توزیع سری‌های توانی است، هرگاه تابع جرم احتمال آن به صورت زیر باشد:

$$P(X=x) = \frac{a(x)\theta^x}{c(\theta)} \quad x=0,1,\dots$$

که در آن $a(x) > 0$ و $c(\theta)$ تابعی مثبت، متناهی و مشتق‌پذیر از θ است. توزیع‌های پواسون و دوجمله‌ای منفی، نمونه‌ای از توزیع‌های سری توانی می‌باشند.

۳-۳-۲. مدل آماسیده صفر

مدل‌های آماسیده صفر یک پارامتر احتمالی اضافی را برای مقدار صفر که نمی‌تواند به‌طور کامل توسط فرض مدل برآورد شود، معرفی می‌کند و برای مدل‌هایی با بیش پراکنش و صفر زیاد به‌کار می‌رود که تابع احتمالی آن به صورت زیر است:

$$f_{ZID}(y) = \phi_0 I(y=0) + (1 - \phi_0) f_D(y|\theta)$$

- $f_D(y|\theta)$: تابع احتمال از توزیع D با پارامتر θ ;

- $f_{ZID}(y)$: مدل آماسیده صفر از توزیع D با یک پارامتر اضافی ϕ_0 برای احتمال آماسیده صفر.

۴-۳-۲. مدل رگرسیونی سری‌های توانی طولی

مدل رگرسیون سری‌های توانی آماسیده صفر زمانی به کار می‌رود که مدل سری‌های توانی برازش شده، مقدار صفر را کم برآورد کند یا به عبارت دیگر تعداد صفر زیاد رخ دهد. در این حالت یک پارامتر اضافه مانند ϕ_{it} به مدل اضافه می‌شود و توزیع آمیخته حاصل به این صورت است که احتمال ϕ_{it} را به رخداد صفر و احتمال $1-\phi_{it}$ را به رخداد توزیع سری‌های توانی می‌دهد. چون در مثال مورد بررسی تعداد صفر زیاد رخ می‌دهد در این بخش به بررسی مدل رگرسیون سری‌های توانی آماسیده صفر طولی پرداخته می‌شود. تابع جرم احتمال سری‌های توانی آماسیده صفر به شرط اثر تصادفی b_i به صورت زیر در نظر گرفته می‌شود:

$$P(Y_{it} = y_{it}) = \begin{cases} \phi_{it} + (1 - \phi_{it}) \frac{a(0)}{c(\lambda_{it})} & y_{it} = 0 \\ \phi_{it} (1 - \phi_{it}) \frac{a(y_{it}) \lambda_{it}^{y_{it}}}{c(y_{it})} & y_{it} = 1, 2, \dots \end{cases}$$

فرض کنید بردار پاسخ \mathbf{Y} شامل داده‌های K گروه مستقل $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_K^T)^T$ باشد که $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{iT_i})^T$ است. در این مدل داریم:

$$\boldsymbol{\lambda}_i = (\lambda_{i1}, \dots, \lambda_{iT_i})^T$$

$$\boldsymbol{\phi} = (\phi_{i1}, \dots, \phi_{iT_i})^T$$

در نتیجه مدل رگرسیون سری‌های توانی با پاسخ‌های طولی آماسیده صفر به این صورت در نظر گرفته می‌شود:

$$Y_{it} \sim ZIPS(\lambda_{it}, \phi_{it})$$

$$\text{Log}(\lambda_{it}) = \mathbf{B}'_{it} \boldsymbol{\beta}_t + \sigma b_i$$

$$\text{Logit}(\phi_{it}) = \mathbf{G}'_{it} \boldsymbol{\gamma}_t \quad i = 1, \dots, K, \quad t = 1, \dots, T$$

که در آن \mathbf{B}_{it} و \mathbf{G}_{it} بردار مربوط به متغیرهای تبیینی شامل عوامل تشکیل‌دهنده خطر و اثرگذار روی تعداد ادعای خسارت می‌باشند. از سوی دیگر \mathbf{b}_i به عنوان اثر تصادفی مربوط به فرد i به صورت توزیع گاما برای توزیع شرطی پواسون به شرط اثر تصادفی و توزیع بتا برای توزیع دو جمله‌ای منفی به شرط اثر تصادفی در نظر گرفته می‌شود (Hausman et al., 1984). اگر $\psi = (\gamma^T, \beta^T, \sigma)^T$ برداری مرکب از پارامترها باشد آنگاه لگاریتم تابع درست‌نمایی برای مدل رگرسیونی سری‌های توانی آماسیده صفر به این صورت به دست می‌آید:

$$LogL_{ZIPS}(\psi; y) = \sum_{i=1}^K Log \int_{-\infty}^{+\infty} [\prod_{t=1}^{T_i} Pr(Y_{it} = y_{it} | b_i)] g(b_i) db_i$$

که در آن:

$$\begin{aligned} Pr(Y_{it} = y_{it} | b_i) &= [(\phi_{it} + (1 - \phi_{it}) \frac{a(0)}{c(\lambda_{it})})^{1-u_{it}} ((1 - \phi_{it}) \frac{a(y_{it}) \lambda_{it}^{y_{it}}}{c(\lambda_{it})})^{u_{it}}] \\ &= [(\frac{\exp(G_{it} \gamma_t)}{1 + \exp(G_{it} \gamma_t)} + \frac{1}{1 + \exp(G_{it} \gamma_t)} \frac{a(0)}{c(\exp(B_{it} \beta_t + \sigma \beta_i))})^{1-u_{it}} \times \\ &(\frac{1}{1 + \exp(G_{it} \gamma_t)} \frac{a(y_{it}) \exp(B_{it} \beta_t + \sigma \beta_i)^{y_{it}}}{c(\exp(B_{it} \beta_t + \sigma \beta_i))})^{u_{it}}] \end{aligned}$$

همچنین u_{it} به این صورت تعریف می‌شود:

$$u_{it} = \begin{cases} 0 & y_{it} = 0 \\ 1 & y_{it} > 0 \end{cases}$$

به دو دلیل به دست آوردن برآورد ماکسیم درست‌نمایی برای لگاریتم تابع نمایی بالا سخت و پیچیده است و نمی‌توان لگاریتم تابع نمایی ماکسیم شود.

- دلیل اول: وجود جمله مجموع توابع نمایی (جمله اول در فرمول بالا) که ماکسیم کردن آن بسیار سخت و پیچیده است؛

- دلیل دوم: از آنجایی که در توزیع سری توانی آماسیده صفر، مقادیر صفری که پاسخ Y_{it} اختیار می‌کند از دو منبع می‌باشند، به طوری که یکی از آنها از توزیع سری توانی تولید شده و دیگری از هیچ توزیع خاصی تولید نمی‌شود بلکه همیشه مقدار آن صفر (صفر مطلق) بوده است، این موضوع سبب می‌شود که نتوان در به دست آوردن ماکسیم لگاریتم تابع

نمایی، این صفرها را تشخیص داد و این موضوع باعث پیچیدگی و سختی در ماکسیم کردن لگاریتم تابع درستنمایی می شود. برای رفع این مشکل از الگوریتم EM مشابه حالت مقطعی استفاده می شود. ابتدا متغیر تصادفی Z_{it} به این صورت تعریف می شود: $Z_{it} = 1$ وقتی که Y_{it} صفر (صفر مطلق از توزیع آماسیده صفر) است و $Z_{it} = 0$ وقتی که Y_{it} از توزیع سری توانی تولید شده باشد. لازم به ذکر است که بردار مربوط به متغیر تصادفی $(Z_{1t}, Z_{2t}, \dots, Z_{nt})'$ ، به عنوان متغیر پنهان و گم شده در نظر گرفته می شود و باید توسط الگوریتم EM ابتدا این بردار به دست آید و سپس با مشخص شدن منبع صفرهای موجود در تابع درستنمایی، لگاریتم تابع درستنمایی به راحتی ماکسیم می شود. بنابراین تابع درستنمایی برای داده های کامل $(\mathbf{y}, \mathbf{z}, \mathbf{b})$ به این صورت خواهد بود:

$$\begin{aligned} \text{Log}L_c(\psi; \mathbf{y}, \mathbf{z}, \mathbf{b}) &= \text{logf}(\mathbf{b}; \psi) + \text{log}(\mathbf{y}, \mathbf{z} | \mathbf{b}; \psi) \\ &= \sum_{i=1}^k \text{log}\phi(b_i) + \sum_{i=1}^k \sum_{t=1}^{T_i} \{ [z_{it} \mathbf{G}_{it} \gamma - \text{log}(1 + e^{\mathbf{G}_{it} \gamma})] \\ &+ (1 - z_{it}) \times [y_{it} (\mathbf{B}_{it} \beta + \sigma b_i) - \text{logc}(\exp(\mathbf{B}_{it} \beta + \sigma b_i)) + \text{log}(a(y_{it}))] \} \end{aligned}$$

که لگاریتم تابع درستنمایی به دو قسمت مجزا تقسیم می شود که قسمت اول تابعی از δ و قسمت دوم تابعی از β می باشند. لگاریتم تابع درستنمایی در این حالت ساده تر است.

۲-۳-۵. مدل های شمارشی برای تعداد ادعای خسارت طولی در بیمه شخص ثالث در این بخش مدل های مطرح شده، روی داده های مربوط به تعداد ادعای خسارت در بیمه شخص ثالث مورد بررسی قرار می گیرد. در این مدل ها تعداد ادعای خسارت بیمه گذار t ام در زمان t ام در داده های صنعت بیمه کشور ایران برای سال های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ می باشد. همچنین متغیرهای تبیینی مورد علاقه عبارت اند از:

سن اتومبیل (x_{i1t}) ، نوع اتومبیل (x_{i2}) و محل رانندگی (x_{i3}) که در این متغیرها مدل های برازش شده عبارت اند از:

مدل اول: مدل رگرسیون پواسون

$$Y_{it} | b_i \sim \text{Poisson}(\lambda_{it})$$

$$\text{log}\lambda_{it} = \beta_0 + \beta_{it} x_{i1t} + \beta_2 x_{i2} + \beta_3 x_{i3} + \sigma b_i$$

- مدل دوم: مدل رگرسیون دوجمله‌ای منفی

$$Y_{it} | b_i \sim NB(k, p_{it}) \quad p_{it} = \frac{k}{k + \lambda_{it}}$$

$$\log \lambda_{it} = \beta_0 + \beta_{it} x_{i1t} + \beta_2 x_{i2} + \beta_3 x_{i3} + \sigma b_i$$

- مدل سوم: مدل رگرسیون پواسون آماسیده صفر

$$Y_{it} | b_i \sim ZIP(\phi_{it}, \lambda_{it})$$

$$\log \lambda_{it} = \beta_0 + \beta_{it} x_{i1t} + \beta_2 x_{i2} + \beta_3 x_{i3} + \sigma b_i$$

$$\log it \frac{\phi_{it}}{1 - \phi_{it}} = \gamma_0 + \gamma_{1t} x_{i1t} + \gamma_2 x_{i2} + \gamma_3 x_{i3}$$

- مدل چهارم: مدل رگرسیون دوجمله‌ای منفی آماسیده صفر

$$Y_i | b_i \sim ZINB(k, p_{it}, \phi_{it}) \quad p_{it} = \frac{k}{k + \lambda_{it}}$$

$$\log \lambda_{it} = \beta_0 + \beta_{it} x_{i1t} + \beta_2 x_{i2} + \beta_3 x_{i3} + \sigma b_i$$

$$\log it \frac{\phi_{it}}{1 - \phi_{it}} = \gamma_0 + \gamma_{1t} x_{i1t} + \gamma_2 x_{i2} + \gamma_3 x_{i3}$$

در این مدل‌ها b_i اثر تصادفی با توزیع گاما برای توزیع شرطی پواسون و بتا برای دوجمله‌ای منفی در نظر گرفته شده است که در آن پاسخ‌های Y_{it} به شرط اثر تصادفی b_i ، از هم مستقل در نظر گرفته می‌شوند و σ پارامتر مربوط به این اثر تصادفی است. همچنین عوامل تشکیل دهنده خطر که قابل مشاهده و اندازه‌گیری نمی‌باشند در b_i قرار می‌گیرند. نحوه برآورد پارامترها با استفاده از نرم‌افزار R و بسته کامپیوتری VGAM است.

۳. روش تحقیق

۳-۱. فرضیه‌های تحقیق

تعداد ادعای خسارت در داده‌های بیمه شخص ثالث، گروه سن اتومبیل برای سال ۱۳۸۸ در ۱۱ رده، (۰-۱)، (۱-۲)، ...، (۹-۱۰) و بیشتر از (۱۰)، برای سال ۱۳۸۹ در ۱۰ رده (۰-۱)، (۱-۲)، (۲-۳)،

۲،...، ۱۰-۹ و بیشتر از ۱۰) و برای سال ۱۳۹۰ در ۹ رده (۳-۲، ۴-۳،...، ۱۰-۹ و بیشتر از ۱۰)، نوع اتومبیل در ۳ رده (پژو، پراید، پیکان) و محل رانندگی در ۳ رده (کم ترافیک، ترافیک متوسط، پر ترافیک) به عنوان متغیرهای مورد علاقه، در این مقاله مورد بررسی قرار گرفته‌اند. همچنین تعداد ادعای خسارت به عنوان یک متغیر هدف شمارشی در نظر گرفته می‌شود، به طوری که دارای توزیع احتمالی پواسون است. بر اساس این توزیع احتمالی می‌توان مدل‌های رگرسیونی ارائه نمود که ارتباط بین تعداد ادعای خسارت با متغیرهای گروه سن اتومبیل، نوع اتومبیل و محل رانندگی مورد بررسی قرار می‌گیرد.

۲-۳. سؤالات تحقیق

سؤالات تحقیق را می‌توان به این شرح خلاصه نمود:

آیا می‌توان تعداد ادعاهای خسارت در طول چند سال آینده را با استفاده از مدل‌های آماسیده صفر پیش‌بینی نمود؟

- مهم‌ترین عوامل مؤثر روی تعداد متوسط ادعای خسارت در داده‌های بیمه شخص ثالث، چه عواملی است؟

- مسئله وجود تعداد صفر زیاد در داده‌های مربوط به تعداد ادعای خسارت بیمه شخص ثالث چیست؟ چه مشکلاتی در تجزیه و تحلیل این داده‌ها ایجاد می‌کند؟ تحلیل این داده‌ها به چه صورت انجام می‌شود؟ همچنین مهم‌ترین عوامل مؤثر روی تعداد متوسط ادعای خسارت با تعداد صفر زیاد در داده‌های بیمه شخص ثالث در طول سال‌های مختلف، چه عواملی است؟

- بررسی مسئله بیش پراکنش روی تعداد ادعای خسارت چیست؟ چه اثری روی تعداد ادعای خسارت دارد؟ تحلیل این داده‌ها به چه صورت انجام می‌شود؟

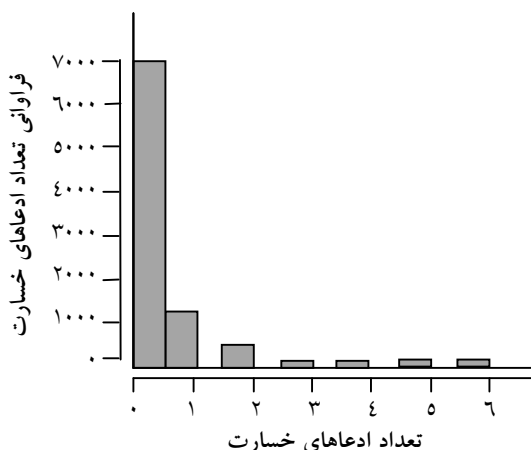
- چه مدل‌های آماری مناسبی برای متوسط تعداد ادعای خسارت با مسئله بیش پراکنش و تعداد صفر زیاد در داده‌های بیمه شخص ثالث در طول سال‌های مختلف، می‌توان ارائه نمود؟

۴. جمع‌آوری داده‌ها و اطلاعات

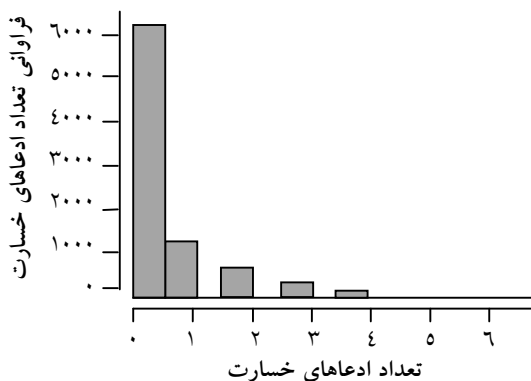
داده‌های مربوط به بیمه اتومبیل بیمه مرکزی ج.ا.ا در سال ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ در نظر گرفته می‌شود که ۸۵۶۶ فرد را پوشش می‌دهد. متغیر مورد علاقه تعداد مراجعه‌های یک فرد به شرکت بیمه برای گزارش خسارت در طی سه سال و متغیرهای تبیینی نیز مانند:

گروه سن اتومبیل برای سال ۱۳۸۸ در ۱۱ رده، (۰-۱، ۱-۲، ۲-۳، ۳-۴، ۴-۵، ۵-۶، ۶-۷، ۷-۸، ۸-۹، ۹-۱۰) و بیشتر از ۱۰، برای سال ۱۳۸۹ در ۱۰ رده (۰-۱، ۱-۲، ۲-۳، ۳-۴، ۴-۵، ۵-۶، ۶-۷، ۷-۸، ۸-۹، ۹-۱۰) و بیشتر از ۱۰ و برای سال ۱۳۹۰ در ۹ رده (۰-۱، ۱-۲، ۲-۳، ۳-۴، ۴-۵، ۵-۶، ۶-۷، ۷-۸، ۸-۹) و بیشتر از ۱۰، نوع اتومبیل در ۳ رده (پژو، پراید و پیکان) و محل رانندگی در ۳ رده (کم‌ترافیک، ترافیک متوسط و پرترافیک) مورد بررسی قرار گرفته‌اند. در نمودارهای ۱، ۲ و ۳، نمودار میله‌ای فراوانی ادعاهای خسارت برای سال‌های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ نشان داده شده است.

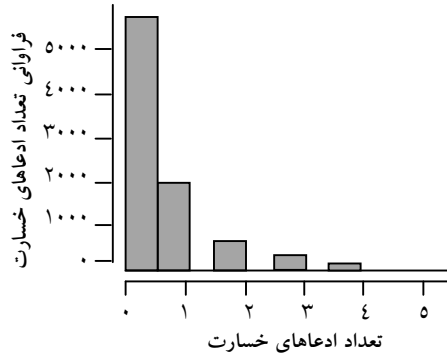
نمودار ۱. نمودار میله‌ای فراوانی ادعای خسارت برای سال ۱۳۸۸



نمودار ۲. نمودار میله‌ای فراوانی ادعای خسارت برای سال ۱۳۸۹



نمودار ۳. نمودار میله‌ای فراوانی ادعای خسارت برای سال ۱۳۹۰



در نمودارهای ۱، ۲ و ۳ ملاحظه می‌شود که فراوانی تعداد ادعای خسارت دارای صفر زیاد می‌باشند، همچنین با در نظر گرفتن جدول ۱ نیز اطلاعات مربوط به فراوانی تعداد ادعای خسارت برای سال‌های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ ارائه گردیده است.

جدول ۱. فراوانی ادعای خسارت برای سال‌های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰

فراوانی ادعاها در سال ۱۳۹۰	فراوانی ادعاها در سال ۱۳۸۹	فراوانی ادعاها در سال ۱۳۸۸	تعداد تصادفات
۶۸۴۳	۶۹۱۳	۷۰۲۰	۰
۱۱۶۳	۱۱۷۹	۱۲۹۹	۱
۲۲۸	۲۹۳	۱۹۵	۲
۱۸۸	۱۰۹	۲۸	۳
۴۰	۶۲	۱۶	۴
۳	۷	۵	۵
۲	۳	۲	۶
۸۴۶۷	۸۵۶۶	۸۵۶۵	مجموع

همانطور که ملاحظه می‌شود در سال ۱۳۸۸ درصد کسانی که هیچ ادعای خسارتی نداشته‌اند ۸۲٪، کسانی که یک ادعای خسارت داشته‌اند ۱۵٪، کسانی که دو ادعای خسارت داشته‌اند ۳/۲٪، کسانی که سه ادعای خسارت داشته‌اند ۴/۰٪ و کسانی که چهار

و بیشتر از چهار ادعای خسارت داشته‌اند $0/3$ ٪ است. در سال ۱۳۸۹ درصد کسانی که هیچ ادعای خسارتی نداشته‌اند $79/7$ ٪، کسانی که یک ادعا داشته‌اند $13/7$ ٪، کسانی که دو ادعای خسارت داشته‌اند $3/4$ ٪، کسانی که سه ادعای خسارت داشته‌اند $2/4$ ٪ و کسانی که چهار و بیشتر از چهار ادعای خسارت داشته‌اند $0/8$ ٪ است. در سال ۱۳۹۰ نیز درصد کسانی که هیچ ادعای خسارتی نداشته‌اند $78/5$ ٪، کسانی که یک ادعا داشته‌اند $15/5$ ٪، کسانی که دو ادعای خسارت داشته‌اند $3/3$ ٪، کسانی که سه ادعای خسارت داشته‌اند $2/2$ ٪ و کسانی که چهار و بیشتر از چهار ادعای خسارت داشته‌اند $0/5$ ٪ است.

۵. تجزیه و تحلیل اطلاعات

در جدول ۲ نیز بر اساس لگاریتم متوسط ادعای خسارت سن اتومبیل برای سال‌های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ در سطح معنی‌داری $0/1$ در مدل‌های رگرسیونی پواسون آماسیده صفر و دوجمله‌ای منفی آماسیده صفر معنی‌دار است. هر چقدر سن اتومبیل بیشتر باشد، لگاریتم متوسط خسارت بیشتر خواهد بود. یعنی سن اتومبیل با متوسط خسارت، رابطه لگاریتمی دارد و هرچقدر سن اتومبیل بیشتر باشد متوسط تعداد خسارت‌ها بیشتر خواهد بود. همچنین نوع اتومبیل نیز در هیچ‌کدام از مدل‌ها معنی‌دار نشده است. ولی محل رانندگی در سطح معنی‌داری $0/1$ در مدل رگرسیونی پواسون آماسیده صفر و دوجمله‌ای منفی آماسیده صفر معنی‌دار است. هرچقدر محل رانندگی، در جاهای کم ترافیک باشد لگاریتم متوسط خسارت کمتر خواهد بود. اثر تصادفی نیز در سطح معنی‌داری $0/05$ در مدل رگرسیونی پواسون آماسیده صفر و دوجمله‌ای منفی آماسیده صفر معنی‌دار می‌باشد که نشان‌دهنده این است که عواملی در مدل وجود داشته که تحت کنترل و اندازه‌گیری ما نبوده است. بر اساس لوجیت نسبت صفرهای ساختاری نیز سن اتومبیل در سال‌های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ و نوع اتومبیل در هیچ سطحی معنی‌دار نبوده و بنابراین در هیچ‌کدام از مدل‌ها اثرگذار نمی‌باشند. اما محل رانندگی در سطح $0/05$ در مدل رگرسیونی پواسون آماسیده صفر و دوجمله‌ای منفی آماسیده صفر معنی‌دار است. هر چقدر محل رانندگی در جاهای کم ترافیک باشد لوجیت نسبت صفرهای ساختاری کمتر

خواهد بود. این موضوع به ما نشان می‌دهد که سن اتومبیل و نوع آن در صفرهای ساختاری که به واسطه وقوع خسارت و گزارش نکردن آن به وجود آمده است، تأثیری ندارد و عامل مؤثر در به وجود آمدن این نوع صفرها محل رانندگی بوده است. البته ممکن است متغیرهای تبیینی دیگری نیز در به وجود آمدن نسبت صفرهای ساختاری نقش داشته باشند ولی در این مقاله فقط اثر این سه متغیر مورد بررسی قرار گرفته است. نتایج به دست آمده به نحوی توانایی مدل‌های آماسیده صفر را نیز نشان می‌دهد. در مدل رگرسیونی دوجمله‌ای و دوجمله‌ای منفی آماسیده صفر به دلیل اینکه پارامتر پراکنش در سطح معنی داری ۰/۰۵ معنی دار است و همچنین به دلیل وجود صفرهای ساختاری زیاد براساس معیارهای AIC و BIC این مدل به عنوان بهترین مدل انتخاب می‌شود؛ یعنی اینکه برازش بهتری نسبت به مدل‌های دیگر به داده‌های تعداد ادعاهای خسارت دارد.

جدول ۲. نتایج برازش مدل‌های رگرسیونی به داده‌ها با پاسخ‌های طولی

Poisson	NB	ZIP	ZINB	پارامتر	مدل
-۱/۵۶۲۹ (۰/۰۸۹۱)***	-۱/۵۶۳۱ (۰/۰۵۲۳)***	-۰/۷۸۶۲ (۰/۰۸۲۶)***	-۱/۵۷۰۱ (۰/۱۰۰۶)***		عرض از مبدأ
۰/۰۰۴۵ (۰/۰۰۳۵)	۰/۰۰۵۷ (۰/۰۰۴۸)	۰/۰۱۴۳ (۰/۰۰۷۳)*	۰/۰۱۱۴ (۰/۰۰۶۸)*		۱۳۸۸ سن اتومبیل (بیشتر از ۱۰ سال)
۰/۰۰۵۸ (۰/۰۰۴۱)	۰/۰۰۶۹ (۰/۰۰۴۶)	۰/۰۱۸۳ (۰/۰۰۹۴)*	۰/۰۱۲۷ (۰/۰۰۶۶)*		۱۳۸۹ سن اتومبیل (بیشتر از ۱۰ سال)
۰/۰۰۵۱ (۰/۰۰۴۹)	۰/۰۰۷۲ (۰/۰۰۵۱)	۰/۰۱۷۱ (۰/۰۰۸۷)*	۰/۰۱۳۳ (۰/۰۰۷۰)*		۱۳۹۰ سن اتومبیل (بیشتر از ۱۰ سال)
-۰/۰۳۷۷ (۰/۰۲۸۱)	-۰/۰۳۸۰ (۰/۰۳۱۱)	-۰/۰۲۱۱ (۰/۰۶۲۸)	-۰/۰۰۵۶ (۰/۰۴۲۵)		نوع اتومبیل (پژو)
-۰/۰۳۳۸ (۰/۰۲۸۴)	-۰/۰۳۳۸ (۰/۰۳۱۴)	-۰/۱۱۲۷ (۰/۰۶۰۴)*	-۰/۰۱۱۲ (۰/۰۰۵۹)*		محل رانندگی (کم ترافیک)
۰/۰۰۳۳ (۰/۰۰۲۶)	۰/۰۰۲۴ (۰/۰۰۱۹)	۰/۰۰۳۲ (۰/۰۰۱۵)**	۰/۰۱۴۲ (۰/۰۰۶۶)**	δ	اثر تصادفی
-	-	-۱/۱۲۳۲ (۰/۳۰۵۵)***	-۰/۵۴۲۳ (۰/۱۳۲۵)***		عرض از مبدأ
-	-	۰/۰۱۷۳ (۰/۰۱۶۴)	۰/۱۷۰۶ (۰/۲۱۳۳)		۱۳۸۸ سن اتومبیل (بیشتر از ۱۰ سال)
-	-	۰/۰۱۳۹ (۰/۰۱۰۶)	۰/۱۶۸۶ (۰/۱۰۱۳)		۱۳۸۹ سن اتومبیل (بیشتر از ۱۰ سال)
-	-	۰/۰۱۲۲ (۰/۰۱۷۳)	۰/۱۲۱۶ (۰/۱۱۶۲)		۱۳۹۰ سن اتومبیل (بیشتر از ۱۰ سال)
-	-	۰/۰۳۴ (۰/۱۳۰۲)	۰/۷۹۹۴ (۰/۶۳۴۳)		نوع اتومبیل (پژو)
-	-	-۰/۳۳۸۲ (۰/۱۲۹۹)**	-۰/۹۸۳ (۰/۴۰۱۳)**		محل رانندگی (کم ترافیک)
-	۱/۱۳۳۴ (۰/۰۵۶۲)**	-	۱/۰۲۵۴ (۰/۰۱۱۴)**		پارامتر پراکنش
۵۰۰۱/۷۶	۴۹۱۷/۲۵	۴۹۳۵/۶۱	۴۹۱۱/۵۸		منفی لگاریتم نمایی
۱۰۰۱۵/۵۲	۹۸۴۸/۵۶	۹۸۹۳/۱۲	۹۸۴۷/۱۶		معیار AIC
۱۰۰۶۳/۲۱	۹۸۸۴/۳۲	۹۹۴۱/۱۱	۹۸۹۲/۴۱		معیار BIC

* معنی داری در سطح ۰/۱، ** معنی داری در سطح ۰/۰۵، *** معنی داری در سطح ۰/۰۰۱

عبارت داخل پرانتز انحراف معیار برآورد پارامترهاست.

۶. بحث و نتیجه‌گیری

در این تحقیق به بررسی و مقایسه مدل‌های شمارشی طولی آماسیده صفر برای پاسخ‌های مربوط به تعداد ادعای خسارت در بیمه شخص ثالث پرداخته شد. این مدل‌ها نقش بسیار زیادی در تعیین عوامل تشکیل‌دهنده خطر روی تعداد ادعای خسارت ایفا می‌کنند و بر اساس این مدل‌ها، به شرکت‌های بیمه و آگاهی کافی در مورد عوامل تشکیل‌دهنده خطر می‌دهند به طوری که برای چند سال آینده تعداد ادعای خسارت مربوط به بیمه شخص ثالث قابل پیش‌بینی است. در این داده‌ها، سن اتومبیل در طول سال‌های ۱۳۸۸، ۱۳۸۹ و ۱۳۹۰ و محل رانندگی به عنوان عوامل مؤثر و تشکیل‌دهنده خطر روی تعداد ادعای خسارت در این سه سال مطرح شده است. از این مدل‌ها می‌توان برای داده‌های بیمه بدنه اتومبیل نیز استفاده کرد.

منابع

1. Boucher, J.P., Denuit, M. and Guille'n, M., 2008. Models of insurance claim counts with time dependence based on generalisation of poisson and negative binomial distributions. *Journal of the Variance*, 2(1), pp. 135–162.
2. Boucher, J.P. and Guille'n, M., 2009. A survey on models for panel count data with applications to insurance. *Journal of the Applied Mathematics*, 3(2), pp. 280-281.
3. Dempster, A.P., Larid, N.M. and Rubin, D.B., 1977. Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Statist*, 39, pp. 1-38.
4. Hall, D.B., 2000. Zero-Inflated poisson and binomial regression with random effects: A case study. *Journal of the Biometrics*, 56, pp. 1030-1039.
5. Hausman, J., Hall, B. and Griliches, Z., 1984. Econometric Models for Count Data with Application to the Patents-R and D Relationship. *Econometrica*, 52, pp. 909–938.
6. Heilbron, D., 1994. Zero-Altered and other regression models for count data with added zeros. *Journal of the Biometrical*, 36, pp. 531–547.
7. Lambert, D., 1992. Zero-Inflated poisson regression, with an application to defects in manufacturing. *Journal of the Technometrics*, 34, pp. 1-14.
8. Thomas, H. and Samson, D., 1987. Linear models as aids in insurance decision making: The estimation of automobile insurance claims. *Journal of the Business*, 15, pp. 247–256.
9. Yip, K.C.H. and Yau, K.K.W., 2005. On modeling claim frequency data in general insurance with extra zeros. *Journal of the Mathematics and Economics*, 36, pp. 153-163.

