

مدل‌های رگرسیونی هاردل، پوآسون آماسیده و دوجمله‌ای منفی آماسیده برای تحلیل داده‌های شمارشی با صفر زیاد

احسان بهرامی سامانی^۱

تاریخ دریافت: ۱۳۹۶/۷/۲۳

تاریخ پذیرش: ۱۳۹۷/۶/۲۶

چکیده:

در این مقاله، مدل رگرسیونی هاردل برای پاسخ‌های شمارشی با تعداد صفر زیاد معرفی می‌شود. یک شیوه برآورد ماکسیمم درست‌نمایی برای برآورد پارامترهای مدل استفاده شده است. کاربردی از مدل معرفی شده در داده‌های بیم‌سنجی ارائه شده است. در این مثال، تعداد زیادی ادعای خسارت برابر صفر وجود دارد که کاربرد مدل با پاسخ شمارشی آماسیده صفر را روشن می‌سازد. مدل‌های رگرسیونی شمارشی مختلفی برای مدل‌بندی چنین پاسخ‌های شمارشی در این مقاله معرفی شده‌اند. از جمله این مدل‌های معرفی شده می‌توان به مدل رگرسیونی پوآسون هاردل و مدل رگرسیونی دوجمله‌ای منفی هاردل اشاره کرد.

واژه‌های کلیدی: توزیع آماسیده صفر، پاسخ‌های سری توانی، مدل رگرسیونی هاردل، تعداد ادعای خسارت.

۱ مقدمه

داده و همچنین مدل دوجمله‌ای آماسیده صفر و مدل رگرسیونی آماسیده صفر با پاسخ‌های طولی را مطرح کرد. مدل هاردل برای داده‌های شمارشی آماسیده صفر در [۵] مطالعه شده است. [۸] برای مدل‌سازی پاسخ‌های شمارشی از مدل رگرسیونی پوآسون و مدل رگرسیونی دوجمله‌ای منفی برای پاسخ‌های مقطعی استفاده کردند. [۶] با استفاده از توزیع‌های چندمتغیره طولی پوآسون و دوجمله‌ای منفی به محاسبه حق بیمه پرداختتند. همچنین [۷] برای تعیین حق بیمه‌ها در بیمه اتومبیل در طول چند سال از توزیع‌های آماسیده صفر استفاده کردند.

در بسیاری از مطالعه‌ها ممکن است با مجموعه‌ای از داده‌های شمارشی روبرو شویم. به‌عنوان مثال در مطالعه‌های مربوط به علوم بیم‌سنجی، تعداد ادعای خسارت مربوط به بیمه شخص ثالث در یک سال مشخص به‌عنوان یک پاسخ شمارشی از اهمیت زیادی برخوردار است، به‌طوری که این پاسخ ممکن است دارای تعداد صفر زیاد باشد. از سویی دیگر بررسی و تحلیل این پاسخ از اهمیت بسیاری در جهت شناسایی عوامل تشکیل‌دهنده خطر دارا است. همچنین به‌واسطه ساختار پیچیده این داده‌ها، تحلیل و مدل‌سازی آنها کار بسیار پیچیده‌ای است که این موضوع سبب شده است بسیاری از محققان به تحلیل داده‌های گسسته روی آورند که برخی از این تحقیقات در ادامه بیان می‌شود.

در این مقاله، خانواده‌ای از مدل‌های آماسیده و مدل‌های هاردل تحت عنوان مدل‌های رگرسیونی آماسیده صفر و هاردل برای معرفی می‌شود. از آنجایی که توزیع‌های دوجمله‌ای، پوآسون و دوجمله‌ای منفی جز خانواده توزیع‌های سری توانی هستند، خانواده‌ای از مدل‌هایی مانند مدل‌های آماسیده صفر پوآسون و پوآسون هاردل و همچنین مدل‌های آماسیده صفر دوجمله‌ای منفی و دوجمله‌ای منفی هاردل معرفی می‌شوند و برای تشریح توانمندی آنها یک مطالعه شبیه‌سازی انجام می‌گیرد و در

مدل‌های شمارشی آماسیده صفر روشی را برای مدل‌بندی جرم احتمالی صفر در توزیع متغیرها فرض می‌کند. مدل‌بندی صفرها برای داده‌های شمارشی در حالت تک متغیره مطالعه شده است [۳، ۱]. همچنین رگرسیون پوآسون آماسیده صفر چندسطحی ارائه شده است [۴]. [۲] روش مطرح شده در [۳] را تعدیل

^۱ عضو هیئت علمی گروه آمار دانشگاه شهید بهشتی تهران، تهران، ایران

۲.۲ توزیع هاردل

تابع جرم احتمال این توزیع به صورت زیر است:

$$f_{ZIF}(y) = \pi I(y = 0) + (1 - \pi) \frac{f(y)}{1 - f(0)} I(y > 0).$$

که در آن پارامتر آمیختگی π است، زمانی که متغیر پاسخ مقدار صفر را اختیار کند. همچنین f یک تابع جرم احتمال سری توانی است، زمانی که متغیر پاسخ مقدار صحیح و نامنفی اختیار کند.

۳.۲ توزیع سری توانی آماسیده صفر

فرض شود که Y یک متغیر تصادفی گسسته شمارشی با تعداد صفر زیاد باشد. توزیع سری توانی آماسیده صفر را با نماد $Y \sim ZIPS(\pi, \theta)$ نشان داده و تابع جرم احتمال آن به صورت زیر است:

$$f_{ZIPS}(y) = \pi I(y = 0) + (1 - \pi) f_{PS}(y|\theta) I(y > 0),$$

که در آن

$$f_{PS}(y|\theta) = \frac{a(y)\theta^y}{c(\theta)}, \quad y = 0, 1, 2, \dots$$

یا به عبارت دیگر:

$$p(Y=y) = \begin{cases} \pi + (1-\pi) \frac{a(0)}{c(\theta)}, & y = 0, \\ (1-\pi) \frac{a(y)\theta^y}{c(\theta)}, & y > 0 \end{cases}$$

$a(\cdot)$ که در آن تابعی مثبت، $c(\cdot)$ تابعی مثبت، متناهی و مشتق پذیر از θ است.

۴.۲ توزیع سری توانی هاردل

فرض شود که Y یک متغیر تصادفی گسسته شمارشی با تعداد صفر زیاد باشد. توزیع سری توانی هاردل را با نماد $Y \sim HPS(\pi, \theta)$ نشان داده و تابع جرم احتمال آن به صورت زیر است

$$f_{HPS}(y) = \pi I(y = 0) + I(y > 0) (1 - \pi) f_{PS}(y|\theta) / f_{PS}(0|\theta).$$

یا به عبارت دیگر

$$p(Y=y) = \begin{cases} \pi, & y = 0, \\ (1-\pi) \frac{a(y)\theta^y}{a(0)}, & y > 0 \end{cases}$$

نهایت این مدل‌ها روی داده‌های بیم‌سنجی برازش داده شده و نتایج مهمی استخراج می‌گردد.

در بخش دوم به معرفی توزیع سری توانی آماسیده صفر و سری توانی هاردل پرداخته می‌شود. در بخش سوم مدل رگرسیونی مربوط به توزیع‌های بیان شده، معرفی می‌گردد و توابع درست‌نمایی آنها معرفی می‌شود. در بخش چهارم مطالعه شبیه‌سازی به منظور بررسی عملکرد مدل‌های رگرسیونی بیان شده، انجام می‌شود. در بخش پنجم، کاربرد مدل‌های بیان شده روی داده‌های بیمه بیان می‌شود و در نهایت در بخش ششم به نتیجه‌گیری پرداخته می‌شود.

۲ برخی توزیع‌های آماری آماسیده صفر

در تحلیل داده‌های شمارشی

در این بخش با مدل آماسیده صفر و مدل هاردل آشنا شده سپس توزیع سری توانی آماسیده صفر معرفی می‌شود.

۱.۲ توزیع آماسیده صفر

یکی از توزیع‌های قابل استفاده برای متغیرهای گسسته شمارشی با تعداد صفر زیاد، توزیع آماسیده صفر (ZI) است. توزیع آماسیده صفر به دلیل زیاد بودن تعداد صفر در مقادیر متغیر گسسته شمارشی یک پارامتر آمیختگی را برای مقدار صفر در نظر می‌گیرند. فرض می‌شود که Y یک متغیر تصادفی گسسته شمارشی با تعداد صفر زیاد باشد. توزیع آماسیده صفر را با نماد $Y \sim ZIF(\pi_0, \theta)$ نشان داده و تابع جرم احتمال آن به صورت زیر است:

$$f_{ZIF}(y) = \pi_0 I(y = 0) + (1 - \pi_0) f_F(y|\theta) I(y > 0),$$

که در آن $f_F(y|\theta)$ تابع احتمال توزیع F با پارامتر θ ، $f_{ZIF}(y)$ توزیع آماسیده صفر از توزیع F با یک پارامتر آمیختگی π_0 و $I(y = 0)$ تابع نشانگر جرم در صفر است.

گسسته شمارشی y و یک احتمال $(1 - \pi)$ را به تابع جرم احتمال پوآسون با پارامتر μ برای مقادیر مثبت پاسخ متغیر گسسته y نسبت می‌دهد که این تابع جرم احتمال به صورت زیر است:

$$Pr(Y = y) = \begin{cases} \pi + (1 - \pi)e^{-\mu}, & y = 0, \\ (1 - \pi)\frac{e^{-\mu}\mu^y}{y!}, & y = 1, 2, \dots \end{cases}$$

تابع درست‌نمایی برای متغیر پاسخ گسسته شمارشی Y که دارای توزیع ZIP است، به صورت زیر محاسبه می‌شود.

فرض کنید Y_1, \dots, Y_n یک نمونه تصادفی از توزیع $ZIP(\pi, \mu)$ باشد، مدل رگرسیون پوآسون آماسیده صفر به صورت است

$$\text{logit}(\pi_i) = B_i' \beta,$$

$$\log(\mu_i) = W_i' \gamma,$$

که در آن B_i و W_i سطر i -ام ماتریس‌های طرح (بردارهای متغیرهای کمکی) و همچنین β و γ بردار پارامترهای رگرسیونی مدل هستند که باید برآورد شوند. از سویی دیگر برای برآورد پارامترها تابع درست‌نمایی این مدل به صورت زیر است:

$$L(\mu, \pi | y) = \prod_{i=1}^n [(\pi_i + (1 - \pi_i)e^{-\mu_i}) I(y_i = 0) \\ ((1 - \pi_i)\frac{e^{-\mu_i}\mu_i^{y_i}}{y_i!}) I(y_i \geq 1)].$$

۳.۳ مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر

فرض کنید Y_1, \dots, Y_n یک نمونه تصادفی از توزیع $ZINB(\pi, \mu, \frac{\delta}{1+\delta})$ باشد، لذا تابع جرم احتمال آن به صورت زیر در نظر گرفته می‌شود:

$$p(Y = y) = \begin{cases} \pi + (1 - \pi)(\frac{\delta}{1+\delta})^\mu, & y = 0 \\ (1 - \pi)\frac{\Gamma(y + \mu)}{y! \Gamma(\mu)} (\frac{\delta}{1+\delta})^\mu (\frac{1}{1+\delta})^y, & y = 1, 2, \dots \end{cases}$$

که در آن δ پارامتر پراکنش نامیده می‌شود.

مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر به صورت زیر است:

$$\text{logit}(\pi_i) = B_i' \beta,$$

$$\log(\mu_i) = W_i' \gamma,$$

۳ مدل‌های رگرسیونی با پاسخ آماسیده صفر

در این بخش به معرفی مدل‌های رگرسیونی با پاسخ آماسیده صفر پرداخته می‌شود، مدل‌هایی مانند مدل رگرسیونی سری توانی آماسیده صفر و مدل رگرسیونی هاردل از جمله این مدل‌ها هستند.

۱.۳ مدل رگرسیونی سری توانی آماسیده صفر

فرض کنید Y_1, \dots, Y_n یک نمونه تصادفی از توزیع $ZIPS(\pi_i, \theta_i)$ در نظر گرفته شود، از سویی دیگر با در نظر گرفتن بردار پارامترهای $\theta = (\theta_1, \dots, \theta_n)'$ و $\pi = (\pi_1, \dots, \pi_n)'$ مدل رگرسیونی سری توانی آماسیده صفر به صورت زیر معرفی می‌شود:

$$\text{logit}(\pi_i) = B_i' \beta,$$

$$\log(\theta_i) = W_i' \gamma,$$

که در آن B_i و W_i سطر i -ام ماتریس‌های طرح (بردارهای متغیرهای کمکی) و همچنین β و γ بردار پارامترهای رگرسیونی مدل هستند که باید برآورد شوند. از سویی دیگر برای برآورد پارامترها تابع درست‌نمایی این مدل به صورت زیر است:

$$\log L(\beta, \gamma, y) = \sum_{i=1}^n (1 - I(y_i > 0)) \log[\pi_i + (1 - \pi_i)\frac{a(\circ)}{c(\theta_i)}] \\ + \sum_{i=1}^n I(y_i > 0) \log[(1 - \pi_i)\frac{a(y_i)\theta_i^{y_i}}{c(\theta_i)}], \\ = \sum_{i=1}^n (1 - I(y_i > 0)) \log[\frac{e^{B_i' \beta}}{1 + e^{B_i' \beta}} + \frac{1}{1 + e^{B_i' \beta}} \frac{a(\circ)}{c(e^{W_i' \gamma})}], \\ + \sum_{i=1}^n I(y_i > 0) \log[\frac{1}{1 + e^{B_i' \beta}} \frac{a(y_i)(e^{W_i' \gamma})^{y_i}}{c(e^{W_i' \gamma})}].$$

از رایج‌ترین این مدل‌ها می‌توان به مدل رگرسیونی پوآسون آماسیده صفر و مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر اشاره نمود.

۲.۳ مدل رگرسیونی پوآسون آماسیده صفر

وقتی تعداد صفر زیاد برای توزیع پوآسون رخ دهد، تابع جرم احتمال آن به صورت توزیع آمیخته در نظر گرفته می‌شود، این توزیع آمیخته یک پارامتر آمیختگی π را به مقدار صفر متغیر

هاردل ابتدا از این توزیع‌ها، داده تولید می‌شود. برای این منظور طی یک بررسی شبیه‌سازی برای اندازه‌های نمونه‌ای ۵۰، ۱۰۰ و ۱۰۰۰ به شیوه زیر، نمونه تصادفی تولید می‌شود و این کار را ۱۰۰۰ بار تکرار خواهد شد.

برای تولید داده، مراحل زیر به ترتیب انجام دهید:

الف) مرحله اول: ابتدا یک عدد تصادفی مانند U از توزیع یکنواخت روی بازه $(0,1)$ تولید کنید.

ب) مرحله دوم: متغیر تبیینی X را از توزیع نرمال استاندارد تولید کنید و مقایر واقعی پارامترهای مدل‌ها را به صورت (β_0, β_1) و $(\gamma_0, \gamma_1) = (0, 0.5)$ ، $(0, 0.5)$ و $\delta = 1$ در نظر بگیرید.

ج) مرحله سوم: متغیر پاسخ را بر اساس یکی از چهار مدل زیر تولید کنید:

ج-۱-مدل اول: مدل رگرسیونی پواسون آماسیده صفر:

$$Y \sim ZIP(\mu, \pi),$$

$$\mu = \exp(\gamma_0 + \gamma_1 X),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 X.$$

ج-۲-مدل دوم: مدل رگرسیونی پواسون هاردل:

$$Y \sim \text{HurdlePois}(\mu, \pi),$$

$$\mu = \exp(\gamma_0 + \gamma_1 X),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 X.$$

ج-۳-مدل سوم: مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر:

$$Y \sim ZINB\left(\pi, \mu, \frac{\delta}{1+\delta}\right),$$

$$\mu = \exp(\gamma_0 + \gamma_1 X),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 X.$$

ج-۴-مدل چهارم: مدل رگرسیونی دوجمله‌ای منفی هاردل:

$$Y \sim \text{HurdleNB}\left(\pi, \mu, \frac{\delta}{1+\delta}\right),$$

$$\mu = \exp(\gamma_0 + \gamma_1 X),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 X.$$

که در آن B_i و W_i سطر i -ام ماتریس‌های طرح (بردارهای متغیرهای کمکی) و همچنین β و γ بردار پارامترهای رگرسیونی مدل هستند که باید برآورد شوند. از سویی دیگر برای برآورد پارامترها تابع درست‌نمایی این مدل به صورت زیر است:

$$L(\mu, \pi | y) = \prod_{i=1}^n \left(\pi_i + (1-\pi_i) \left(\frac{\delta}{1+\delta} \right)^{\mu_i} \right) I(y_i = 0) \\ \times \left((1-\pi_i) \frac{(y_i + \mu_i)}{y_i!} \left(\frac{\delta}{1+\delta} \right)^{\mu_i} \left(\frac{1}{1+\delta} \right)^{y_i} \right) I(y_i \geq 1).$$

۴.۳ مدل رگرسیونی سری توانی هاردل

فرض کنید Y_1, \dots, Y_n یک نمونه تصادفی از توزیع $HPS(\pi_i, \theta_i)$ در نظر گرفته شود، از سویی دیگر با در نظر گرفتن بردار پارامترهای $\theta = (\theta_1, \dots, \theta_n)'$ و $\pi = (\pi_1, \dots, \pi_n)$ مدل رگرسیون سری‌های توانی هاردل را به صورت زیر است

$$\text{logit}(\pi_i) = B_i' \beta,$$

$$\log(\theta_i) = W_i' \gamma,$$

که در آن B_i و W_i سطر i -ام ماتریس‌های طرح (بردارهای متغیرهای کمکی) و همچنین β و γ بردار پارامترهای رگرسیونی مدل هستند که می‌بایست برآورد شوند. از سویی دیگر برای برآورد پارامترها تابع درست‌نمایی این مدل به صورت زیر است:

$$\log L(\beta, \gamma, y) = \sum_{i=1}^n I(y_i = 0) \log[\pi_i] \\ + \sum_{i=1}^n I(y_i > 0) \log\left[(1-\pi_i) \frac{a(y_i)\theta_i^{y_i}}{a(0)} \right], \\ = \sum_{i=1}^n I(y_i = 0) \log\left[\frac{e^{B_i' \beta}}{1+e^{B_i' \beta}} \right] \\ + \sum_{i=1}^n I(y_i > 0) \log\left[\frac{1}{1+e^{B_i' \beta}} \frac{a(y_i)(e^{W_i' \gamma})^{y_i}}{a(0)} \right].$$

۴ مطالعه شبیه‌سازی

در این بخش، طی چند مطالعه شبیه‌سازی به بررسی عملکرد مدل‌های بیان شده در بخش سوم پرداخته می‌شود. برای یک مطالعه شبیه‌سازی از مدل‌های رگرسیونی پواسون آماسیده صفر، مدل رگرسیونی پواسون هاردل، مدل‌های رگرسیونی دوجمله‌ای منفی آماسیده صفر و در نهایت مدل رگرسیونی دوجمله‌ای منفی

جدول ۱. نتایج شبیه‌سازی مربوط به برآورد پارامترهای چهار مدل بیان شده

پارامتر	مقدار واقعی	$n = 50$		$n = 100$		$n = 1000$	
		برآورد	انحراف معیار	برآورد	انحراف معیار	برآورد	انحراف خطا
<i>ZIP</i>							
β_0	۰/۰۰۰	۰/۰۷۲	۰/۸۷۱	۰/۰۶۴	۰/۷۷۲	۰/۰۰۸	۰/۰۳۴
β_1	۰/۵۰۰	۰/۴۳۰	۰/۷۸۷	۰/۵۱۰	۰/۶۲۴	۰/۵۰۳	۰/۰۲۵
γ_0	۰/۰۰۰	۰/۰۸۹	۰/۶۲۵	۰/۰۷۶	۰/۴۲۳	۰/۰۶۳	۰/۰۶۳
γ_1	۰/۵۰۰	۰/۴۰۰	۰/۶۲۳	۰/۵۳۰	۰/۵۶۶	۰/۴۹۰	۰/۰۴۵
<i>HurdlePois</i>							
β_0	۰/۰۰۰	۰/۰۲۵	۰/۲۱۰	۰/۰۱۹	۰/۲۰۰	۰/۰۰۸	۰/۰۴۲
β_1	۰/۵۰۰	۰/۴۰۱	۰/۱۹۹	۰/۴۸۵	۰/۱۸۷	۰/۴۹۸	۰/۰۵۶
γ_0	۰/۰۰۰	۰/۱۵۰	۰/۲۳۲	۰/۱۲۹	۰/۱۱۲	۰/۰۲۶	۰/۰۷۸
γ_1	۰/۵۰۰	۰/۴۸۰	۰/۲۶۵	۰/۴۹۲	۰/۱۰۵	۰/۴۹۵	۰/۰۴۵
<i>ZINB</i>							
β_0	۰/۰۰۰	-۰/۳۲۸	۰/۱۷۹	-۰/۰۲۵	۰/۰۹۵	-۰/۰۲۰	۰/۰۲۵
β_1	۰/۵۰۰	۰/۴۱۸	۰/۲۳۹	۰/۴۵۹	۰/۱۳۸	-۰/۰۱۷	۰/۰۳۵
γ_0	۰/۰۰۰	-۰/۲۵۵	۰/۱۸۵	-۰/۰۲۰	۰/۱۱۲	-۰/۰۰۰۲	۰/۰۴۵
γ_1	۰/۵۰۰	۰/۴۴۲	۰/۲۸۱	۰/۴۴۸	۰/۰۸۸	۰/۴۹۵	۰/۰۵۶
δ	۱/۰۰۰	۱/۲۵۰	۰/۲۱۶	۱/۱۲۵	۰/۱۰۲	۱/۰۲۶	۰/۰۵۵
<i>HurdleNB</i>							
β_0	۰/۰۰۰	۰/۱۴۹	۰/۲۱۶	۰/۰۲۶	۰/۱۲۳	۰/۰۲۷	۰/۰۷۹
β_1	۰/۵۰۰	۰/۴۶۴	۰/۲۲۱	۰/۴۸۰	۰/۰۸۷	۰/۴۹۹	۰/۰۳۱
γ_0	۰/۰۰۰	۰/۱۵۰	۰/۱۳۲	۰/۰۲۹	۰/۱۰۲	۰/۰۱۶	۰/۰۵۸
γ_1	۰/۵۰۰	۰/۴۸۲	۰/۱۶۴	۰/۴۹۶	۰/۱۰۵	۰/۴۹۸	۰/۰۴۵
δ	۱/۰۰۰	۰/۸۱۷	۰/۲۱۷	۰/۸۷۶	۰/۱۰۲	۰/۹۲۱	۰/۰۵۲

برآورد پارامترهای مدل در هر چهار مدل به مقدار واقعی پارامتر نزدیک می‌شود. از سوی دیگر برآورد پارامترهای این چهار مدل برآوردهای سازگاری هستند.

۵ تحلیل داده‌های بیم‌سنجی

تعداد ادعای خسارت در داده‌های بیمه شخص ثالث، گروه سن اتومبیل برای سال ۱۳۸۶ در ۲ رده، (۰-۷ و بیشتر از ۷)، نوع اتومبیل در ۳ رده (پژو، پراید، پیکان) و محل رانندگی در ۳ رده (کم‌ترافیک، ترافیک متوسط، پرترافیک) به‌عنوان متغیرهای

(د) مرحله چهارم: اگر $U < \pi$ ، آن‌گاه مقدار متغیر پاسخ Y صفر در نظر گرفته می‌شود، در غیر این صورت برای تولید پاسخ از توزیع پواسون برای مدل اول، پواسون بریده شده در صفر برای مدل دوم، توزیع دوجمله‌ای منفی برای مدل سوم و در نهایت توزیع دوجمله‌ای منفی بریده شده در صفر برای مدل چهارم استفاده می‌شود.

نتایج شبیه‌سازی مربوط به برآورد پارامترهای چهار مدل بیان شده در جدول ۱ گردآوری شده است. با توجه به نتایج جدول ۱، هر چهار مدل توانایی مناسبی در برآورد پارامترهای مربوط به مدل خود را دارند و با افزایش اندازه نمونه مقدار

۳-مدل سوم: در این مدل برای تعداد ادعای خسارت در بیمه شخص ثالث (پاسخ گسسته) مدل رگرسیونی دوجمله‌ای منفی آماسیده صفر در نظر گرفته می‌شود:

$$Y \sim ZINB\left(\pi, \mu, \frac{\delta}{1+\delta}\right),$$

$$\mu = \exp(\gamma_0 + \gamma_1 x_1 + \gamma_2 x_{11} + \gamma_3 x_{12} + \gamma_4 x_{21} + \gamma_5 x_{22}),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 x_1 + \beta_2 x_{11} + \beta_3 x_{12} + \beta_4 x_{21} + \beta_5 x_{22}.$$

۴-مدل چهارم: در این مدل فرض می‌شود برای تعداد ادعای خسارت در بیمه اتومبیل (پاسخ گسسته) مدل رگرسیونی هاردل دوجمله‌ای منفی در نظر گرفته می‌شود:

$$Y \sim \text{HurdleNB}\left(\pi, \mu, \frac{\delta}{1+\delta}\right),$$

$$\mu = \exp(\gamma_0 + \gamma_1 x_1 + \gamma_2 x_{11} + \gamma_3 x_{12} + \gamma_4 x_{21} + \gamma_5 x_{22}),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 x_1 + \beta_2 x_{11} + \beta_3 x_{12} + \beta_4 x_{21} + \beta_5 x_{22}.$$

۲.۵ نتایج

نتایج برآورد پارامترهای مدل‌های ZIP و HurdlePois در جدول ۴ گردآوری شده است. سن ماشین در سطح معنی داری ۰/۰۵ در هر دو مدل در سال ۸۶ روی تعداد ادعای خسارت معنی دار است. هرچه قدر سن ماشین بیشتر باشد، در این صورت لگاریتم متوسط تعداد ادعای خسارت بیشتر خواهد بود، یعنی سن اتومبیل با متوسط تعداد ادعای خسارت رابطه لگاریتمی دارد و هرچه قدر سن اتومبیل بیشتر باشد متوسط تعداد خسارت‌ها بیشتر خواهد بود. همچنین نوع اتومبیل نیز در هیچکدام از مدل‌ها معنی دار نشده است، ولی محل رانندگی در سطح معنی داری ۰/۰۵ در این دو مدل معنی دار است. بر اساس لوجیت نسبت صفرهای زیاد مربوط به تعداد ادعای خسارت، سن اتومبیل و نوع اتومبیل در هیچ سطحی معنی دار نبوده و بنا بر این در هیچ کدام از مدل‌ها اثرگذار نیستند. اما محل رانندگی در سطح ۰/۰۵ در این دو مدل معنی دار می‌باشد. هر چه قدر محل رانندگی در جاهای کم ترافیک باشد لوجیت نسبت صفرهای زیاد مربوط به تعداد ادعای خسارت کمتر خواهد بود.

مورد علاقه، در این مقاله مورد بررسی قرار گرفته‌اند. تعداد ادعای خسارت به‌عنوان یک متغیر شمارشی در نظر گرفته می‌شود، به‌طوری که برای متغیر تعداد ادعای خسارت توزیع‌های بیان شده در بخش دوم، در نظر گرفته می‌شود. در جدول ۲ متغیرهای تبیینی معرفی شده‌اند. در جدول ۳ اطلاعات مربوط به فراوانی تعداد ادعای خسارت برای سال‌های ۱۳۸۶ ارائه گردیده است. همان‌طور که ملاحظه می‌شود در سال ۱۳۸۶ کسانی که هیچ ادعای خسارتی نداشته‌اند، ۷۲ درصد؛ کسانی که یک ادعای خسارت داشته‌اند، ۱۴ درصد؛ کسانی که دو ادعای خسارت داشته‌اند، ۶ درصد؛ کسانی که سه ادعای خسارت داشته‌اند، ۴ درصد؛ کسانی که بیش از ۳ ادعای خسارت داشته‌اند، ۴ درصد بوده است.

۱.۵ مدل‌های برازش شده روی داده‌ها

در این بخش مدل‌های بیان شده در بخش سوم روی داده‌های تعداد ادعای خسارت برازش داده می‌شود. این مدل‌ها به شرح زیر هستند:

۱-مدل اول: در این مدل برای تعداد ادعای خسارت در بیمه شخص ثالث (پاسخ گسسته) مدل رگرسیونی پوآسون آماسیده صفر در نظر گرفته می‌شود:

$$Y \sim ZIP(\mu, \pi),$$

$$\mu = \exp(\gamma_0 + \gamma_1 x_1 + \gamma_2 x_{11} + \gamma_3 x_{12} + \gamma_4 x_{21} + \gamma_5 x_{22}),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 x_1 + \beta_2 x_{11} + \beta_3 x_{12} + \beta_4 x_{21} + \beta_5 x_{22}.$$

۲-مدل دوم: در این مدل برای تعداد ادعای خسارت در بیمه شخص ثالث (پاسخ گسسته) مدل رگرسیونی هاردل پوآسونی در نظر گرفته می‌شود:

$$Y \sim \text{HurdlePois}(\mu, \pi),$$

$$\mu = \exp(\gamma_0 + \gamma_1 x_1 + \gamma_2 x_{11} + \gamma_3 x_{12} + \gamma_4 x_{21} + \gamma_5 x_{22}),$$

$$\text{logit}(\pi) = \beta_0 + \beta_1 x_1 + \beta_2 x_{11} + \beta_3 x_{12} + \beta_4 x_{21} + \beta_5 x_{22}.$$

جدول ۲. متغیرهای تبیینی مورد استفاده

متغیرهای تبیینی	سن اتومبیل	نوع اتومبیل: پژو	نوع اتومبیل: پراید	محل رانندگی: کم ترافیک	محل رانندگی: ترافیک متوسط
سال ۸۶	x_1	x_{11}	x_{12}	x_{21}	x_{22}

نسبت به مدل‌های دیگر به داده‌های تعداد ادعاهای خسارت دارد.

۶ بحث و نتیجه‌گیری

در این مقاله به معرفی مدل رگرسیونی آماسیده و هاردل سری توانی پرداخته شد و انواع این مدل‌ها مورد بررسی قرار گرفت. این مدل‌ها نقش عمده‌ای در تعیین عوامل مؤثر بر تعداد ادعای خسارت و شدت خسارت ایفا کرده و شرکت‌های بیمه با استفاده از این مدل‌ها می‌توانند عوامل اثرگذار روی تعداد ادعای خسارت را شناسایی کرده و تعداد ادعای خسارت و حق بیمه در چند سال آینده را به راحتی پیش‌بینی کنند. در این داده‌ها، عواملی همچون نوع اتومبیل، سن اتومبیل و مکان رانندگی به‌عنوان عوامل اثرگذار روی تعداد ادعای خسارت و شدت خسارت مطرح شده است. مدل‌های مذکور را می‌توان با در نظر گرفتن گمشدگی در داده‌های تعداد ادعای خسارت بررسی کرد.

این موضوع به ما نشان می‌دهد که سن اتومبیل و نوع آن در صفرهای زیاد مربوط به تعداد ادعای خسارت 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. نتایج به‌دست آمده به نحوی توانایی مدل‌های آماسیده صفر و هاردل را نیز نشان می‌دهد.

نتایج برآورد پارامترهای مدل‌های ZINB و HurdleNB در جدول ۵ گردآوری شده است. نتایج یکسانی نیز در این دو مدل نسبت به دو مدل قبل رخ می‌دهد. در مدل‌های رگرسیونی هردو مدل ZINB و HurdleNB به دلیل این که پارامتر پراکنش در سطح معنی داری ۰/۰۵ معنی دار است و همچنین به دلیل وجود صفرهای ساختاری زیاد بر اساس معیارهای AIC این مدل به‌عنوان بهترین مدل انتخاب می‌شود یعنی این که برازش بهتری

جدول ۳. درصد فراوانی ادعای خسارت برای سال‌های ۱۳۸۶

تعداد تصادفات	درصد فراوانی ادعاها در سال ۱۳۸۶
۰	۷۲
۱	۱۴
۲	۶
۳	۴
۴	۲
۵	۱
۶	۱
۷	۰
مجموع	۱۰۰

جدول ۴. نتایج مربوط به مدل‌های مورد نظر (HurdlePois, ZIP)

HurdlePois		ZIP		مدل
انحراف معیار	برآورد	انحراف معیار	برآورد	
μ				
۱/۸۹۰	۰/۲۰۷	۰/۴۱۵	۰/۱۸۱	عرض از مبدأ
۰/۵۱۷	۱.۰۳۰	۰/۰۷۷	۱.۰۲۹	سن ماشین
۱/۵۱۹	۱/۱۳۰	۱/۲۲۳	۱/۱۲۸	نوع ماشین: پژو
۱/۰۳۲	۱/۰۷۳	۱/۱۴۶	۱/۰۷۳	نوع ماشین: پراید
۱.۰۷۶	۱.۵۸۸	۰.۴۲۳	۱.۵۸۸	محل رانندگی: کم ترافیک
۱.۰۸۸	۱.۶۲۹	۰.۴۴۶	۱.۶۵۷	محل رانندگی: ترافیک متوسط
π				
۰/۳۰۹	-۰/۴۰۳	۰/۳۳۴	-۰/۳۴۰	عرض از مبدأ
۰/۲۱۵	۰/۲۳۴	۰/۲۴۴	۰/۲۶۸	سن ماشین
۰/۳۳۷	۰/۶۶۵	۰/۶۱۲	۰/۶۴۷	نوع ماشین: پژو
۰/۵۵۱	۰/۵۰۹	۰/۵۰۴	۰/۷۳۴	نوع ماشین: پراید
۰/۰۰۱	-۰/۶۱۱	۰/۰۰۱	-۰/۶۴۳	محل رانندگی: کم ترافیک
۰/۰۰۱	-۰/۴۱۱	۰/۱۰۳	-۰/۵۴۳	محل رانندگی: ترافیک متوسط
۱۸۵۷/۱۴۷		۱۸۶۵/۳۲۳		معیار AIC

جدول ۵. نتایج مربوط به مدل‌های مورد نظر (HurdleNB, ZINB)

HurdleNB		ZINB		مدل
انحراف معیار	برآورد	انحراف معیار	برآورد	
μ				
۱/۸۷۹	-۱/۰۷۹	۱/۲۷۶	-۲/۳۴۸	عرض از مبدأ
۰.۵۱۷	۱.۲۴۱	۰.۶۲۳	۲.۹۰۲	سن ماشین
۱/۵۱۹	۱/۲۳۹	۰/۱۵۹	۰/۲۶۶	نوع ماشین: پژو
۱/۰۳۲	۱/۲۸۹	۱/۴۰۳	۱/۸۰۰	نوع ماشین: پراید
۰/۰۶۴	۱/۸۷۶	۰/۰۴۵	۱/۴۲۷	محل رانندگی: کم ترافیک
۰/۰۸۲	۲/۲۳۰	۰/۰۳۰	۱/۲۰۵	محل رانندگی: ترافیک متوسط
π				
۰/۱۰۸	-۰/۴۵۳	۰/۱۳۵	-۰/۵۴۲	عرض از مبدأ
۰/۱۴۹	۰/۲۵۸	۰/۲۱۳	۰/۲۷۰	سن ماشین
۰/۴۳۶۷	۰/۶۷۹۴	۰/۶۳۴	۰/۷۰۹	نوع ماشین: پژو
۰/۰۰۱	-۰/۸۰۴	۰/۰۰۱	-۰/۸۹۰	نوع ماشین: پراید
۰/۱۰۰	-۰/۷۹۹	۰/۱۰۳	-۰/۸۷۵	محل رانندگی: کم ترافیک
۰/۰۱۶	۱/۱۳۰	۰/۰۲۱	۱/۰۲۰	محل رانندگی: ترافیک متوسط
۸۰۴/۱۱۲		۸۱۷/۲۸۹		معیار AIC

مراجع

- [1] Greene, W. (2005). Functional form and heterogeneity in models for count data, *Foundations and Trends in Econometrics*, **1(2)**, 113-218.
- [2] Hall, D. B. (2000). Zero-Inflated Poisson and binomial regression with random effects: A case study, *Biometrics*, **56**, 1030-1039.
- [3] Lambert, D. (1992). Zero-inflated Poisson regression, with an application to defects in manufacturing, *Technometrics*, **34**, 1-14.
- [4] Lee, A., Wang, K., Scott, J., Yan, K., and McLachlan, G. (2006). Multi-level zero-inflated Poisson regression modelling of correlated count data with excess zeros", *Statistical Methods in Medical Research*, **15**, 47-61.
- [5] Min, Y. and Agrestic, A. (2005). Random effect models for repeated measures of zero inflated count data, *Statistical Model*, **5**, 5-19.
- [6] Boucher, J.,P., Denuit, M. and Guillee' n, M. (2008). Models of insurance claim counts with time dependence based on generalization of Poisson and negative binomial distributions, *Journal of the Variance*, **1(2)**, 135-162.
- [7] Boucher, J. P. and Guillee' n, M. (2009). A survey on models for panel count data with applications to insurance, *Journal of the Applied Mathematics*, **3(2)**, 280-281.
- [8] Thomas, H. and Samson, D. (1987). Linear models as aids in insured decision making: The estimation of automobile insurance claims, *Journal of the Business*, **15**, 247-256 .