

# برآورد پارامترهای مدل رگرسیون لوژستیک به کمک ماکسیمم آنتروپی تعمیم یافته

مهسا مرکانی<sup>۱</sup>، منیژه صانعی طبس<sup>۲</sup>، حبیب نادری<sup>۳</sup>، حامد احمدزاده<sup>۴</sup>، جواد جمالزاده<sup>۵</sup>

تاریخ دریافت: ۱۴۰۰/۰۷/۲۵

تاریخ پذیرش: ۱۴۰۰/۱۲/۲۶

## چکیده:

هنگام کار با یک مجموعه داده رگرسیونی ممکن است برخی شرایط برقرار نباشند و محدودیت‌هایی برای اجرای مدل رگرسیون به وجود آیند که ما را در استفاده از روش حداقل مربعات با مشکل مواجه می‌کند. روش ماکسیمم آنتروپی تعمیم یافته با زیربنای رگرسیونی قادر است پارامترهای مدل رگرسیونی را بدون در نظر گرفتن هیچ شرطی روی توزیع احتمال خطاها برآورد کند. پیش از این توانمندی این روش در حجم نمونه‌های کم بررسی و تأیید شده است. هنگامی که متغیر پاسخ یک متغیر کیفی است روش رگرسیون لوژستیک به کار می‌رود در این پژوهش ابتدا روش ماکسیمم آنتروپی تعمیم یافته را برای یک مدل رگرسیونی لوژستیک معرفی کردیم. یک نمونه تصادفی از مشتریان بانک جمع‌آوری شده و در این بررسی برای برآورد پارامترهای مدل از مدل رگرسیون لوژستیک دودویی و با استفاده از دو روش ماکسیمم آنتروپی تعمیم یافته و درست‌نمایی ماکسیمم تجزیه و تحلیل و کار آماری انجام گرفت و در نهایت دو روش ذکر شده را باهم مقایسه می‌کنیم. با استناد به مقدار آماره میانگین مربعات خطا برای پیش‌گویی تقاضای مشتری برای افتتاح حساب بلند مدت که از رگرسیون لوژستیک با استفاده از دو روش ماکسیمم آنتروپی تعمیم یافته و درست‌نمایی ماکسیمم به دست آمده، مشخص شد که روش برآورد ماکسیمم آنتروپی تعمیم یافته دارای نتایج دقیق‌تری است.

**واژه‌های کلیدی:** رگرسیون لوژستیک، ماکسیمم آنتروپی تعمیم یافته، درست‌نمایی ماکسیمم، میانگین مربعات خطا

## ۱ مقدمه

نوع از رگرسیون رابطه بین چند متغیر توضیحی با یک متغیر پاسخ دودویی<sup>۶</sup> را مورد تجزیه و تحلیل قرار می‌دهد. در اینجا فقط دو حالت ممکن برای متغیر پاسخ وجود دارد، موفقیت یا عدم موفقیت، وقوع یک رخداد یا عدم وقوع آن و ... .

روش درست‌نمایی ماکسیمم<sup>۸</sup> ( $ML$ ) محبوب بسیاری از محققین بوده و برای برآورد بسیاری از مسائل به کار می‌رود. نقطه شروع این تاریخچه را در کارهای پیرسون برای اندازه‌گیری نیکویی برازش یک مدل و روش گشتاور پیرسون باید جستجو کرد [۱۳] و [۱۴]. این منطق توسط فیشر با کارش بر روی روش  $ML$  ادامه پیدا کرد [۴] و [۵]. اما وقتی پذیره‌های زیربنایی رگرسیونی برقرار نباشد باید به دنبال یک روش جایگزین برای برآورد باشیم. جاج و گلان پارامترهای مدل رگرسیونی

وقتی یک مدل رگرسیونی می‌خواهیم به یک مجموعه داده برازش دهیم باید یکسری پذیره‌های زیربنایی شامل (۱) نرمال بودن باقیمانده‌ها (۲) صفر بودن میانگین باقیمانده‌ها (۳) همگنی واریانس باقیمانده‌ها (۴) استقلال باقیمانده‌ها برقرار باشند که عدم وجود هر یک از این شرایط ما را در به کار بردن روش رگرسیونی حداقل مربعات با مشکل مواجه می‌کند.

رگرسیون لوژستیک<sup>۶</sup> برای تحلیل روابط بین متغیر پاسخ و یک یا چند متغیر توضیحی به کار می‌رود، این نوع رگرسیون به عنوان جایگزینی برای روش طبقه‌بندی فیشر و تحلیل تشخیصی به کار رفته است. این

<sup>۱</sup> دانشجوی کارشناسی ارشد دانشگاه سیستان و بلوچستان

<sup>۲</sup> عضو هیئت علمی دانشگاه سیستان و بلوچستان

<sup>۳</sup> عضو هیئت علمی دانشگاه سیستان و بلوچستان (نویسنده مسئول: h.h.naderi@gmail.com)

<sup>۴</sup> عضو هیئت علمی دانشگاه سیستان و بلوچستان

<sup>۵</sup> عضو هیئت علمی دانشگاه سیستان و بلوچستان

<sup>۶</sup> Logistic Regression

<sup>۷</sup> Binary

<sup>۸</sup> Maximum Likelihood

## ۲ معرفی روش ماکسیمیم آنتروپی تعمیم یافته در مدل رگرسیونی لوژستیک

روش  $GME$  برای اولین بار با چارچوب رگرسیونی مطرح و به کار گرفته شد. این روش ضمن دوباره پارامتری کردن مدل خطی مربوطه پارامترها و خطاهای مجهول را به فرم احتمالاتی که باید برآورد شوند در نظر می‌گیرد. بدین ترتیب توزیع احتمال هر یک از عناصر مجهول برآورد گردیده و مجموع آنتروپی‌های توزیع ضرایب و خطاهای متناظر ماکسیمیم می‌شود. مدل رگرسیون خطی زیر با  $T$  مشاهده و  $K$  متغیر توضیحی در نظر بگیرد:

$$Y = X\beta + e, \quad (1)$$

که در آن  $Y$  یک بردار  $T$  بعدی از متغیر تصادفی مشاهده شده،  $X$  یک ماتریس  $T \times K$  از متغیرهای توضیحی و  $\beta = (\beta_1, \dots, \beta_K)$  یک بردار  $K$  بعدی از پارامترهای مجهول است که از روی داده‌ها برآورد می‌شوند و  $e$  نیز بردار  $T$  بعدی از خطاهای تصادفی غیرقابل مشاهده هستند. برای تعیین تکیه‌گاه خطا، گلان و همکاران از قانون  $3\sigma$  پولکیشم<sup>۱۱</sup> استفاده کرده‌اند کران پایین را  $-3\sigma$  و کران بالا را  $+3\sigma$  قرار داده که مقدار  $\sigma$  نیز انحراف معیار بردار  $Y$  است ([۱۵]). بنابراین در ادامه به کمک تکیه‌گاه پیشین، برآوردهای  $GME$  ضرایب رگرسیونی و جملات خطا به صورت،

$$\beta_k^{GME} = \sum_{m=1}^M \hat{p}_{km} z_{km}, \quad k = 1, 2, \dots, K. \quad (2)$$

$$e_t^{GME} = \sum_{j=1}^J \hat{w}_{tj} v_{tj}, \quad t = 1, 2, \dots, I. \quad (3)$$

حاصل می‌شوند. ملاحظه می‌شود که جواب‌های  $\beta_k^{GME}$  و  $e_t^{GME}$  هر دو به ضرایب لاگرانژ  $\lambda_t$  وابسته بوده و هیچ فرم بسته‌ای برای آن‌ها وجود ندارد و باید مقدار آن‌ها به روش‌های عددی محاسبه گردد ([۱]). در مواردی که متغیر پاسخ کیفی است رگرسیون لوژستیک موردنیاز است. حال برای برآورد ضرایب رگرسیون لوژستیک چه باید کرد؟ در تحلیل رگرسیون لوژستیک، همیشه یک متغیر پاسخ و معمولاً مجموعه‌ای از متغیرهای توضیحی وجود دارند که ممکن است دودویی،

را به کمک روش ماکسیمیم آنتروپی تعمیم یافته<sup>۹</sup> ( $GME$ ) در مسائل اقتصادی برآورد کردند ([۱۰]). سپس گلان و همکاران در سال ۱۹۹۶ روش  $GME$  را معرفی کردند و نشان دادند که برای استفاده از این روش نیازی به برقراری پذیره‌های یادشده نیست و پس‌از آن محققان مختلف در مسائل مختلف از این روش استفاده کردند که از آن جمله می‌توان به موارد زیر اشاره کرد:

سیاولینو و کالکانی در پژوهشی که دارای حجم نمونه کم و هم خطی چندگانه بود، روش برآورد  $GME$  را استفاده کردند ([۳]). سربوونجیتا و همکاران در مسائل اقتصادی از روش  $GME$  بهره بردند ([۱۷]). یادھیس و همکاران در پژوهشی که دارای حجم نمونه کم بود، روش  $GME$  را استفاده کردند ([۱۸]). سانگ و ویت مدل اقتصادسنجی نوین را برای برآورد و پیش‌بینی تقاضای گردشگری ارائه دادند و به روندهای تحقیق مربوط به پیش‌بینی تقاضای گردشگری اشاره کردند ([۱۶]). هدف از انجام این پژوهش مقایسه برآورد روش  $GME$  و برآورد روش  $ML$  برای مدل رگرسیون لوژستیک است. با توجه به اهمیت و کاربرد روش  $GME$ ، در این مطالعه تلاش گردید تا قابلیت این روش با روش  $ML$  با استفاده از آماره میانگین مربعات خطا<sup>۱۰</sup> ( $MSE$ ) مشخص شود. فرض کنید از مشتریان یک بانک در خصوص اینکه آیا مجدداً تصمیم به افتتاح حساب سپرده بلند مدت در این بانک دارند یا خیر پرسیده شده است، برخی جوابشان مثبت و برخی جواب منفی داده‌اند و بنا به هر دلیلی ترجیح داده‌اند در این بانک حساب سپرده بلند مدت نداشته باشند. پس متغیر پاسخ یک متغیر کیفی دوحالتی است. نمونه‌ای به حجم ۳۹۹ نفر به‌طور تصادفی شامل اطلاعاتی از قبیل سن، شغل، وضعیت تأهل، میزان تحصیلات و وام (منظور این است که آیا مشتری وام گرفته است یا خیر). به‌عنوان متغیرهای توضیحی و تقاضای مشتری برای افتتاح حساب بلند مدت به‌عنوان متغیر پاسخ گرفته شد. با توجه به اینکه متغیر پاسخ، یک متغیر دودویی است رگرسیون لوژستیک باید به کار گرفته شود و برای برآورد پارامترهای مدل از روش برآورد  $GME$  استفاده شد. همچنین با برآورد پارامترهای مدل با استفاده از روش  $ML$  و استناد به آماره  $MSE$  مشخص شد که روش  $GME$  از دقت بالاتری برخوردار است. لذا در ادامه با توجه به اینکه مدل رگرسیون این پژوهش از نوع لوژستیک است روش برآورد مربوطه برای مدل لوژستیک شرح داده شده است.

<sup>9</sup>Generalized Maximum Entropy

<sup>10</sup>Mean Square Error

<sup>11</sup>Pullkishm

دارند به طوری که مقادیری مثبت هستند و نیز هر  $e_{ij}$  را می نویسیم،

$$e_{ij} = w_1(-1) + w_2(0) + w_3(1) = w_1(-1) + w_3(1)$$

نقش  $H$ ، تعیین مقدار  $H$  و مثال های تئوری و بحرانی که در آن ها  $H$  چگونه انتخاب شده است در مقاله گلان و همکاران (۱۹۹۶) آورده شده است.

$$\max(-P' \ln P - W' \ln W) \quad (4)$$

تحت  $(JK)$  محدودیت گشتاوری،

$$\sum_{i=1}^T y_{ij} x_{ik} = \sum_{i=1}^T x_{ik} p_{ij} + \sum_{i=1}^T \sum_{h=1}^H x_{ik} \nu_h w_{ijh} \quad (5)$$

$$k = 1, 2, \dots, K, j = 1, 2, \dots, J$$

و همچنین با توجه به اینکه  $w_{ijh}$  و  $p_{ij}$  باید در شرط تابع احتمال بودن صدق کنند این شرایط را داریم،

$$\begin{cases} \sum_{j=1}^J p_{ij} = 1, & i = 1, 2, \dots, T \\ \sum_{h=1}^H w_{ijh} = 1, & i = 1, 2, \dots, T, j = 1, 2, \dots, J \end{cases} \quad (6)$$

حال معادله لاگرانژ را به صورت زیر تشکیل می دهیم،

$$\begin{aligned} L = & - \sum_{i=1}^T \sum_{j=1}^J p_{ij} \ln(p_{ij}) - \sum_{i=1}^T \sum_{j=1}^J \sum_{h=1}^H w_{ijh} \ln(w_{ijh}) \\ & - \sum_{k=1}^K \sum_{j=1}^J \lambda_{jk} \left[ \sum_{i=1}^T y_{ij} x_{ik} - \sum_{i=1}^T x_{ik} p_{ij} - \sum_{i=1}^T \sum_{h=1}^H x_{ik} \nu_h w_{ijh} \right] \\ & + \sum_{i=1}^T \alpha_i \left[ \sum_{j=1}^J p_{ij} - 1 \right] + \sum_{i=1}^T \sum_{j=1}^J \gamma_{ij} \left[ \sum_{h=1}^H w_{ijh} - 1 \right] = 0 \quad (7) \end{aligned}$$

با مشتق گیری از معادله (۷) نسبت به  $p_{ij}$  رابطه زیر به دست آمده می آید.

$$\frac{\partial L}{\partial p_{ij}} = -\ln p_{ij} - 1 - \sum_{k=1}^K \lambda_{jk} x_{ik} + \alpha_i = 0 \quad (8)$$

$$-\ln p_{ij} = 1 - \alpha_i + \sum_{k=1}^K \lambda_{jk} x_{ik}$$

در نتیجه،

$$p_{ij} = e^{-(1 - \alpha_i + \sum_{k=1}^K \lambda_{jk} x_{ik})} \quad (9)$$

از طرفی با استفاده از محدودیت  $\sum_{j=1}^J p_{ij} = 1$  در رابطه (۹) داریم،

$$\sum_{j=1}^J e^{-(1 - \alpha_i + \sum_{k=1}^K \lambda_{jk} x_{ik})} = 1$$

$$e^{\alpha_i - 1} \sum_{j=1}^J e^{-\sum_{k=1}^K \lambda_{jk} x_{ik}} = 1$$

$$e^{\alpha_i - 1} = \frac{1}{\sum_{j=1}^J e^{-\sum_{k=1}^K \lambda_{jk} x_{ik}}}$$

کمی یا ترکیبی از آن ها باشند. به علاوه لازم نیست متغیرهای دودویی به طور واقعی دوتایی باشند. برخی از متغیرهای توضیحی در رگرسیون لوژستیک می توانند به عنوان متغیرهای کمکی مورد استفاده قرار گیرند تا پژوهشگران بتوانند با ثابت نگه داشتن یا کنترل آماری این متغیرها اثرات دیگر متغیرهای توضیحی را بهتر ارزیابی کنند.

با این که رگرسیون لوژستیک در مقایسه با رگرسیون خطی پیش فرض های کمتری دارد (به عنوان مثال پیش فرض های همگنی واریانس و نرمال بودن خطاها وجود ندارد)، رگرسیون لوژستیک نیازمند موارد زیر است:

۱. هم خطی چندگانه کامل وجود نداشته باشد.
۲. خطاهای خاص نباید وجود داشته باشد (یعنی همه متغیرهای پیش بین مرتبط وارد شوند و پیش بین های نامربوط کنار گذاشته شوند).
۳. متغیرهای توضیحی باید در مقیاس فاصله ای یا سطح نسبی اندازه گیری شده باشند (هرچند که متغیرهای دودویی نیز می توانند مورد استفاده قرار گیرند).

دو استفاده مهم از رگرسیون لوژستیک عبارتند از:

- ۱- پیش بینی اعضای گروه: در رگرسیون لوژستیک احتمال موفقیت قبل از احتمال شکست محاسبه می شود، بنابراین نتایج به فرم نسبت شانس به دست آمده می آیند.
- ۲- شناسایی و اطلاع از روابط و شدت آن در بین متغیرها.

در این بخش روش  $GME$  برای مدل رگرسیون لوژستیک را بیان می کنیم. همچنان که قبلاً متذکر شدیم با توجه به اینکه وزن های مجهول پارامترها و تکیه گاه خطاهای مدل رگرسیونی از یکدیگر مستقل اند [۶]، [۷] و [۸] می توان با حل مسئله بهینه سازی ماکسیم آنتروپی زیر، پارامترها و خطاهای مجهول را برآورد کرد. اطلاع محدودیت گشتاوری

را در مسئله آنتروپی باید به کار ببندیم پس توابع غیر قابل مشاهده و مجهول  $p$  و  $e$  خواص تابع احتمال را باید داشته باشند. عناصر  $p$  که به فرم احتمالات هستند و اما اعداد حقیقی مقدار  $e_{ij}$  روی بازه  $[-1, 1]$  تغییر می کند. چون این خطاها به فرم احتمال نیستند، یک

ابزار مفهومی بازنویسی پارامتری که توسط جاج (۱۹۹۱) و جاج و گلان (۱۹۹۳) برای تبدیل  $e_{ij}$  ها به مقادیر احتمالی که روی بازه  $(0, 1)$  تغییر

کنند را به کار می گیریم پس روی بازه  $[-1, 1]$  متغیر تصادفی گسسته کران دار با مقادیر ممکن اعداد حقیقی  $\nu_{ij} = (\nu_{ij1}, \nu_{ij2}, \dots, \nu_{ijH})$  با بعد  $H \geq 2$  تعریف می کنیم. با وزن های متناظر مجهول  $w_{ij}$  که  $w_{ij} = (w_{ij1}, w_{ij2}, \dots, w_{ijH})$  که دارای خواص تابع احتمال هستند و باید  $\sum_{h=1}^H w_{ijh} = 1$  و بنابراین  $e_{ij} = \sum_{h=1}^H w_{ijh} \nu_h$  مثلاً اگر قرار

دهیم  $H = 3$  آنگاه  $\nu = (-1, 0, 1)$  پس نقاط  $w_1$  و  $w_2$  و  $w_3$  وجود

و برآورد تابع  $p_{ij}$  به صورت زیر است:

$$\hat{p}_{ij} = \frac{e^{-\sum_{k=1}^K \lambda_{jk} x_{ik}}}{\sum_{j=1}^J e^{-\sum_{k=1}^K \lambda_{jk} x_{ik}}} \quad (10)$$

به طور مشابه برای به دست آمده آوردن  $\hat{w}_{ijh}$  نیز از معادله لاگرانژ (۷) نسبت به  $w_{ijh}$  مشتق می گیریم که نتایج،

$$\frac{\partial l}{\partial w_{ijh}} = -\ln w_{ijh} - 1 - \sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h + \gamma_{ij} = 0$$

در نتیجه،

$$w_{ijh} = e^{-(1-\gamma_{ij} + \sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h)} \quad (11)$$

به دست آمده می آیند. با استفاده از محدودیت  $\sum_{h=1}^H w_{ijh} = 1$  در رابطه (۱۱) داریم،

$$\sum_{h=1}^H e^{-(1-\gamma_{ij} + \sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h)} = 1$$

$$e^{-1+\gamma_{ij}} \sum_{h=1}^H e^{-\sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h} = 1$$

$$e^{-1+\gamma_{ij}} = \frac{1}{\sum_{h=1}^H e^{-\sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h}} \quad (12)$$

و برآورد تابع  $p_{ij}$  به صورت زیر است:

$$\hat{w}_{ijh} = \frac{e^{-\sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h}}{\sum_{h=1}^H e^{-\sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h}} \quad (13)$$

و بنابراین مقدار  $e_{ij}$  به دست آمده می آید،

$$e_{ij} = \sum_{h=1}^H \hat{w}_{ijh} \nu_h \quad (14)$$

اگر قرار دهیم،

$$e_j = \nu_j w_j = \begin{pmatrix} \nu_{j1} & 0 & \dots & 0 \\ 0 & \nu_{j2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \nu_{jT} \end{pmatrix} \begin{pmatrix} w_{j1} \\ w_{j2} \\ \vdots \\ w_{jT} \end{pmatrix}$$

که در آن،

$$\sum_{h=1}^H \nu_{ijh} w_{ijh} = e_{ij} \quad i = 1, 2, \dots, T, \quad j = 1, 2, \dots, J$$

جواب های (۱۰) و (۱۳) به فرم نمایی هستند. پس بنابراین  $\hat{w}_{ij}$  و  $\hat{p}_{ij}$  همیشه مثبت هستند. عوامل نرمال سازی ما را مطمئن می سازد که  $\hat{p}_{ij}$

و  $\hat{w}_{ij}$  دارای خواص تابع احتمال هستند و اطلاعاتی در مورد توزیع احتمالات برای پارامترهای مجهول فراهم می کند. پارامترهای مجهول بازنویسی شده  $\hat{p}_{ij}$  و  $\hat{e}_{ij}$  فقط بر اساس اطلاعات روابط گشتاوری و نه هیچ پذیره اضافه دیگری پایه گذاری شده اند. چون  $\lambda_j$  های لاگرانژ یکتا نیستند یک روش نرمال سازی رایج این است که قرار دهیم  $\lambda_1 = \beta_1 = 0$  که

$$0 = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}_{K \times 1}$$

روش  $GME$  نسبت به روش ماکسیمیم آنتروپی و درست نمایی ماکسیمیم لوژیستیک چند جمله ای پذیره ها و محدودیت های کمتری دارد. در معادله لاگرانژ (۷) می توان قرار داد  $\beta_j = -\lambda_j$ . یک انتخاب برای نقاط  $\nu_j$  می تواند  $\pm \frac{1}{\sqrt{T}}$  باشد به عنوان مثال اگر  $T = 100$  و  $H = 3$  آنگاه  $\nu_{ij} = (0, 1, 0, 0, 1)$

$$\hat{p}_{ij} = \frac{e^{-\sum_{k=1}^K \lambda_{jk} x_{ik}}}{\sum_{j=1}^J e^{-\sum_{k=1}^K \lambda_{jk} x_{ik}}}$$

اگر قرار دهیم  $\beta_{jk} = -\lambda_{jk}$

$$\hat{p}_{ij} = \frac{e^{\sum_{k=1}^K \beta_{jk} x_{ik}}}{\sum_{j=1}^J e^{\sum_{k=1}^K \beta_{jk} x_{ik}}}$$

و همچنین،

$$\hat{w}_{ijh} = \frac{e^{-\sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h}}{\sum_{j=1}^J e^{-\sum_{k=1}^K \lambda_{jk} x_{ik} \nu_h}}$$

برای  $\beta_{jk} = -\lambda_{jk}$  به دست آمده می آوریم.

$$\hat{w}_{ijh} = \frac{e^{\sum_{k=1}^K \beta_{jk} x_{ik} \nu_h}}{\sum_{j=1}^J e^{\sum_{k=1}^K \beta_{jk} x_{ik} \nu_h}}$$

و همچنین،

$$\hat{e}_{ij} = \sum_{h=1}^H \hat{w}_{ijh} \nu_h$$

می توان قرار داد،

$$\hat{y}_{ij} = \hat{p}_{ij} + \hat{e}_{ij}$$

ضرایب رگرسیونی  $\beta_j = -\lambda_j$  برای تعیین هر یک از این ضرایب می توان از نسبت بخت ها کمک گرفت. در رگرسیون لوژیستیک معمولی

که

خطا هستند.  $AG_i$  متغیر ظاهری است،  $AG_i = 1$  اگر سن عامل مؤثر برای افتتاح حساب بلند مدت باشد، در غیر این صورت  $AG_i = 0$ . شغل  $JO_i$ ، وضعیت تأهل  $MA_i$ ، میزان تحصیلات  $ED_i$  و وام  $LO_i$  وجود دارند.

$$\ln\left(\frac{p_{ij}}{p_i}\right) = -\lambda_j x_i \Rightarrow \ln\left(\frac{p_{ij}}{p_i}\right) = -\sum_{k=1}^K \lambda_{jk} x_{ik} = \sum_{k=1}^K \beta_{jk} x_{ik}$$

$$j = 1, 2, \dots, J, i = 1, 2, \dots, T, k = 1, 2, \dots, K$$

پس  $JK$  مجهول داریم و چون  $i = 1, 2, \dots, T$  پس تعداد  $T$  دستگاه  $J$  معادله‌ای داریم ([۱۲]).

عیار ارزیابی

در این مطالعه به منظور ارزیابی کارایی برآوردگر برای پیش‌بینی از میانگین مربعات خطا استفاده می‌کنیم. که برابر رابطه زیر است:

$$MSE = \frac{1}{T} \sum_{i=1}^T (y_i - \hat{y}_i)^2$$

که در آن  $T$  تعداد مشاهدات است و  $(y_i - \hat{y}_i)^2$  نشان‌دهنده خطای پیش‌بینی است.

### ۳ مثال کاربردی

در این پژوهش به منظور مقایسه روش  $GME$  و  $ML$  از داده‌های واقعی که مربوط به مشتریان بانک بودند استفاده شده است. داده‌های این پژوهش از بانکی در کشور پرتغال گردآوری شده است که مربوط به سال ۲۰۱۴ است. جامعه آماری شامل مشتریان بانک است که حجم این داده‌ها حدود ۴۰۰۰ بوده و نمونه‌ای به حجم ۳۹۹ نفر به طور تصادفی گرفته شده است. بر اساس مبانی نظری مطرح شده در مقدمه از متغیرهای توضیحی از قبیل سن ( $Age$ )، شغل ( $Job$ )، وضعیت تأهل ( $Maritalstatus$ )، میزان تحصیلات ( $Educationlevel$ ) و وام (منظور این است که آیا مشتری وام گرفته است یا خیر). به عنوان متغیرهای مؤثر در پیش‌بینی عملکرد مشتریان در راستای افتتاح حساب بلند مدت استفاده گردید و تقاضای مشتری برای افتتاح حساب بلند مدت ( $termdeposit$ ) به عنوان متغیر پاسخ گرفته شد. متغیرهایی که کیفی بودند مانند شغل، وضعیت تأهل، میزان تحصیلات و وام به صورت متغیرهای ظاهری (۰ و ۱) تعریف شدند.

در اینجا مدل رگرسیونی لوژیستیک به صورت زیر تعریف می‌شود:

$$\text{Logit}(DEP_i) = \beta_0 + \beta_1 AG_i + \beta_2 JO_i + \beta_3 MA_i + \beta_4 ED_i + \beta_5 LO_i + \epsilon_i$$

که در آن  $DEP_i = 1$  اگر مشتریان افتتاح حساب بلند مدت داشته باشند و صفر در غیر این صورت.  $\beta_i$  متغیرهای پیش‌گو و  $\epsilon_i$  جملات

در این پژوهش از نرم‌افزار گمز برای تحلیل استفاده شده است. نتایج برازش مدل رگرسیونی لوژیستیک با استفاده از دو روش برآورد  $GME$  و  $ML$  برای سه حجم نمونه مختلف ۱۸۰، ۲۷۰، ۳۹۹ نفر در جدول (۱) آورده شده است.

طبق جدول (۱) مشاهده می‌شود که در حجم‌های مختلف نمونه، مقدار آماره  $MSE$  روش برآورد  $GME$  از  $ML$  کمتر است و لذا منجر به نتایج دقیق‌تری می‌شود. همچنین مشاهده می‌شود که در حجم‌های نمونه کمتر، روش برآورد  $GME$  کارایی بیشتری نسبت به  $ML$  دارد به طوری که تفاوت آماره  $MSE$  دو روش برآورد خیلی بیشتر از حجم‌های نمونه بزرگ‌تر است. ظاهراً هر چه حجم نمونه بیشتر شده، مقدار  $MSE$  هر دو روش  $GME$  و  $ML$  کمی افزایش پیدا کرده است اما در حجم نمونه‌های بالاتر دیگر شاهد این افزایش نیستیم و در واقع این مقادیر برای حجم نمونه‌های خیلی بالاتر کاهش می‌یابند.

### ۴ بحث و نتیجه‌گیری

این تحقیق سعی بر بررسی یک روش جایگزین روش برآورد  $ML$  دارد، یعنی روش  $GME$  در مدل رگرسیونی لوژیستیک، که در این تحقیق از داده‌های واقعی موجود در بانکی در کشور پرتغال به منظور بررسی احتمال افتتاح حساب بلند مدت توسط مشتریان این بانک استفاده گردید. بدین منظور از اطلاعاتی از قبیل سن، شغل، وضعیت تأهل، میزان تحصیلات و وام (منظور این است که آیا مشتری وام گرفته است یا خیر). استفاده گردید. با استفاده از نرم‌افزار گمز، دو روش برآورد  $GME$  و  $ML$  برای برآورد پارامترهای این مدل به کار گرفته شدند. در نهایت با توجه به مقدار آماره  $MSE$  نتیجه گرفته شد که روش برآورد  $GME$  دارای نتایج قابل اعتمادتری است، پیشنهاد شده است که روش  $GME$  ممکن است نسبت به روش  $ML$  برتری داشته باشد، نتیجه حاصل از این پژوهش با نتیجه گرفته شده [۲]، [۹] و [۱۱] کاملاً همسو است.

نام متغیرهای توضیحی	شرح	نوع متغیرهای توضیحی
سن	سن مشتری	عددی
شغل	شغل مشتری	طبقه بندی: مدیر کارگر کارآفرین خدمتکار مدیریت بازنشسته آزاد خدماتی دانشجو تکنسین بیکار ناشناخته
وضعیت تأهل	وضعیت تأهل مشتری	طبقه بندی: مجرد مطلقه متاهل
تحصیلات	سطح تحصیلات مشتری	طبقه بندی: بی سواد زیر دیپلم بالتر از دیپلم ناشناخته
وام	منظور این است که آیا مشتری وام گرفته است یا خیر	طبقه بندی: بله خیر

## مراجع

[۱] صانعی طبس، م. (۱۳۹۴). اصل ماکسیمیم آنتروپی تعمیم یافته مرتبه  $\alpha$  و برآورد پارامترهای مدل رگرسیونی به روش ماکسیمیم آنتروپی تعمیم یافته. رساله دکتری (آمار ریاضی)، دانشکده علوم ریاضی، دانشگاه فردوسی مشهد.

[2] Ayusuk, A. and Autchariyapanitkul, K. (2017). Factors Influencing Tourism Demand to Revisit Pha Ngan Island Using Generalized Maximum Entropy. *Thai Journal of Mathematics*, 187-196.

جدول ۱: برآورد پارامتر مدل رگرسیون لوژستیک برای افتتاح حساب بلند مدت

$n$	۱۸۰		۲۷۰		۳۹۹	
	$ML$	$GME$	$ML$	$GME$	$ML$	$GME$
$\beta_0$	۳/۴۲۶	۱/۴۲۱	۱/۶۸۴	۰/۴۵۷	۱/۵۰۱	۰/۵۶۰
$\beta_1$	-۰/۰۱۹	۰/۰۰۲	-۷/۹۵۴	۰/۰۱۶	-۷/۰۳۶	۰/۰۱۶
$\beta_2$	-۰/۴۰۹	-۰/۰۶۸	-۱/۵۰۲	-۰/۰۵۲	-۱/۴۶۲	-۰/۰۲۹
$\beta_3$	-۰/۲۶۵	-۰/۳۱۳	۰/۵۹۵	-۰/۲۸۹	-۱/۳۱۰	-۰/۳۱۹
$\beta_4$	-۱/۱۵۹	-۰/۰۱۸	-۷/۱۵۱	۰/۰۶۸	۳/۱۷۷	۰/۰۲۱
$\beta_5$	-۱/۵۰۰	-۱/۲۲۳	-۱/۳۶۵	-۰/۹۷۰	-۱/۳۴۵	-۱/۱۳۲
$MSE$	۰/۳۶۱	۰/۱۷۲	۰/۴۲۵	۰/۲۷۳	۰/۴۵۳	۰/۳۹۳

- [3] Ciavolino, E. and Calcagni, A. (2016). A Generalized Maximum Entropy (GME) estimation approach to fuzzy regression model. **38**, 51-63.
- [4] Fisher, R. A. (1912). On an absolute criterion for fitting frequency curves. *Messenger of Mathematics*, **41**, 155-160.
- [5] Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London Series A*, **222**, 309-368.
- [6] Golan, A. (2002). Information and Entropy Econometrics Editor's View. *Journal of Econometrics*, **107**, 1-15.
- [7] Golan, A. (2008). *Information and Entropy Econometrics A Review and Synthesis*. New York, Now Publishers.
- [8] Golan, A., Judge, G. and M. Perloff, J. (1996). A maximum entropy approach to recovering information from multinomial response data. *Journal of the American Statistical Association*, **91**, 841-853.
- [9] Golan, A., Moretti, E. and Perloff, J. M. (2001). A Small Sample Estimator for the Sample Selection Model. *CUDARE Working Papers in University of California, Berkeley*.
- [10] Judge, G. G. and Golan, A. (1992). Recovering information in the case of ill-posed inverse problems with noise. *Mimeo Department Of Agricultural And Natural Resources*, University of California, Berkeley, CA.
- [11] Maneejuk, P. (2021). On regularization of generalized maximum entropy for linear models. *Soft Computing*, **25**, 7867-7875.
- [12] Mittelhammer, R. C., Judge, G. G. and Miller, D. J. (2000). *Econometric Foundations*. New York, Cambridge University Press.
- [13] Pearson, K. (1894). Contribution to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London Series A*, **185**, 71-110.
- [14] Pearson, K. (1900). On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *Philosophical Magazine Series*, **50**, 157-175.

- [15] Pukelsheim, F. (1994). The three-sigma Rule. *The American Statistician*, **48**, 88-91.
- [16] Song, H. and Witt, S. F. (2000). *Tourism demand modelling and forecasting: Modern econometric approaches*. Routledge.
- [17] Sriboonchitta, S., Liu, J. and Sirisrisakulchai. (2015). Willingness-to-pay estimation using generalized maximum-entropy: A case study, *International Journal of Approximate Reasoning*, **60**, 1-7.
- [18] Thi Binh An, D., Tsuchida, J. and Yadohisa, H. (2021). K-means generalized maximum entropy estimation for structural equation modeling. *Behaviormetrika*, **48**, 103–115.



## Estimation of Logistic Regression Model Parameters Using Generalized Maximum Entropy

Mahsa Markani<sup>1</sup> and Manije Sanei Tabas<sup>2</sup> and Habib Naderi<sup>3</sup> and Hamed Ahmadzadeh<sup>4</sup> and Javad Jamalzadeh<sup>5</sup>

### Abstract:

When working with a regression data set, some conditions may not be met and there may be limitations to running the regression model, which makes it difficult for us to use the least squares method. The generalized maximum entropy method with regression infrastructure is able to estimate the parameters of the regression model without considering any conditions on the probability distribution of errors. The capability of this method in small sample sizes has already been investigated and confirmed. When the response variable is a qualitative variable, the logistic regression method is used. In this study, we first introduced the generalized maximum entropy method for a logistic regression model. A random sample of bank customers was collected and in this study to estimate the model parameters from the binary logistic regression model using two generalized maximum entropy methods and maximum likelihood analysis and statistical work was performed and finally compared the two methods. Based on the mean square error statistics for predicting customer demand for long-term account opening, which was obtained from logistic regression using the generalized maximum entropy and maximum likelihood methods, it was found that the generalized maximum entropy estimation method has accurate results.

**Keywords:** Entropy, Generalized maximum entropy, Logistic regression, Logit, Maximum likelihood.

---

<sup>1</sup> Master student, University of Sistan and Baluchestan

<sup>2</sup> Assistant Professor, Dept.of statistics, Faculty of Math, University of Sistan and Baluchestan, Daneshgah Ave., Zahedan, Iran.

<sup>3</sup>Assistant Professor, Dept.of statistics, Faculty of Math, University of Sistan and Baluchestan, Daneshgah Ave., Zahedan, Iran.

<sup>4</sup>Assistant Professor, Dept.of statistics, Faculty of Math, University of Sistan and Baluchestan, Daneshgah Ave., Zahedan, Iran.

<sup>5</sup>Assistant Professor, Dept.of statistics, Faculty of Math, University of Sistan and Baluchestan, Daneshgah Ave., Zahedan, Iran.