



A Credibility Approach to Predicting Healthcare Costs: A Case Study of Dental Costs

Ehsan Jahanbani¹ MSc, Amirteymor Payandeh Najafabadi^{1*} PhD, Khaled Masoumifard¹ PhD

¹ Faculty of Mathematical Sciences, Shahid Beheshti University, Tehran, Iran

Correspondence to: Amirteymor Payandeh Najafabadi, Email: amirpayandeh@gmail.com

Received: March 16, 2022

Revised: June 25, 2022

Accepted: July 24, 2022

Online Published: August 16, 2022

Abstract

Introduction: The increase in the cost of medical services and the aging of Iran's population have created problems for the medical insurance industry, which has turned this field into one of the riskiest fields. In this situation, one of the most basic needs of this field is to fit a suitable model to the treatment losses, so that the insurer can check the trend of the treatment costs, have an accurate estimate of the next year's losses, and calculate the corresponding insurance premium.

Methods: Based on an empirical fundamental study, the loss data related to group treatment insurance policies issued from the middle of 2017 to the middle of 2019 was collected by Day insurance company and using mixed distributions, dental claims were modeled and the loss of the next year was predicted using the credibility method.

Results: Due to the heterogeneity in the data, mixed distribution will be more suitable due to its high flexibility. The insured people were divided into 2 high-risk and low-risk categories according to their characteristics, and a distribution was given to each of these categories. The general distribution of damages is considered as a mixed distribution of these distributions, and the mixed weights of this distribution were estimated using logistic regression.

Conclusion: One of the most basic issues in all types of insurances is determining their insurance premiums. In this research, using the theory of credibility, the next year's insurance premium is calculated for each of the insured according to their risk parameter. In fact, credibility theory combines the existing insurance premium with past losses and provides the adjusted premium.

Keywords: Finite Mixed Distribution, Complementary Treatment, Logistic Regression, Dentistry

Highlights

1. By using credibility, the existing insurance premium can be combined with the past damages and the adjusted insurance premium can be considered.
2. In case of a more accurate estimate of the treatment costs, a fair insurance premium can be calculated for each insured according to its risk parameters, and in addition to increasing the attractiveness for the policyholder, the insurance company will suffer less loss.

Citation:

Jahanbani E, Payandeh Najafabadi A, Masoumifard K. A credibility approach to predicting healthcare costs: a case study of dental costs. Iran J Health Insur. 2022;5(2):96-107.



رویکرد باورمندی به پیشگویی هزینه‌های درمان: مطالعه موردی بررسی هزینه‌های دندانپزشکی

احسان جهان‌بانی¹ MSc، امیر تیمور پاینده نجف‌آبادی^{1*} PhD، خالد معصومی‌فرد¹ PhD

¹ گروه بیم‌سنجی، دانشکده علوم ریاضی، دانشگاه شهید بهشتی، تهران، ایران

* نویسنده مسئول: امیر تیمور پاینده نجف‌آبادی، پست الکترونیک: amirtpayandeh@gmail.com

انتشار آنلاین: ۱۴۰۱/۰۵/۲۵

پذیرش: ۱۴۰۱/۰۵/۰۲

تصحیح: ۱۴۰۱/۰۴/۰۴

دریافت: ۱۴۰۰/۱۲/۲۵

چکیده

مقدمه: افزایش هزینه‌های خدمات پزشکی و سالمند شدن جمعیت ایران، مشکلاتی را برای صنعت بیمه درمان به وجود آورده که این حوزه را به یکی از پرمخاطره‌ترین حوزه‌ها تبدیل کرده است. در این شرایط یکی از اساسی‌ترین نیازهای این حوزه، برآزش مدلی مناسب به خسارت‌های درمان است، به طوری که بیمه‌گر بتواند روند هزینه‌های درمان را بررسی کند، برآورد دقیقی از خسارت‌های سال آینده داشته باشد و حق بیمه متناسب با آن را محاسبه کند.

روش بررسی: با مطالعه بنیادی تجربی، داده‌های خسارت مربوط به بیمه‌نامه‌های درمان گروهی صادره از اواسط سال ۱۳۹۷ تا اواسط ۱۳۹۹ شرک بیمه دی جمع‌آوری شد و با استفاده از توزیع‌های آمیخته، خسارت‌های دندانپزشکی مدل‌بندی شد و با استفاده از روش باورمندی خسارت سال آینده پیش‌بینی شد.

یافته‌ها: باتوجه به ناهمگنی موجود در داده‌ها، توزیع‌های آمیخته به دلیل انعطاف‌پذیری بالا برآزش مناسب‌تری خواهد داشت. افراد بیمه شده باتوجه به مشخصات آنها در ۲ دسته پرریسک و کم‌ریسک تقسیم شدند و به هر یک از این دسته‌ها یک توزیع برآزش داده شده است. توزیع کلی خسارت‌ها را یک توزیع آمیخته از این توزیع‌ها در نظر گرفته که وزن‌های آمیختگی این توزیع با استفاده از رگرسیون لجستیک برآورد شد.

نتیجه‌گیری: یکی از اساسی‌ترین مسائل در انواع بیمه‌ها، تعیین حق بیمه آنهاست. در این پژوهش با استفاده از نظریه باورمندی، برای هر یک از بیمه شده‌ها باتوجه به پارامتر ریسک آنها حق بیمه سال آینده محاسبه می‌شود. در واقع باورمندی، حق بیمه موجود را با خسارت‌های گذشته ترکیب کرده و حق بیمه تعدیل شده را ارائه می‌دهد.

واژگان کلیدی: توزیع آمیخته متناهی، درمان تکمیلی، رگرسیون لجستیک، دندانپزشکی

نکات ویژه

۱. با استفاده از باورمندی، می‌توان حق بیمه موجود را با خسارت‌های گذشته ترکیب و حق بیمه تعدیل شده را در نظر گرفت.
۲. در صورت برآورد دقیق‌تر از هزینه‌های درمان، می‌توان حق بیمه منصفانه‌ای برای هر بیمه شده باتوجه به پارامترهای ریسک آن محاسبه و علاوه بر افزایش جذابیت برای بیمه‌گذار، شرکت بیمه نیز متحمل زیان کمتری خواهد شد.

مقدمه

بهداشت و درمان، به همراه داشته است. افزایش سالمندان، واقعیت انکارناپذیر جوامع کنونی است. این پدیده هشدار به جوامع برای توجه بیشتر به مسائل این جمعیت رو به رشد است. مسائلی چون نبود حمایت‌های اجتماعی، نداشتن شغل و نقش اجتماعی، مخارج زندگی و به‌ویژه هزینه‌های بهداشتی درمانی سرسام‌آور و غیره از جمله مواردی است که لزوم توجه بیشتر به این قشر از جامعه

طی پنجاه سال اخیر توسعه اقتصادی اجتماعی، کاهش باروری و به دنبال آن کاهش رشد جمعیت و افزایش امید به زندگی، منجر به ایجاد تغییرات هنجاری در ساختار جمعیت جهان شده است، به نحوی که طی این مدت تعداد سالمندان به‌طور کلانی افزایش یافته است. این امر هزینه‌های زیادی را، به‌خصوص در بخش

از مهم‌ترین وظایف شرکت‌های بیمه، پیش‌بینی میزان و تعداد خسارت‌ها در سال آتی است. روش‌های مختلفی مانند سری‌های زمانی، روش بیزی و نظریه باورمندی برای انجام چنین پیش‌بینی‌های وجود دارد. نظریه باورمندی روشی غیر پارامتریک است در حالی که روش بیزی و سری‌های زمانی پارامتریک هستند. همه این روش‌ها برای پیش‌بینی تعداد خسارت‌ها در سال آتی قابل اطمینان هستند اما با توجه به آسان بودن اجرای نظریه باورمندی، استفاده از این روش پیشنهاد می‌شود. باورمندی به‌عنوان ابزاری عملی، عمدتاً حاصل کار اکچوئرها آمریکای شمالی است، با این وجود ریشه در نظریه ریسک و برآوردهای آماری دارد. نظریه باورمندی تلاش در راستای ایجاد مدلی است که توسط آن بتوان درباره محیطی که در آن دستگاه‌های بیمه عمل می‌کنند، اطلاع حاصل کرد. نظریه باورمندی سابقه طولانی در علم بیم‌سنجی دارد که اساس آن به موبرای (۱۹۱۴) برمی‌گردد [۷]. ویتنی (۱۹۱۸) [۸] شکل جالبی از میانگین وزنی متوسط خسارت‌های وارده بر یک دسته و سایر دسته‌ها را به‌منظور پیش‌بینی خسارت‌های وارده در سال‌های آتی ارائه کرد. کفر (۱۹۲۹) [۹] برای اولین بار استفاده از روش بیزی برای نرخ‌گذاری در بیمه‌های عمر گروهی را پیشنهاد داد. به دنبال آن بیلی (۱۹۵۰) [۱۰] رابطه نظریه باورمندی خطی و روش بیزی را نشان داد. در مسائل پیشرفته نظریه باورمندی، برای برآورد حق بیمه، روش‌های بیزی به‌کار گرفته می‌شود که به روش‌های باورمندی بیزی موسوم هستند. روش‌های بیزی با ارائه بولمن و بولمن-استراب (۱۹۶۰) [۱۱] وارد علم بیم‌سنجی شد. میلر و هیگمن در سال ۱۹۵۷ از نظریه باورمندی در توزیع تجمعی خسارت‌ها استفاده کردند، به‌علاوه پین کویت (۱۹۹۷) استفاده از نظریه باورمندی را در بیمه اتومبیل مورد مطالعه قرار داد [۱۲]. چون باورمندی به‌الگوهای می‌پردازد که پاسخ‌هایی مطابق با محیط و پیرامون خود ارائه می‌دهد، بیشتر در علم بیم‌سنجی نفوذ کرده است. نظریه باورمندی تلاش در راستای ایجاد مدلی است که توسط آن بتوان درباره محیطی که در آن دستگاه‌های بیمه عمل می‌کنند، اطلاع حاصل کرد. نظریه باورمندی مجموعه‌ای از ابزار کمی است که به بیمه‌گر اجازه می‌دهد حوادث آتی را برای یک ریسک یا گروهی از ریسک‌ها نرخ‌بندی کند، به این معنا که حق بیمه‌های آتی را بر اساس حوادث گذشته تعدیل می‌کند.

را نمایان می‌سازد. در سال ۲۰۲۰ تخمین زده شد یک میلیارد نفر (۱۳ درصد) از جمعیت جهان، افراد سالخورده ۶۵ سال و بالاتر هستند. برآورد شده که تا سال ۲۰۳۰، از هر ۶ نفر در جهان، یک نفر ۶۰ ساله یا بیشتر خواهد بود. انتظار می‌رود این رقم در سال ۲۰۵۰ تقریباً سه برابر شده و به حدود یک و نیم میلیارد نفر (۱۶ درصد از جمعیت جهان) برسد. اگرچه در حال حاضر توسعه یافته‌ترین کشورها، سالخورده‌ترین نمودار جمعیت را دارند، اما سریع‌ترین سالخوردگی جمعیت در کشورهای کمتر توسعه یافته اتفاق می‌افتد. این پدیده با کاهش نرخ باروری و عمر طولانی شکل گرفته است. با کاهش تولد و افزایش طول عمر، سهم افراد سالخورده از کل جمعیت در حال افزایش است [۱].

جمعیت ایران نیز در حال حاضر در مرحله انتقال ساختار سنی از جوانی به سالخوردگی است. هم‌اکنون درصد کمتری از جمعیت کشور در سن سالخوردگی هستند، اما با توجه به کاهش سریع باروری در چند دهه اخیر افزایش تعداد و درصد سالمندان در سال‌های آینده، برنامه‌ریزی آینده‌محور برای حل مشکلات این قشر از جمعیت ضروری است [۲]. موضوع درمان در تمام کشورها، به‌خصوص کشورهای در حال توسعه، یکی از مسائل مهم است و دولت‌ها می‌کوشند شیوه‌های خدمات‌رسانی و ایجاد امکانات لازم در زمینه‌های درمانی را بهبود بخشند [۳]. یکی از راه‌های مناسب برای فراهم ساختن پشتوانه بسیار مطلوب، به‌منظور حمایت از وضعیت اقتصادی خانواده و گذر از موقعیت‌های اضطراری ناشی از بیماری‌ها و حوادث، تهیه بیمه‌نامه درمان است [۴]. در حقیقت، بیمه سازوکاری است که به افراد ریسک‌پذیر اجازه می‌دهد از ریسک و فقدان اطمینانی که با آن مواجه هستند، رها شوند. مردم با پرداخت حق بیمه تعیین شده، زیان و مخارج ناشی از بیماری را کاهش داده و به این ترتیب رفاه خود را افزایش می‌دهند [۳ و ۴]. سالمند شدن جمعیت ایران و افزایش هزینه‌های بازنشستگی، تقاضا برای خدمات درمانی را بسیار افزایش داده است. واضح است که تمام این موارد، متضمن هزینه برای صنعت بیمه کشور خواهد بود. بنابراین مدل‌بندی هزینه‌های درمان، نیاز اصلی صنعت بیمه کشور بوده و بیمه‌گر اجتماعی باید با رویکرد مالی و بر اساس خصوصیات هر جامعه، روند هزینه‌های درمان را به‌دقت بررسی کرده و بتواند مدل مناسبی را برای پیش‌بینی هزینه‌های درمان در آینده ارائه دهد [۵ و ۶].

توزیع‌های آمیخته

امروزه تحلیلگران داده، به لطف روش‌های محاسباتی قوی و همچنین طیف گسترده‌تری از توزیع‌های مدل‌سازی، قادر به توصیف، تخمین، پیش‌بینی و استنباط در مورد دستگاه‌های پیچیده موردعلاقه هستند. مدل‌های آمیخته تصویری جذاب از این جنبه‌ها فراهم می‌کند. اگرچه مدل‌های آمیخته درون یک خانواده از توزیع‌های پارامتری تقریب‌های سازگاری ارائه می‌کند، اما دارای چالش‌های محاسباتی بسیار پیچیده‌ای است. توزیع‌های آمیخته شامل یک تعداد متناهی یا نامتناهی از عضوهای ممکن است نوع توزیع‌ها از همدیگر متفاوت باشند که این خاصیت ویژگی‌های متفاوت داده‌ها را توصیف می‌کند. بنابراین، آنها توصیف‌های بسیار دقیق‌تری از دستگاه‌های پیچیده مانند نجوم، بوم‌شناسی، بیوانفورماتیک، علوم رایانه، اقتصاد، مهندسی، رباتیک و زیست‌شناسی ارائه می‌کنند. توزیع‌های آمیخته اهمیت زیادی در مدل‌بندی پدیده‌های تصادفی دارد. در دهه گذشته، مدل‌های آمیخته متناهی بسیار مورد توجه قرار گرفته است، زیرا انعطاف‌پذیری مدل‌های آمیخته بسیار بالاست و راه‌های ساده‌ای برای مدل‌بندی داده‌های که شکل توزیع آنها نامشخص است ارائه می‌کند. باتوجه به گفته باکس (۱۹۷۶) [۱۳] که «همه مدل‌ها اشتباه هستند، اما برخی مفید هستند»، مدل‌های آمیخته نیز بسیار مفید هستند، در صورتی که عمدتاً درست نیستند، اما می‌توانند به‌عنوان جایگزین‌های مفیدی برای برآورد چگالی، خوشه‌بندی، ناهمگنی غیرقابل مشاهده و غیره باشند. مدل‌سازی آمیخته را می‌توان برای طیف گسترده‌ای از داده‌ها مانند: تک‌متغیره یا چندمتغیره، پیوسته یا طبقه‌ای، مقطعی، سری زمانی و شبکه‌ها استفاده کرد. معمولاً در مواردی که جامعه آماری ناهمگن و ترکیبی از چند زیرجامعه باشد توزیع‌های آمیخته بسیار کاربرد دارد. مدل‌های آمیخته بیش از ۱۵۰ سال است که در بسیاری از شاخه‌های مدل‌سازی آماری به‌عنوان ابزاری همه‌کاره و چندوجهی برای توانمند ساختن مجموعه توزیع‌های احتمال برای مدل‌بندی داده‌ها مورد استفاده قرار گرفته است. از نظر تاریخی، مفهوم توزیع‌های آمیخته برای اولین بار توسط نیوکام (۱۸۸۶) به‌عنوان الگویی برای محیط‌های خارجی یافت شد و بعداً توسط کارل پیرسون و والتر ولدون مورد مطالعه قرار گرفت. توزیع‌های آمیخته به دلیل انعطاف‌پذیر بودن در توصیف بسیاری از پدیده‌های تصادفی مفید هستند. کاربرد توزیع‌های آمیخته معمولاً در مواردی

است که جامعه آماری ناهمگن و ترکیبی از چند زیرجامعه باشد. در چنین حالتی مشاهداتی که از این جامعه به دست می‌آیند، می‌توانند با احتمال معینی به هر یک از این زیرجامعه‌ها تعلق داشته باشند. برای مدل‌بندی داده‌های بیمه‌ای، باتوجه به پارامتر ریسک آنها که در کلاس‌های مختلفی قرار می‌گیرند، استفاده از توزیع‌های پارامتری شناخته شده کار مشکل و پیچیده‌ای است زیرا با یک جامعه یکسان و همگن سروکار نداریم و پارامتر ریسک برای زیرجامعه‌های مختلف متفاوت است. به همین دلیل جامعه را خوشه‌بندی کرده و برای هر یک از این خوشه‌ها می‌توان توزیع‌های متفاوتی در نظر گرفت. در حوزه بیمه توزیع‌های آماری متفاوتی برای بررسی زیان مورد استفاده قرار می‌گیرد که یکی از آنها توزیع‌های آمیخته است. توزیع‌های آمیخته به دلیل انعطاف‌پذیری بالای آنها در مدل‌بندی داده‌ها در حوزه بیمه مورد توجه بیم‌سازها قرار گرفته و به‌طور گسترده‌ای برای مدل‌بندی داده‌های به‌کار می‌رود که مشاهدات از گروه‌های مختلف آمده باشد.

تعریف: فرض کنید که F خانواده از توابع توزیع باشد. می‌گوییم Y یک متغیر تصادفی با یک توزیع آمیخته متناهی است. اگر تابع توزیع آن به فرم زیر باشد:

$$F(y) = \sum_{j=1}^k w_j F_j(y|\theta_j) \quad (1)$$

که $F_j \in F$ تابع توزیع از k جامعه قابل تشخیص است و $w = (w_1, w_2, \dots, w_k)$ ، $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ و $0 \leq w_j \leq 1$ و وزن‌های آمیختگی به‌طوری که $\sum_{j=1}^k w_j = 1$.

رگرسیون لجستیک

رگرسیون یک رابطه آماری بین دو یا بیش از دو متغیر است که هر تغییری در مقدار متغیر مستقل این معادله، متغیر وابسته را تحت تأثیر قرار می‌دهد. اما وقتی که متغیر پاسخ (وابسته) کمی نباشد و یک متغیر گسسته باشد مثلاً بله یا خیر، درست یا غلط، رأی‌داده یا نداده، کالای خریداری شده یا نه؟ از رگرسیون لجستیک استفاده می‌شود. یکی از روش‌های دسته‌بندی در مبحث یادگیری ماشین رگرسیون لجستیک است. در این روش رگرسیون، از مفهوم و شیوه محاسبه نسبت بخت استفاده می‌شود. فرض کنید فرم مدل

نکته حائز اهمیت در اینجا این است که داده‌ها را با بایستی بر اساس مقیاسی یکسان در معیارهای اندازه‌گیری فاصله به کار برد (داده‌ها را نرمال سازی کرد). محاسبه فاصله باید براساس داده‌های بدون مقیاس صورت بگیرد تا ویژگی یا متغیرهای خاص با واحد اندازه‌گیری بزرگ، باعث اریبی در مقدار فاصله بین نقطه‌ها نشود، در غیر این صورت ممکن است نتایج خوشه‌بندی بستگی زیادی به متغیری داشته باشد که دارای مقدارهای بزرگ‌تر و در نتیجه واریانس بزرگ‌تری خواهد داشت.

روش بررسی

برای مدل‌بندی داده‌های بیمه‌ای، با توجه به پارامتر ریسک بیمه شده که در کلاس‌های مختلفی قرار می‌گیرند، استفاده از توزیع‌های پارامتری تک‌عضوی مناسب نیست زیرا پارامتر ریسک برای کلاس‌های مختلف متفاوت است و با جامعه همگن و یکسانی سروکار نداریم. توزیع‌های آمیخته به دلیل انعطاف‌پذیری بالای آنها، در توصیف بسیاری از پدیده‌های تصادفی کاربرد دارند. معمولاً در مواردی که جامعه ناهمگن است برای مدل‌بندی داده‌ها، توزیع‌های آمیخته بسیار مفید هستند. در این پژوهش، با توجه به ناهمگنی موجود در خسارت‌های دندانپزشکی که به دلیل عواملی چون سن، جنسیت، شغل و غیره است، نمی‌توان با توزیع‌های پارامتری تک‌عضوی مدل‌بندی کرد، به همین خاطر با استفاده از توزیع‌های آمیخته متناهی داده‌ها را مدل‌بندی کرده و با توجه به مدل برازش داده شده، حق بیمه باورمندی برای سال آینده محاسبه می‌شود.

باورمندی

یکی از اساسی‌ترین مسائل در انواع بیمه‌ها، تعیین حق بیمه آنهاست. نظریه باورمندی در بیمه، به‌عنوان ابزاری برای تعیین و تعدیل حق بیمه به‌کار می‌رود. در واقع باورمندی حق بیمه موجود را با خسارت‌های گذشته ترکیب کرده و حق بیمه تعدیل شده را ارائه می‌دهد. اما در مدل‌های بیمه هرچقدر حجم نمونه بیشتر باشد، باورمندی افزایش می‌یابد. چون باورمندی به الگوهای می‌پردازد که پاسخ‌هایی مطابق با محیط و پیرامون خود ارائه می‌دهد، بیشتر در علوم اکچوئری نفوذ کرده است. نظریه باورمندی مجموعه‌ای از ابزار کمی است که به بیمه‌گر اجازه می‌دهد حوادث آتی را برای یک ریسک یا گروهی از ریسک‌ها نرخ‌بندی کند. به این معنا که

رگرسیون خطی ساده به صورت

$Z = \beta_0 + \beta_1 X + \epsilon$ باشد همان‌طور که دیده می‌شود این رابطه، معادله یک خط است که جمله خطا یا همان ϵ به آن اضافه شده است. در این حالت اگر Z مقدار برآورد برای متغیر وابسته باشد، می‌توان آن را میانگین مشاهدات برای متغیر وابسته به ازای مقدار ثابت متغیر مستقل در نظر گرفت. اگر مقدار متغیر وابسته Z ، باینری (دو وضعیت) و شامل صفر و ۱ باشد مشخص است که دارای توزیع برنولی است و امید ریاضی آن به صورت $z = E(Z|X=x) = P(Z=1|X=x) = p(x)$ محاسبه می‌شود. برای مشخص کردن مدل رابطه بین متغیر وابسته و مستقل به‌جای رابطه خطی، به تابعی احتیاج داریم که در حدود ۰ تا ۱ تغییر کند:

$$p(x) = \hat{z} = E(Z = 1|X = x) = \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}} \quad (2)$$

به‌منظور برآورد پارامترهای این مدل، می‌توان از تبدیل لوجیت استفاده کرد.

$$g(x) = \ln\left(\frac{p(x)}{1-p(x)}\right) = \ln(e^{b_0 + b_1 x}) = b_0 + b_1 x \quad (3)$$

با استفاده از تابع درست‌نمایی و حداکثرسازی آن می‌توان پارامترها را برآورد کرد. با این کار به یک دستگاه معادلات می‌رسیم که متأسفانه برای حل آن روش تحلیلی وجود ندارد و باید به کمک روش‌های عددی برآورد را انجام داد. در این پژوهش هدف این است که بیمه شده‌ها را به دو خوشه پریسک و کم‌ریسک تقسیم کرده و به هر یک از این خوشه‌ها یک توزیع مناسب برازش داده شود، یکی از روش‌های خوشه‌بندی K-means است. مسئله مهمی که در این رابطه به وجود می‌آید، نرمال‌سازی داده‌ها در خوشه‌بندی است زیرا باید ویژگی‌ها در محاسبه فاصله بدون مقیاس باشند تا بزرگی واحد اندازه‌گیری هر بُعد باعث اریبی مقدار تابع فاصله به سمت یک ویژگی نشود. برای اندازه‌گیری شباهت یا فاصله بین مشاهدات از فاصله اقلیدسی یا منهن استفاده کرد که در زیر نحوه محاسبه بین دو نقطه p بعدی x و y ارائه شده است:

$$x = (x_1, \dots, x_p) \quad y = (y_1, \dots, y_p)$$

$$D_{Euc}(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2} \quad (4)$$

$$D_{Man}(x, y) = \sum_{i=1}^p |x_i - y_i| \quad (5)$$

محاسبه شود. بنابراین تابع درست‌نمایی برای توزیع آمیخته که در رابطه ۱ ارائه شد به صورت زیر خواهد بود:

$$\mathcal{L}(\theta, \mathbf{p} | \mathbf{y}) = \prod_{i=1}^n \sum_{j=1}^k w_j f(y_i | \theta_j) \quad (7)$$

همان‌طور که ملاحظه می‌شود بسط تابع درست‌نمایی فوق، k^n جمله خواهد داشت، به همین دلیل محاسبه تابع درست‌نمایی، تابع توزیع پسین و برآوردگرهای بیز، حتی برای حالتی که توزیع پیشین مزدوج برای (θ, \mathbf{p}) در نظر گرفته شود، کار بسیار مشکلی است. برای حل این مشکل یک سری نمادهای جدید در جدول ۱ ارائه شده که با استفاده از این نمادهای جدید فرم ساده‌تری برای تابع درست‌نمایی فوق پیدا کرده و با توجه به آن می‌توان برآوردگر باورمندی را برای توزیع‌های آمیخته محاسبه کرد.

قضیه ۱) فرض کنید که نمونه تصادفی Y_1, Y_2, \dots, Y_n از توزیع آمیخته k عضوی رابطه ۱ آمده است. تابع درست‌نمایی برای توزیع آمیخته رابطه ۱ یک فرمول بازگشتی به صورت زیر خواهد بود:

$$L_k(\psi | Y) = \sum_{i=1}^n w_k^i (1 - w_k)^{n-i} \sum_{r=1}^{(i)} L_{k-1}(\psi(-k) | \widetilde{Y}_{B_{ir}^c}) \prod_{k \in B_{ir}} f_k(y | \theta_k) \quad (8)$$

که $L_{k-1}(\psi(-k) | \widetilde{Y}_{B_{ir}^c})$ تابع درست‌نمایی برای تابع چگالی

$$g^*(y | \psi(-k)) = \frac{w_1}{1-w_k} f_1(y | \theta_1) + \frac{w_2}{1-w_k} f_2(y | \theta_2) + \dots + \frac{w_{k-1}}{1-w_k} f_{k-1}(y | \theta_{k-1})$$

خواهد بود.

هدف ما در این پژوهش این است که بیمه شده‌ها را در دو دسته پرریسک و کم‌ریسک تقسیم کنیم. برای این منظور فرض بر این است که متغیر پاسخ Z_i دو مقدار ۰ و ۱ را اختیار می‌کند به این معنی که $Z_i=1$ یعنی بیمه شده i ام متعلق به جامعه پرریسک و $Z_i=0$ یعنی بیمه شده i ام متعلق به جامعه کم‌ریسک خواهد بود. در تحلیل عوامل مؤثر بر پرریسک یا کم‌ریسک بودن بیمه شده‌ها، متغیرهای همچون جنسیت، شغل، سن و غیره ممکن است تأثیرگذار باشد، برای این منظور اگر $X_i=(X_{i1}, X_{i2}, \dots, X_{ik})$ نشان‌دهنده بردار متغیرهای مستقل باشد. آنگاه احتمال اینکه بیمه شده به جامعه پرریسک متعلق باشد با استفاده از رگرسیون لجستیک به صورت زیر قابل محاسبه است:

$$w_1 = w_1(\mathbf{x}, \boldsymbol{\alpha}) = P(Z_i = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp\{\beta_{.1} + \boldsymbol{\beta}'_1 \mathbf{x}\}}{1 + \exp\{\beta_{.1} + \boldsymbol{\beta}'_1 \mathbf{x}\}} \quad (9)$$

حق بیمه‌های آتی را بر اساس حوادث گذشته تعدیل می‌کند. اگر حوادث یک بیمه‌گذار یا گروهی از بیمه‌گذاران در دسترس باشد، دانش آماری ما ممکن است ما را به استفاده از میانگین نمونه‌ای یا برخی برآوردگرهای ناریب دیگر متقاعد کند. اما نظریه باورمندی به ما می‌گوید بهینه آن است که تنها قسمتی از وزن را به این حوادث بدهیم و بقیه وزن را به برآوردگرهایی که از اطلاعات دیگری به دست می‌آیند داده شود. فرض کنید که Y_1, Y_2, \dots, Y_n شدت خسارت‌های یک فرد بیمه شده در n سال گذشته باشد. شرکت بیمه معمولاً یک نرخ بیمه‌ای μ برای کل داشتن برآورد کرده است. باورمندی یک دستورالعمل برای ترکیب کردن داده‌های گذشته و μ است. بر اساس دیدگاه باورمندی، یک برآوردگر سازگار برای حق بیمه به صورت زیر قابل محاسبه است:

$$\pi_c = \xi Y + (1 - \xi) \mu \quad (6)$$

که $0 < \xi < 1$ عامل باورمندی، Y میانگین مشاهدات گذشته و μ میانگین توزیع پیشین است. از دیدگاه بیمه‌گر، اگر گذشته بیمه‌گذار پایدار باشد، یعنی واریانس خسارت‌ها کم باشد، به نظر می‌رسد که از Y برای پیش‌بینی نتایج دوره بعد استفاده شود. همچنین اگر گذشته افراد دارای تغییرپذیری بیشتری باشد، بهتر است که برای پیش‌بینی نتایج دوره بعد، کمتر از Y استفاده شود و وزن میانگین توزیع پیشین μ بیشتر باشد. فرض کنید که شرکت بیمه‌ای با توجه به مشاهدات مربوط به خسارت‌های n سال گذشته یک بیمه شده با پارامتر ریسک Θ می‌خواهد حق بیمه انفرادی را برای سال $n+1$ ام به دست آورد، در این نوع حق بیمه با توجه به پارامتر ریسک هر فرد Θ ، برای هر فرد یک حق بیمه خاص به دست می‌آید. حق بیمه انفرادی (خالص) با پارامتر ریسک Θ برابر با $[\mu(\Theta) = E[Y_{n+1} | \Theta]]$ است. این نوع حق بیمه را حق بیمه عادلانه نیز می‌نامند و بهترین حق بیمه ممکن است، اما مشکل آنجاست که θ و به دنبال آن $\mu(\theta)$ نامعلوم است. اگر $\mu(\theta)$ یک برآوردگر دلخواه برای $\mu(\theta)$ و $\mu(\theta)$ $[E[\mu(\theta) | Y]]$ برآوردگر بیز برای $\mu(\theta)$ (حق بیمه انفرادی) باشد می‌توان نشان داد که برآوردگر بیز $\mu(\theta)$ تحت تابع زیان درجه دوم نسبت به هر برآوردگر دلخواه $\mu(\theta)$ بهتر است. اما در بسیاری از موارد محاسبه برآوردگر بیز کاری بسیار مشکل و دشوار است به همین خاطر از برآوردگر باورمندی استفاده می‌شود. برای محاسبه برآوردگر باورمندی، ابتدا باید تابع درست‌نمایی را محاسبه کرد که بتوان توزیع پسین را به دست آورد و در نهایت برآوردگر باورمندی

محاسبه است:

$$E_k(Y_{n+1}|Y_1, Y_2, \dots, Y_n) = E(E(Y_{n+1}|\theta_k, \psi(-k))|\bar{Y} \in \cup_{i=1}^k PoP_i) = \sum_{i=1}^n \sum_{r=1}^n E(Y_{n+1}|\bar{Y}_{B_{ir}} \in PoP_k, \bar{Y}_{B_{ir}}^c \notin PoP_k) P(\bar{Y}_{B_{ir}} \in PoP_k, \bar{Y}_{B_{ir}}^c \notin PoP_k) \left| \bar{Y} \in \cup_{i=1}^k PoP_i \right. \\ = \sum_{i=1}^n \sum_{r=1}^n w_k^i (1-w_k)^{n-i} E(E(Y_{n+1}|\theta_k, \psi(-k))|\bar{Y}_{B_{ir}} \in PoP_k, \bar{Y}_{B_{ir}}^c \notin PoP_k) \\ = \sum_{i=1}^n \sum_{r=1}^n w_k^i (1-w_k)^{n-i} [w_k E(\mu(\theta_k)|\bar{Y}_{B_{ir}} \in PoP_k) + (1-w_k) E(\mu(\psi(-k))|\bar{Y}_{B_{ir}}^c \notin PoP_k)] = \\ = \sum_{i=1}^n \sum_{r=1}^n w_k^i (1-w_k)^{n-i} [w_k E(Y_{n+1}|\bar{Y}_{B_{ir}} \in PoP_k) + (1-w_k) E_{k-1}(Y_{n+1}|\bar{Y}_{B_{ir}}^c \notin PoP_k)] \quad (14)$$

مراحل محاسبه برآوردگر باورمندی برای توزیع‌های آمیخته در

یافته‌ها

داده‌های مورد استفاده در این مطالعه شامل خسارت‌های مربوط به مجموعه بیمه‌نامه‌های درمان گروهی صادره طی دو سال از اواسط سال ۱۳۹۷ تا اواسط ۱۳۹۹ است. به طوری که در هر قرارداد بنا بر نوع قرارداد که شامل متغیرهای سقف تعهد و فرانشیز است، توسط بیمه‌شدگان که اطلاعات سن، جنسیت (زن و مرد)، وضعیت (اصلی و تحت تکفل)، استان محل زندگی و گروه شغلی (اتحادیه/اصناف/خانه کارگر/تعاونی، دانشگاه‌های علوم پزشکی و مراکز درمانی، وزارت کشور، استانداری‌ها، شهرداری‌ها، بخشداری و نظایر آن، محیط‌های اداری (به‌غیر از طبقه ۲ و ۳) و محیط‌های تولیدی و صنعتی) آنها در اختیار است، خسارت‌هایی در خصوص خدمات مختلف ثبت شده، اما هدف در این پژوهش مدل‌بندی و پیش‌بینی خسارت‌های مربوط به دندانپزشکی است. با اینکه ۴۴ درصد از خسارت‌های دندانپزشکی مربوط به آقايان و ۵۶ درصد از آنها مربوط به خانم‌هاست، برای بررسی تفاوت بین خسارت‌های دندانپزشکی آقایان و خانم‌ها با استفاده از روش تحلیل واریانس (ANOVA)، پی-مقدار برابر با 2×10^{-16} شده که در سطح $\alpha = 0.05$ فرض صفر رد می‌شود. بدین معنا که بین خسارت‌های دندانپزشکی آقایان و خانم‌ها به‌طور معناداری تفاوت وجود دارد. به‌منظور تجزیه و تحلیل متغیرها، ابتدا آمار توصیفی از داده‌ها ارائه می‌شود که نتایج آن در جدول ۳ آمده است. در این جدول مقادیر حداکثر، حداقل، میانگین، میانه و انحراف معیار هر یک از متغیرهای جنسیت، شغل و نسبت بیمه شده برای خسارت‌های مربوط به دندانپزشکی است.

جدول ۳. متغیرهای توصیفی پژوهش

با نگاهی به شاخص سلامت دندان گروه‌های سنی مختلف در ایران،

بنابراین اگر توزیع برازش داده شده به جامعه کم‌ریسک $f(x|\theta_1)$ و به جامعه پرریسک $f(x|\theta_2)$ باشد، آنگاه توزیع آمیخته برازش داده شده به صورت زیر خواهد بود:

$$f(x|\theta, \theta_r) = \left(1 - \frac{\exp\{\beta_{-1} + \beta'x\}}{1 + \exp\{\beta_{-1} + \beta'x\}}\right) f(x|\theta_1) + \frac{\exp\{\beta_{-1} + \beta'x\}}{1 + \exp\{\beta_{-1} + \beta'x\}} f(x|\theta_r) \quad (10)$$

قضیه ۲) فرض کنید که نمونه تصادفی Y_1, Y_2, \dots, Y_n از توزیع آمیخته k عضوی رابطه ۱ آمده است به طوری که وزن‌های آمیختگی در این توزیع با استفاده از رگرسیون لجستیک به صورت $w_k = w_k(x, \alpha) = \frac{\exp\{\beta_{-k} + \beta'x\}}{1 + \sum_{j=1}^k \exp\{\beta_{-j} + \beta'x\}}$ به دست آمده باشد. در این حالت برآوردگر باورمندی به صورت زیر قابل محاسبه است:

$$E_k(Y_{n+1}|Y_1, Y_2, \dots, Y_n) = \sum_{i=1}^n \sum_{r=1}^n w_k^i (1-w_k)^{n-i} [w_k E(Y_{n+1}|\bar{Y}_{B_{ir}} \in PoP_k) + (1-w_k) E_{k-1}(Y_{n+1}|\bar{Y}_{B_{ir}}^c \notin PoP_k)] \quad (11)$$

اثبات) ابتدا باید تابع توزیع متغیر تصادفی Y_1 محاسبه شود

$$G_{Y_1|\psi}(t) = P(Y_1 \leq t|\psi) = \sum_{j=1}^k P(Y_1 \leq t|Y_1 \in PoP_j, \psi) P(Y_1 \in PoP_j|\psi) = \sum_{j=1}^k w_j F_j(t) \quad (12)$$

باتوجه به تابع توزیع فوق، تابع چگالی به صورت $g(y|\psi) = w_1 f_1(y|\theta_1) + w_2 f_2(y|\theta_2) + \dots + w_k f_k(y|\theta_k)$ می‌توان تابع چگالی فوق را به صورت زیر بازنویسی کرد:

$$g(y|\psi) = w_k f_k(x|\theta_k) + (1-w_k) f^*(y|\psi(-k)) \quad (13)$$

که

$$f^*(y|\psi(-k)) = \frac{w_1}{1-w_k} f_1(y|\theta_1) + \frac{w_2}{1-w_k} f_2(y|\theta_2) + \dots + \frac{w_{k-1}}{1-w_k} f_{k-1}(y|\theta_{k-1})$$

باتوجه به تابع چگالی فوق برآوردگر باورمندی به صورت زیر قابل

جدول ۱ | لیست علائم و نمادها

نماد	تعریف
K	تعداد عضوهای مدل آمیخته
N	تعداد مشاهدات
\bar{Y}	(Y_1, Y_2, \dots, Y_n)
S^n	$\{1, 2, \dots, n\}$
B_{ir}	R امین زیرمجموعه A عضوی از S^n
B_{ir}^c	متمم مجموعه B_{ir}
S_i^n	$= \{B_{ir}; r = 1, 2, \dots, \binom{n}{i}\}$
$\bar{Y}_{B_{ir}}$	$= (Y_{k_1}, Y_{k_2}, \dots, Y_{k_i}), k_1, k_2, \dots, k_n \in B_{ir}$
$\bar{Y}_{B_{ir}^c}$	$= (Y_1, Y_{k_2}, \dots, Y_{k_{n-i}}), k_1, k_2, \dots, k_{n-i} \in B_{ir}^c$
Ψ	$= (\theta_1, \theta_2, \dots, \theta_K)$
$\Psi(-1)$	$= (\theta_1, \theta_2, \dots, \theta_{l-1}, \theta_{l+1}, \dots, \theta_K)$
$E_R(\cdot)$	امید ریاضی برای توزیع آمیخته با k عضو

بیمه با توافق طرفین مشخص می‌شود. همچنین سقف تعهدات و پوشش‌های بیمه تکمیلی انفرادی به صورت ثابت انجام می‌شود. هدف از این مطالعه این است که مدل‌بندی خسارت‌های مربوط به پوشش دندانپزشکی افراد بیمه شده و برآوردگر باورمندی ارائه شود. در نمودار ۱ چگالی مربوط به خسارت‌های دندانپزشکی که بر عدد ۱۰۰۰۰ تقسیم شده‌اند، رسم شده است.

باتوجه به نمودار چگالی مشخص است که یک ناهمگنی در خسارت‌ها وجود دارد که این ناهمگنی به دلیل مشخصه‌های مانند سن، جنسیت، شغل، موقعیت جغرافیایی، نسبت با بیمه شده و... است. بنابراین هدف این است که بیمه‌شده‌ها را باتوجه به مشخصه‌هایی که دارند به دو زیرجامعه پریسک و کم‌ریسک تقسیم کرد. برای این منظور با استفاده از روش K-means مرکز هر یک از خوشه‌ها به دست آمده است که یک نمایش گرافیکی دومتغیره از خوشه‌بندی داده‌ها در شکل ۱ است. در این شکل، داده‌ها به صورت نقاط در یک نمودار دومتغیره و خوشه‌ها به صورت بیضی‌هایی با اندازه‌ها و اشکال مختلف نمایش داده شده است.

حال بعد از اینکه داده‌ها در خوشه‌های پریسک و کم‌ریسک تقسیم شدند باید به هر یک از این زیرجامعه‌ها توزیع مناسب برازش داده شود. برای این منظور ۴ توزیع نرمال، نمایی، گاما و لگ‌نرمال را به هر یک از زیرجامعه‌های پریسک و کم‌ریسک برازش داده و مقدار

می‌توان به این نتیجه رسید که شاخص سلامت دندان‌های شیری برابر با ۱۶/۵ درصد، دندان‌های دائمی در کودکان با رده سنی ۱۲ سال برابر با ۹/۲ درصد، برای گروه سنی ۳۵-۴۴ سال برابر با ۹۸/۱۳ و گروه سنی ۶۵-۷۴ سال برابر با ۲۵/۲۷ درصد است. در گروه سنی ۶۵-۷۴ سال تعداد قابل توجهی از دندان‌های پوسیده درمان شده وجود دارد. باتوجه به موارد گفته شده می‌توان فهمید که هزینه‌های دندانپزشکی یکی از پرهزینه‌ترین خدمات درمانی بوده و همه افراد حداقل چند بار در طول زندگی به دندانپزشکی مراجعه کرده‌اند. این هزینه در کشور ما، بعد از هزینه بستری شدن، در رتبه دوم قرار دارد. به‌گونه‌ای که حدود ۱۵-۲۰ درصد از کل هزینه‌های خدمات درمانی-دارویی را دربر می‌گیرد. همین موضوع اهمیت بیمه تکمیلی دندانپزشکی را نشان می‌دهد. خدمات درمانی تعهدات دندانپزشکی سازمان‌های بیمه‌گر شامل خدمات ویزیت (معاینه و تشخیص، تجویز دارو، اعزام به رادیولوژی و درخواست آزمایش)، رادیوگرافی تک‌دندان (گرافی پری اپیکال یا بایت وینگ)، کشیدن دندان قدامی، کشیدن دندان خلفی، کشیدن دندان عقل، جراحی دندان نهفته در نسج نرم، جراحی دندان نهفته در نسج سخت، جرم‌گیری بالای لثه، جرم‌گیری زیر لثه، بروساژ هر فک (نوعی جرم‌گیری)، پر کردن تک‌سطحی و ترمیم دندان است. بیمه درمان تکمیلی گروهی بین نهاد متقاضی و بیمه به صورت قراردادی منعقد می‌شود. شرایط و پوشش‌های این

جدول ۲ | به اختصار توضیح داده شده است

الگوریتم محاسبه برآوردگر باورمندی برای توزیع‌های آمیخته متناهی

ورودی: Y مجموعه داده‌های خسارت n سال گذشته یک بیمه شده خاص، مرکز خوشه‌ها
خروجی: محاسبه وزن‌های آمیختگی، برازش توزیع آمیخته، محاسبه برآوردگر باورمندی

(۱) انتخاب $\mu_1, \mu_2, \dots, \mu_k$ به عنوان مرکز خوشه‌ها

(۲) فاصله تک‌تک مشاهدات را از مرکز خوشه‌ها $l = 1, 2, \dots, k$ $i = 1, 2, \dots, n$ $(y_i - \mu_l)^2$

(۳) مشخص کردن خوشه هر یک از مشاهدات $l = 1, \dots, k$ $i = 1, 2, \dots, n$ $Z_i = \arg \min (y_i - \mu_l)^2$

(۴) محاسبه مرکز خوشه‌ها با توجه به مشاهداتی که در هر خوشه قرار می‌گیرد $\hat{\mu}_l = \frac{1}{n_l} \sum_{i=1}^n y_i I_{\{Z_i=l\}}$

(۵) تکرار مراحل ۲ تا ۴ تا زمانی که مرکز خوشه‌ها تغییر نکند

(۶) برازش توزیع مناسب به هر یک از خوشه‌ها

(۷) محاسبه وزن‌های آمیختگی w_i با استفاده از تابع لجستیک $w_i = w_i(x, \alpha) = \frac{\exp\{\beta_{0i} + \beta'_i x\}}{1 + \sum_{j=1}^{k-1} \exp\{\beta_{0j} + \beta'_j x\}}$

برازش توزیع نهایی به صورت آمیخته‌ای از توزیع‌های برازش داده شده در گام پنجم با وزن‌های آمیختگی محاسبه شده در گام ششم

(۸) محاسبه برآوردگر باورمندی با استفاده از توزیع آمیخته برازش داده شده

آماره‌های آزمون کولموگروف اسمیرنوف، اندرسون دارلینگ و معیارهای AIC و BIC با همدیگر مقایسه شده‌اند که در نهایت توزیع برازش داده شده به زیر جامعه کم‌ریسک دارای توزیع لگ‌نرمال با پارامترهای $\mu = 4195/15$ و $\sigma = 81319/0$ و توزیع برازش داده شده به زیر جامعه پرریسک دارای توزیع لگ‌نرمال با پارامترهای $\mu = 15,5946$ و $\sigma = 0,9684762$ بوده است. بنابراین توزیع برازش داده شده به خسارت‌های دندانپزشکی به صورت زیر خواهد بود:

هریک از زیر جامعه‌ها (وزن‌های آمیختگی) را برآورد کرد. همان‌طور که قبلاً اشاره شد، متغیرهای کمکی جنسیت، شغل، سن، وضعیت و موقعیت جغرافیایی از جمله عواملی بود که ممکن است در ناهمگنی داده‌ها تأثیرگذار باشد. براساس نتایج مربوط به آزمون معنادار هر یک از این متغیرها در جدول ۴ نشان داده شده که پی-مقدار برای هر یک از این متغیرها در سطح $\alpha=0.05$ معنادار است. بدین معنا که هر یک از این عوامل بر خسارت‌های دندانپزشکی تأثیرگذار خواهد بود. نتایج مربوط برآورد پارامترها و پی-مقدار در جدول ۴ آمده است. همان‌طور که در جدول زیر نشان داده شده است مقدار sig برای هر یک از این متغیرها گزارش شده است. باتوجه به اینکه مقدار sig (پی-مقدار) برای هر یک از این متغیرها کمتر از 0.05 است می‌توان گفت که هر یک از این متغیرها در میزان خسارت یا اینکه بیمه شده متعلق به جامعه پرریسک یا کم‌ریسک است تأثیرگذار خواهد بود.

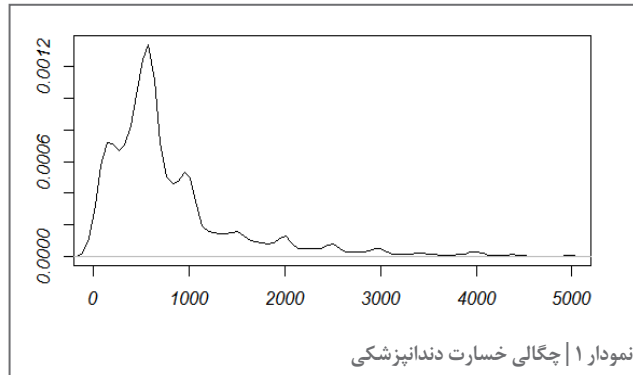
$$f(x|\theta_1, \theta_r) = \left(1 - \frac{\exp\{\beta_{\cdot 1} + \beta'x\}}{\{1 + \exp\{\beta_{\cdot 1} + \beta'x\}\}}\right) LN(15/4, 0/18) + \frac{\exp\{\beta_{\cdot 1} + \beta'x\}}{\{1 + \exp\{\beta_{\cdot 1} + \beta'x\}\}} LN(15/6, 0/96) \quad (15)$$

یا می‌توان گفت که لگاریتم خسارت‌ها دارای توزیع آمیخته متناهی از نرمال خواهد بود:

$$f(x|\theta_1, \theta_r) = \left(1 - \frac{\exp\{\beta_{\cdot 1} + \beta'x\}}{\{1 + \exp\{\beta_{\cdot 1} + \beta'x\}\}}\right) N(15/4, 0/18) + \frac{\exp\{\beta_{\cdot 1} + \beta'x\}}{\{1 + \exp\{\beta_{\cdot 1} + \beta'x\}\}} N(15/6, 0/96) \quad (16)$$

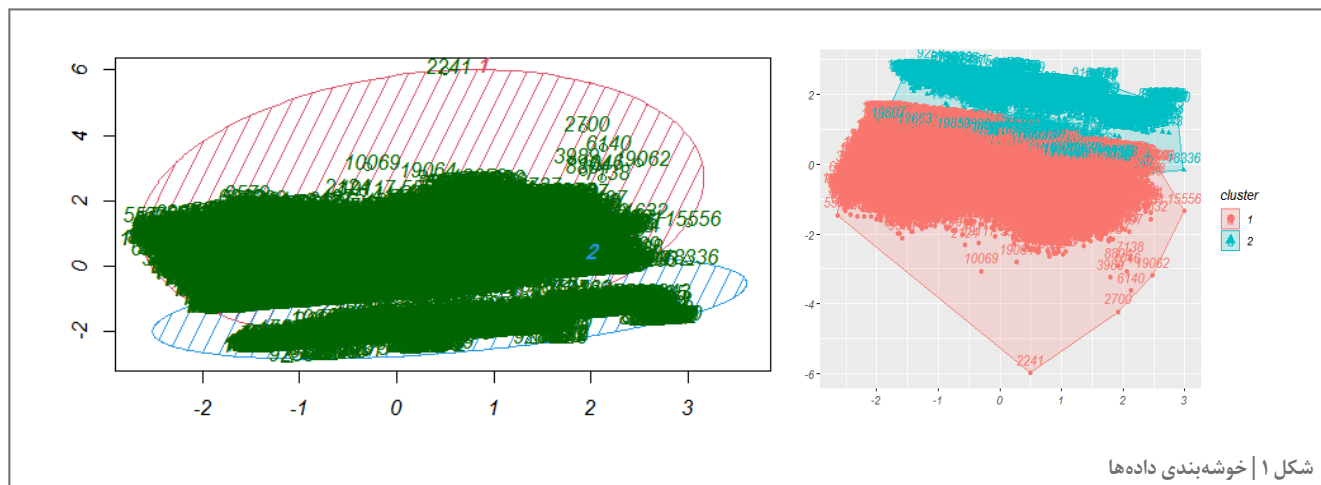
باتوجه به اینکه متغیر پاسخ Y دو مقدار صفر و یک را اختیار می‌کند (Y=1 یعنی بیمه شده متعلق به زیر جامعه پرریسک و Y=0 یعنی بیمه شده متعلق به زیر جامعه کم‌ریسک) احتمال اینکه یک بیمه شده با مشخصات $X=(X_1, X_2, X_3, X_4, X_5)$ متعلق به جامعه پرریسک باشد به صورت زیر قابل محاسبه است:

حال بایستی با استفاده از رگرسیون لجستیک احتمال تعلق به



$$P(Y = 1|X = x) = \frac{\exp\{3.0/81.06 + 4/7949X_1 - 5/16X_2 + 0.0678X_3 - 18/96X_4 + 1/1871X_5\}}{1 + \exp\{3.0/81.06 + 4/7949X_1 - 5/16X_2 + 0.0678X_3 - 18/96X_4 + 1/1871X_5\}} \quad (17)$$

به‌عنوان مثال، فرض کنید که خسارت‌های Y_1, Y_2, \dots, Y_n مربوط به یک بیمه شده با مشخصه‌های جنسیت خانم ($X_1=1$)، بیمه شده اصلی ($X_2=0$)، سن ۶۰ سال ($X_3=60$)، طبقه شغلی ۳ ($X_4=3$) در استان سمنان



بحث

ردنر و والکر [۱۴] به برخی از مشکلات برآوردگر ماکسیمم درست‌نمایی (MLE) برای توزیع‌های آمیخته متناهی اشاره کردند. اولاً نمی‌توان فرض کرد که تابع درست‌نمایی به‌طور کلی دارای کران بالاست، ممکن است واگرا باشد. ثانیاً اغلب کمینه‌های محلی بهینه زیادی برای تابع درست‌نمایی وجود دارد. مک لاجلان و باسفورد [۱۵] این مشکل را اظهار داشتند که عملکرد الگوریتم EM نسبت به مقدار اولیه پارامترهای مدل بسیار حساس است. اگر مقادیر اولیه نزدیک به مقادیر تکین در تابع درست‌نمایی باشد، الگوریتم EM ممکن است واگرا شود و سرعت همگرایی بسیار کند شود. روش دیگر برای برآورد پارامترهای مدل آمیخته، برآورد بیزی است که تابع درست‌نمایی داده‌ها با توزیع پیشین پارامترها، برای استخراج یک استنباط، ترکیب می‌شوند. یکی از مزایای روش بیزی نسبت به روش MLE، استفاده از اطلاعات قبلی است، این مزیت باعث می‌شود که با تعداد کمی از داده‌ها استنباط همچنان قابل انجام باشد. روش MLE، به‌ویژه برای مدل‌هایی که پارامترهای زیادی دارند، وقتی مجموعه داده کوچک است، برآورد مناسب و خوبی نیست. با این حال، در موقعیت‌های خاص، برآورد بیزی پارامترهای یک توزیع آمیخته از محاسبات طولانی برای انجام استنتاج برخوردار است، زیرا انتگرال‌ها باید بر روی یک فضای چندبعدی غیرمستقیم بالقوه انجام شوند. در الگوی بیزی برگر [۱۶]، بساگ و همکاران [۱۷] بیانیه‌های احتمال در مورد پارامترهای نامشخص، توزیع پیشین یا نظر نخبگان در تحلیل‌ها و توصیف‌های سلسله‌مراتبی از مدل‌های محلی - موضعی و سراسر بیان شده است. همچنین در زمینه مدل‌های آمیخته به‌صورت محلی - موضعی نیز می‌توان به مطالعه پاینده و همکاران [۱۸] اشاره کرد که ساختار پیچیده از یک مدل آمیخته را با استفاده از متغیرهای پنهان به یک ساختار ساده‌تر تجزیه کرد. همان‌طور که گفته شد، از زمان پیدایش الگوریتم ML، EM، تاکنون متداول‌ترین روش برای برآزش توزیع‌های آمیخته بوده است. کاربرد الگوریتم EM برای محاسبه برآوردهای ML مدل‌های آمیخته پارامتری زمانی کاربرد دارد

$(X_3=16)$ باشد که از چگالی برآزش داده شده در رابطه (۴) با پیشین‌های $N(\mu_k, b_k^2)$ $k=1,2$ پیروی می‌کند. با توجه به این مفروضات این بیمه شده با احتمال $0.8567543 = P(Y=1 | X=x)$ متعلق به جامعه پرریسک و با احتمال $0.1432457 = P(Y=0 | X=x)$ متعلق به جامعه کم‌ریسک خواهد بود. بنابراین لگاریتم خسارت‌های این بیمه شده توزیع آمیخته‌ای به‌صورت:

$f(y) = 1432457N(15/4195, 0.11319) + 0.856743N15$ خواهد بود. حال برای محاسبه برآوردگر باورمندی فرض کنید که $0.1432457, 0.15/0.0515, 0.15/69760, 0.13/92823, 0.16/19502, 0.14/75258, 0.14/49716, 0.16/84823, 0.16/0.3626, 0.16/54005$ لگاریتم خسارت‌های ده سال گذشته این بیمه شده باشد و پارامترهای توزیع پیشین به این صورت $5/b_2=0.5, \mu_1=8, \mu_2=9, b_1=0$ باشد. با توجه به این مفروضات برآوردگر باورمندی به‌صورت زیر قابل محاسبه است:

$$\theta_{\tau} | \tilde{Y}_{B_{ir}} \in PoP_{\tau} \sim N \left(\frac{b_{\tau}^{\tau} \sum_{k \in B_{ir}} x_k + \mu_{\tau} \sigma_{\tau}^{\tau}}{i b_{\tau}^{\tau} + \sigma_{\tau}^{\tau}}, \frac{\sigma_{\tau}^{\tau} b_{\tau}^{\tau}}{i b_{\tau}^{\tau} + \sigma_{\tau}^{\tau}} \right) \quad (18)$$

$$\theta_{\tau} | \tilde{Y}_{B_{ir}^c} \in PoP_{\tau} \sim N \left(\frac{b_{\tau}^{\tau} \sum_{k \in B_{ir}^c} x_k + \mu_{\tau} \sigma_{\tau}^{\tau}}{(n-i) b_{\tau}^{\tau} + \sigma_{\tau}^{\tau}}, \frac{\sigma_{\tau}^{\tau} b_{\tau}^{\tau}}{(n-i) b_{\tau}^{\tau} + \sigma_{\tau}^{\tau}} \right) \quad (19)$$

$$E_{\tau}(Y_{n+1} | Y_1, Y_2, \dots, Y_n) = \sum_{i=0}^n \sum_{\tau=1}^{\binom{n}{i}} \cdot / \lambda \delta \epsilon^i \cdot / \lambda \epsilon^{n-i} [\cdot / \lambda \delta \epsilon E_{\tau}(Y_{n+1} | \tilde{Y}_{B_{ir}} \in PoP_{\tau}) + \cdot / \lambda \epsilon E_{\tau}(Y_{n+1} | \tilde{Y}_{B_{ir}^c} \notin PoP_{\tau})] \quad (20)$$

همان‌طور که در بالا اشاره شد پیش‌بینی مربوط به خسارت سال $1+\pi$ ام برای شخصی با مشخصات جنسیت خانم که به‌عنوان بیمه شده اصلی بودند در ۶۰ سالگی که در طبقه شغلی ۳ و استان سمنان قرار داشتند برابر $15/39324$ است. در نتیجه عواملی مانند سن، جنسیت، موقعیت جغرافیایی، طبقه شغلی و وضعیت بیمه شده در پیش‌بینی خسارت در نظر گرفته شدند.

جدول ۳ | برآورد پارامترهای رگرسیون لجستیک (یافته‌های پژوهش)

متغیرها	عرض از مبدأ	جنسیت	وضعیت	سن	شغل	موقعیت جغرافیایی
برآورد پارامتر	۳۰/۸۱	۴/۷۹	-۵/۱۶	۰/۰۶۷۸	-۱۸/۹۶	۱/۱۸
پی مقدار	$2 * 10^{-16}$	$2 * 10^{-16}$	$1/25 * 10^{-14}$	$2.91 * 10^{-7}$	$2 * 10^{-16}$	$2 * 10^{-16}$
Odd ratio		$1/2 * 10^2$	$5/83 * 10^{-3}$	$1/0.7 * 10^0$	$5/71 * 10^{-9}$	$3/28 * 10^0$

خسارت‌های مربوط به مجموعه بیمه‌نامه‌های درمان گروهی صادره طی دو سال از اواسط سال ۱۳۹۷ تا اواسط ۱۳۹۹ شرکت بیمه دی مدل‌بندی شد که برای این منظور باتوجه به ناهمگنی موجود در داده‌های خسارت مربوط به بیمه درمان تکمیلی، بیمه‌شده‌ها را در دو دسته پرریسک و کم‌ریسک تقسیم کرده و به هر یک از این زیرجمعه‌ها یک توزیع مناسب برازش داده شد و توزیع کلی خسارت‌ها یک توزیع آمیخته از این توزیع‌ها در نظر گرفته شد که وزن‌های آمیختگی این توزیع را یک تابع لجستیک از متغیرهای مربوط به هر بیمه‌شده در نظر گرفته و در نهایت باتوجه به مدل برازش داده شده برآوردگر، باورمندی برای خسارت‌های دندانپزشکی محاسبه شد.

با استفاده از نتایج این مطالعه می‌توان عوامل مهم و تأثیرگذار مانند سن، جنسیت، طبقه شغلی و... را به مدل اضافه کرد و پیش‌بینی خسارت سال آینده را بر اساس این فاکتورها انجام داد. این پژوهش با ارائه مدل‌های ریاضی، مشکل مربوط به مدل‌بندی خسارت‌ها را مرتفع ساخته و با استفاده از رویکرد باورمندی خسارت‌های سال آینده مربوط به دندانپزشکی پیش‌بینی شد. بنابراین نتایج این مطالعه می‌تواند به بهبود بیمه درمان کمک کرده و مشکلات این حوزه را تا حد زیادی برطرف سازد. به‌علاوه باتوجه به این مدل، برآورد دقیق‌تری از هزینه‌های درمان داشت و حق بیمه منصفانه‌ای برای هر بیمه‌شده باتوجه به پارامترهای ریسک آن محاسبه کرد. به‌این ترتیب علاوه بر اینکه محصول طراحی شده، جذابیت بیشتری برای بیمه‌گذار خواهد داشت و با قیمت مناسب ارائه خواهد شد، شرکت بیمه نیز متحمل زیان کمتری خواهد شد.

تشکر و قدردانی: مقاله حاضر از پایان نامه در مقطع دکترای تخصصی استخراج شده است.

تأییدیه اخلاق: مقاله حاضر با کد IR.SBU.REC.1400.262 مورد تأیید کمیته اخلاق دانشگاه شهید بهشتی قرار گرفته است.

تضاد منافع: کلیه نویسندگان هیچ‌گونه تضاد منافی را اعلام نکردند.

سهم نویسندگان: سهم نویسندگان در این مقاله یکسان است.

منابع مالی: هیچ‌گونه حمایت مالی از پژوهش حاضر صورت نپذیرفته است.

References

1. World Health Organization. World report on ageing and health. World Health Organization; 2015 Oct 22.
2. Sadeghu R. Population and development in Iran: Dimensions

که داده‌ها ناکامل باشند. استفاده از الگوریتم MM توسط هانتز و لانگ [۱۹] و لانگ [۲۰] و همچنین در مطالعه نگوین و مک‌لاچان [۲۱] برای محاسبه برآوردهای ML در توزیع‌های آمیخته، در حالتی که محاسبه گام E در الگوریتم EM مشکل باشد، بسیار مورد توجه قرار گرفته است. همانند الگوریتم EM، برآورد بیزی با استفاده از روش‌های MCMC مبتنی بر کار دمپستر [۲۲] است که اظهار داشت یک مدل آمیخته متناهی را چگونه می‌توان بر اساس داده‌های ناکامل یا داده‌های گم شده ارزیابی کرد. همان‌طور که قبلاً بحث شد محاسبه تابع درست‌نمایی برای توزیع‌های آمیخته برحسب k^n جمله خواهد شد که محاسبه آن کار وقت‌گیر و بسیار مشکلی است. پاینده و ساکی زاده [۲۳] برای تابع درست‌نمایی توزیع‌های آمیخته یک تقریب ارائه کردند که برای حالت‌هایی که تعداد آمیخته‌ها زیاد باشد مناسب نیست. تاکنون برای تابع درست‌نمایی در مدل‌های آمیخته فرم بسته و دقیقی ارائه نشده بود. در این مطالعه با اعمال شرایطی جدید فرم ساده‌تر و بسته‌ای برای تابع درست‌نمایی ارائه شد که این فرم بسته باعث می‌شود به میزان قابل توجهی پیچیدگی محاسبات کاهش پیدا کند و همچنین برای محاسبه برآوردگر باورمندی عوامل مهمی که در پیش‌بینی خسارت تأثیرگذار بودند در نظر گرفته شد. لازم به ذکر است با استفاده از رویکرد به کار گرفته شده در این مطالعه می‌توان هر عامل دیگری را که در محاسبه برآوردگر باورمندی تأثیرگذار باشد، بدون پیچیدگی خاصی به مدل اضافه کرد و اثر آن را کنترل کرد.

نتیجه‌گیری

از مهم‌ترین وظایف شرکت‌های بیمه، پیش‌بینی میزان و تعداد خسارت‌ها در سال آتی است. روش‌های مختلفی مانند سری‌های زمانی، روش بیزی و نظریه باورمندی برای انجام چنین پیش‌بینی‌های وجود دارد. یکی از مزیت‌های روش باورمندی نسبت به بقیه روش‌ها، استفاده از اطلاعات قبلی است. این مزیت باعث می‌شود که با تعداد کمی از داده‌ها استنباط همچنان قابل انجام باشد. اما مسئله‌ای که در این زمینه وجود دارد این است که محاسبه تابع درست‌نمایی و برآوردگر باورمندی برای توزیع‌های آمیخته کاری مشکل و دشوار بود که با اعمال شرایطی جدید فرم ساده و بسته‌ای برای تابع درست‌نمایی ارائه و بنابراین مفروضات برآوردگر باورمندی برای توزیع‌های آمیخته محاسبه شد. سپس باتوجه به روش ارائه شده

- 1976;71(356):791-9. doi: [10.1080/01621459.1976.10480949](https://doi.org/10.1080/01621459.1976.10480949).
14. Redner RA, Walker HF. Mixture densities, maximum likelihood and the EM algorithm. *SIAM Review*. 1984;26(2):195-239. doi: [10.1137/1026034](https://doi.org/10.1137/1026034).
 15. McLachlan GJ, Basford KE. *Mixture models: Inference and applications to clustering*. New York: M. Dekker; 1988.
 16. Berger JO. Prior information and subjective probability. In *Statistical Decision Theory and Bayesian Analysis 1985* (pp. 74-117). Springer, New York, NY. doi: [10.1007/978-1-4757-4286-2_3](https://doi.org/10.1007/978-1-4757-4286-2_3).
 17. Besag J, Green P, Higdon D, Mengersen K. Bayesian computation and stochastic systems. *Statistical Science*. 1995;3-41. doi: [10.1214/ss/1177010123](https://doi.org/10.1214/ss/1177010123).
 18. Najafabadi AT, Barmalzan G, Aghaei S. A weighted model confidence set: applications to local and mixture model confidence sets. arXiv preprint arXiv:1701.05455. 2017 Jan 19. doi: [10.1504/IJMMNO.2017.086797](https://doi.org/10.1504/IJMMNO.2017.086797).
 19. Hunter DR, Lange K. A tutorial on MM algorithms. *Am Stat*. 2004;58(1):30-7. doi: [10.1198/0003130042836](https://doi.org/10.1198/0003130042836).
 20. Lange K. The MM algorithm. In *Optimization 2013* (pp. 185-219). Springer, New York, NY. doi: [10.1007/978-1-4614-5838-8_8](https://doi.org/10.1007/978-1-4614-5838-8_8).
 21. Nguyen HD, McLachlan GJ. Maximum likelihood estimation of triangular and polygonal distributions. *Comput Stat Data Anal*. 2016;102:23-36. doi: [10.1016/j.csda.2016.04.003](https://doi.org/10.1016/j.csda.2016.04.003).
 22. Dempster AP, Laird NM, Rubin DB. Maximum likelihood from incomplete data via the EM algorithm. *J R Stat Soc*. 1977;39(1):1-38. doi: [10.1111/j.2517-6161.1977.tb01600.x](https://doi.org/10.1111/j.2517-6161.1977.tb01600.x).
 23. Najafabadi AT, Sakizadeh M. Designing an Optimal Bonus--Malus System Using the Number of Reported Claims, Steady-State Distribution, and Mixture Claim Size Distribution. arXiv preprint arXiv:1701.05441. 2017 Jan 19.
- and challenges. Tehran: National Population Studies & Comprehensive Management Institute. 2009. [Persian]
 3. Rezaeian M, Asadpour M, Hadavi M. Research barriers from the perspective of academic members and strategies for confronting with these barriers in Rafsanjan University of Medical Sciences, Iran. *Health System Research*. 2013;9(3):269-76. [Persian]
 4. Hakimzadeh SM, Hosseini Shokouh SM, Bahadori M, Tahernezhad K. Research needs assessment and priority setting for health economics: a mixed method study in Iran. *J Mil Med*. 2014;16(1):23-8. [Persian]
 5. Akbarzadeh BAR, M, Esmaceli M, Kimiafar Kh. Medical information management and assessment of direct costs of treatment of lung cancer. *Health Information Management*. 2009;5(210):151-8. [Persian]
 6. Pearson K. Contributions to the mathematical theory of evolution. *Philos Trans R Soc Lond A*. 1894;185:71-110. doi: [10.1098/rsta.1894.0003](https://doi.org/10.1098/rsta.1894.0003).
 7. Feddern HA. Hybridization between the western Atlantic angelfishes, *Holacanthus isabelita* and *H. ciliaris*. *Bulletin of Marine Science*. 1968;18(2):351-82.
 8. Whitney AW. *Theory of experience rating*. 1918.
 9. Keffer R. An experience rating formula. *Transactions of the Actuarial Society of America*. 1929;30:130-9.
 10. Bailey AL. *Credibility Procedures: Laplace's generalization of Bayes' Rule and the combination of collateral knowledge with observed data*. New York State Insurance Department; 1950.
 11. Bühlmann H, Straub E. Glaubwürdigkeit für Schadensätze. *Bulletin of the Swiss Association of Actuaries*. 1970;70(1):111-33.
 12. Norberg R. Credibility theory. *Encyclopedia of Actuarial Science*. 2004;1:398-406. doi: [10.1002/9780470012505.tac068](https://doi.org/10.1002/9780470012505.tac068).
 13. Box GE. Science and statistics. *J Am Stat Assoc*.