

Application of Data Mining through Machine Learning Algorithms to Study Effect of Car Features in Predicting Financial Claim of Motor Third Party Liability Insurance

Mohammadreza Asghari Oskoei¹ Farbod Khanizadeh² Azadeh Bahador³

Received: 2020 April 22

Accepted: 2020 December 2

Abstract

Objective: Risk classification of insurance customers, based on the observable characteristics, can significantly help insurers mitigate losses, classify their customers and prevent adverse selection. This paper aims to study losses occurred in motor Third Party Liability (TPL) insurance and predict customers' risk of loss.

Methodology: With the help of four supervised algorithms namely; decision tree, SVM, naïve Bayes and neural network hidden pattern of data is discovered to classify customers of TPL insurance. Furthermore, the imbalanced dataset was the main challenge for implementing machine learning and data mining techniques which will be discussed throughout the article.

Findings: The dataset contains more than 400,000 observations for five years from an Iranian insurance company. It also has five variables of which four are independent: car type, car group, plate type, car age; and one binary dependent variable: financial loss. Comparing the model performances, decision tree is the most efficient ($F1=0.72\pm 1$).

Conclusions: The model provides prioritization of independent features as follows: car type, plate type, car age, car group. Findings also suggest that to obtain more accurate prediction on claims and high-risk customers, more features concerning drivers' traits are required.

Keywords: Insurance Customer Classification, Decision Tree, Support Vector Machine, Naïve Bayes, Neural Networks.

JEL-Classification: G22, G17, F47.

1. Assistant Professor, Faculty of Mathematics and Computer Sciences, Allameh Tabatabaee'i University. (**Corresponding Author**) oskoei@atu.ac.ir

2. Assistant Professor, Insurance Research Centre. khanizadeh@irc.ac.ir

3. Holder of Automobile Research Desk, Insurance Research Centre. bahador@irc.ac.ir

کاربرد داده‌کاوی با استفاده از الگوریتم‌های یادگیری ماشین برای بررسی تاثیر ویژگی‌های خودرو در پیش‌بینی ریسک خسارت مالی در رشته بیمه شخص ثالث

محمد رضا اصغری اسکویی^۱، فرید خانی‌زاده^۲، آزاده بهادر^۳

تاریخ دریافت: ۱۳۹۹/۰۲/۰۳ تاریخ پذیرش: ۱۳۹۹/۰۷/۱۲

چکیده

هدف: طبقه‌بندی ریسک بیمه‌گذاران بر مبنای ویژگی‌های قابل مشاهده می‌تواند به شرکت‌های بیمه جهت کاهش زیان، شناخت دقیق‌تر مشتریان و جلوگیری از وقوع انتخاب نامساعد در بازار بیمه کمک شایانی کند. هدف این مقاله، بررسی خسارت‌های مالی ایجاد شده در بیمه شخص ثالث و پیش‌بینی ریسک بیمه‌گذاران در احتمال وقوع حادثه می‌باشد.

روش‌شناسی: با استفاده از الگوریتم‌های درخت تصمیم، ماشین بردار پشتیبان، نایو بیز و شبکه عصبی؛ به کشف الگوهای پنهان داده‌ها، در راستای طبقه‌بندی بیمه‌گذاران بیمه شخص ثالث پرداخته شده است. همچنین توزیع نامتعادل داده‌ها در دو گروه خسارت‌دیده و خسارت‌ندیده سبب یک چالش مهم در کاربرد روش‌های یادگیری ماشین و داده‌کاوی است که در این مقاله مورد توجه قرار گرفته است.

یافته‌ها: مجموعه داده متعلق به یکی از شرکت‌های بیمه و حاوی بیش از چهارصد هزار نمونه ثبت شده در پنج سال و شامل چهار متغیر مستقل نوع خودرو، گروه خودرو، نوع پلاک و سن خودرو و یک متغیر وابسته و دو ارزشی خسارت مالی است. با توجه به نتایج بدست آمده بهترین کارکرد و دقت پیش‌بینی (با دقت $F1 = 0.72 \pm 0.01$) مربوط به مدل درخت تصمیم می‌باشد.

نتیجه‌گیری: میزان تأثیرگذاری متغیرها در وقوع خسارت به ترتیب اولویت عبارتند از: نوع خودرو، نوع پلاک، سن خودرو و گروه خودرو. نتایج ارزیابی نشان می‌دهد برای پیش‌بینی دقیق‌تر خسارت و مشتریان پر ریسک به داده‌های بیشتری مرتبط با ویژگی‌های راننده نیاز می‌باشد.

کلید واژه‌ها: دسته‌بندی مشتریان، درخت تصمیم، ماشین بردار پشتیبان، نایو بیز و شبکه‌های عصبی.

طبقه‌بندی موضوعی: G22, G17, F47.

۱. استادیار دانشکده علوم ریاضی و رایانه دانشگاه علامه طباطبائی (نویسنده مسئول) oskoei@atu.ac.ir

۲. استادیار پژوهشکده بیمه و مسئول میز تخصصی طراحی الگوریتم و یادگیری ماشین، khanizadeh@irc.ac.ir

۳. پژوهشگر پژوهشکده بیمه و مسئول میز تخصصی بیمه‌های اتومبیل، bahador@irc.ac.ir

مقدمه

بیمه یکی از اصلی‌ترین ابزارهای مدیریت ریسک است و شرکت‌های بیمه با قبول ریسک باعث ایجاد آرامش در جامعه می‌گردند. به همین خاطر لازم است که شرکت‌های بیمه به ابزارهای تحلیل ریسک قدرتمندی دسترسی داشته باشند تا بتوانند ریسک دریافتی را به خوبی مدیریت کنند. هر بیمه‌گذار یا هر مورد بیمه‌شده، سطح متفاوتی از ریسک را به شرکت بیمه تحمیل می‌نماید. برای اطمینان از اینکه هر بیمه‌گذار حق بیمه منصفانه‌ای را پرداخت می‌کند، بیمه‌گران سطح ریسک بیمه‌گذار را تعیین و او را در یکی از طبقات ریسک قرار می‌دهند که بالتبع هر چه ریسک بیشتر باشد، حق بیمه بیشتر خواهد بود.

در شرایط فعلی، ارزیابی ریسک در صنعت بیمه کشور، براساس تجربیات سایر کشورها صورت می‌گیرد و صنعت بیمه از فقدان الگوریتم‌ها و سیستم‌های خودکاری که با حساسیت قابل قبولی بتوانند میزان ریسک مشتریان مختلف را بررسی و ارزیابی کند رنج می‌برد. این موضوع در رشته بیمه شخص ثالث که یک بیمه اجباری است، شرایط را حادثر می‌کند و در همین راستا ارزیابی ریسک فرد در رشته بیمه شخص ثالث بسیار حائز اهمیت است. به همین دلیل استفاده از ابزارهای داده‌کاوی می‌تواند در سنجش و پیش‌بینی ریسک بیمه‌گذاران بسیار راه‌گشا باشد.

در صنعت بیمه، داده‌کاوی می‌تواند به شرکت‌ها جهت کسب مزیت تجاری کمک کند. به عنوان مثال با به‌کارگیری تکنیک‌های داده‌کاوی، شرکت‌ها می‌توانند با استفاده از داده‌ها در مورد الگوهای خرید و رفتار مشتری، به کسب دانش پرداخته و همچنین درک خود را برای کمک به کاهش تقلب، ارتقای بیمه‌گری و بالابردن مدیریت ریسک افزایش دهند. (حاجی‌حیدری و همکاران، ۱۳۹۰)

شناسایی مشتریان مستلزم تحلیل مشتریان هدف و دسته‌بندی کردن مشتریان است که منجر به یافتن گروه‌هایی از مشتریان سودآور براساس ویژگی‌های آن‌ها می‌شود. طبقه‌بندی ریسک در حقیقت به معنای گروه‌بندی مشتریان با خصوصیات ریسک مشابه است که احتمال بروز خسارت‌های مشابهی دارند. طبقه‌بندی ریسک بیمه‌گذاران بر

مبنای ویژگی‌های قابل مشاهده می‌تواند به شرکت‌های بیمه جهت کاهش زیان، افزایش نرخ پوشش بیمه و جلوگیری از وقوع انتخاب نامساعد در بازار بیمه کمک شایانی کند. (حنفی‌زاده و رستخیز پایدار، ۱۳۹۰)

عدم دسته‌بندی مشتریان باعث شده که مشتریان کم‌ریسک‌تر، خسارات مالی مشتریان پُر ریسک را جبران کنند. از این رو تفاوت چندانی بین مشتریان پُر ریسک و کم‌ریسک وجود ندارد. در واقع در کشور ما به جای فرد، اتومبیل بیمه می‌شود و این امر موجب شده تا بیشتر شرکت‌های بیمه در زمینه بیمه اتومبیل، متحمل زیان شوند (ترکستانی و همکاران، ۱۳۹۵). با استفاده از ابزارهای دسته‌بندی به منظور جداسازی مشتریان و شناخت آن‌ها، می‌توان سیاست‌گذاری‌های مناسبی را برای آن‌ها در نظر گرفت (ایزدپرست، ۱۳۹۰). در این مقاله از روش‌های درخت تصمیم^۱، ماشین بردار پشتیبان^۲، نایو بیز (بیز ساده)^۳ و همچنین شبکه عصبی^۴، برای تحلیل داده‌های بیمه شخص ثالث و طبقه‌بندی نمونه‌های خسارت دیده و خسارت نادیده استفاده نمودیم و دقت مدل‌ها در این چهار روش را مقایسه کردیم.

۱. پیشینه تحقیق

در کشور ما بیمه اتومبیل و به‌ویژه بیمه شخص ثالث از مهمترین رشته‌های بیمه‌ای است که سهم عمده‌ای را در پرتفوی صنعت بیمه به خود اختصاص داده است و از طرفی به دلیل داشتن ضریب خسارت بالا، توجه بیش از پیش به این رشته بیمه‌ای را ضروری می‌نماید. این نوع بیمه در اکثر کشورهای جهان، یکی از مهمترین نوع فعالیت بیمه‌ای محسوب می‌شود و دست‌کم در ایران، حدود نیمی از صنعت بیمه را در اختیار دارد. ضمن اینکه به دلیل وجود این سهم عمده، فرصت مناسبی برای کاوش اطلاعات و

1. Decision Tree
2. Support Vector Machine
3. Naïve Bayes
4. Neural Networks

استخراج الگوهای ناشناخته جهت تصمیمات کلان در صنعت بیمه از این طریق فراهم می‌شود. (کریم‌زادگان مقدم و بهروان، ۱۳۹۴)

عدم توجه به سطح ریسک مشتریان از سوی شرکت‌های بیمه، موجب شده تا مشتریان با سطوح ریسکی متفاوت، حق بیمه یکسانی را پرداخت کنند که از طرفی موجب نارضایتی بیمه‌گذاران و از طرفی دیگر باعث افزایش روزافزون ضریب خسارت و زیان‌ده شدن در این رشته شده است. از نوآوری‌های این تحقیق به نسبت تحقیقات گذشته، می‌توان به تعداد حجم داده‌های استفاده‌شده (۴۱۵,۶۸۷) در این تحقیق اشاره نمود و همچنین بهره‌مندی از دو روش داده کاوی^۱ و الگوریتم‌های یادگیری ماشین^۲ که موجب ارائه ارزیابی دقیقی از میزان تأثیر مشخصات خودرو در ایجاد خسارت شده است. تاکنون در هیچ تحقیقی در داخل کشور خودروهای سواری براساس کیفیت‌شان مورد تحلیل قرار نگرفته بودند و اکثر تحقیقات داخلی در خصوص بیمه بدنه اتومبیل انجام شده و تحقیقات اندکی به بیمه شخص ثالث اختصاص یافته است. در ادامه به برخی پژوهش‌های مرتبط انجام شده در داخل و نیز در جدول (۱) به پژوهش‌های خارجی اشاره می‌شود.

در پژوهش ترکستانی و همکاران (۱۳۹۵)، از شبکه عصبی در راستای پیش‌بینی میزان خسارت بالقوه بیمه‌گذاران و تعیین نرخ بهینه استفاده شده است و نتایج پژوهش نشان می‌دهد که مدل ارائه شده می‌تواند با دقت ۹۱ درصد طبقه خسارتی را تخمین بزند و با دقت ۸۷ درصد میزان خسارت بالقوه بیمه‌گذاران را پیش‌بینی کند. پژوهش کریم‌زادگان مقدم و بهروان (۱۳۹۴)، در خصوص تعرفه‌گذاری پویا در رشته بیمه شخص ثالث می‌باشد که در این مقاله از شبکه‌های عصبی، درخت تصمیم و خوشه‌بندی استفاده شده است که نتایج به‌دست آمده از مدل‌ها با استفاده از ماتریس آشفتگی^۳ و نسبت خسارت، مورد اعتبارسنجی قرار گرفته که نتایج حاکی از امکان

1. Data Mining
2. Machine Learning
3. Confusion Matrix

استفاده از روش ارائه شده در تعرفه گذاری پویا در خصوص بیمه شخص ثالث به صورتی کارآمد را نشان می دهد، به نحوی که نسبت خسارت، کاهش می یابد و ماتریس آشفستگی، صحت ارزیابی را نشان می دهد.

در پژوهش حاجی حیدری و همکاران (۱۳۹۰)، با در نظر گرفتن همزمان مشخصه های بیمه گذار و اتومبیل در رشته بیمه بدنه، چند مدل (درخت تصمیم، شبکه های عصبی، شبکه های بیزین، ماشین بردار پشتیبان، رگرسیون لجستیک و تحلیل تمایزی) را به منظور پیش بینی طبقه خسارتی بیمه گذاران مقایسه کردند. طبق نتایج این پژوهش، مدل های شبکه های عصبی و درخت تصمیم با حدود ۸۲ درصد، بیشترین دقت را در پیش بینی داشتند. حنفی زاده و رستخیز پایدار (۱۳۹۰)، ابتدا عوامل موثر بر ایجاد خسارت در بدنه اتومبیل را در ایران بررسی کردند. پس از مشخص شدن عوامل با استفاده از شبکه های عصبی خودسازمان ده^۱، به خوشه بندی بیمه گذاران براساس ریسک بالقوه آن ها پرداختند. در پژوهشی دیگر، فتح نژاد و ایزدپرست (۱۳۹۰)، با استفاده از تکنیک خوشه بندی k-means و درخت تصمیم بیمه گذاران را خوشه بندی کردند و نتیجه گرفتند که علاوه بر مشخصات اتومبیل، مشخصات رفتاری مشتری نیز در پیش بینی سطح خسارت مشتریان بیمه بدنه اتومبیل تأثیرگذار است. دقت مدل های استفاده شده در این پژوهش حدود ۶۰ درصد بوده است. در پژوهش اصغری اسکویی و قاسم زاده (۱۳۹۵) و اصغری اسکویی (۱۳۹۴) رویکرد انتخاب ویژگی^۲ براساس الگوریتم تکاملی^۳ و کاربرد شبکه عصبی برای پیش بینی سری زمانی به کار رفته است.

پژوهش های خارجی در این موضوع طیف وسیعی را شامل می شود. از جمله ویو و سرنا (Wuyu and Cerna, 2019) برای ارزیابی ریسک مدل سازی و پیش گویی کننده و تحلیل الگوی سطح ریسک در حوزه بیمه اتومبیل از درخت تصمیم و شبکه عصبی استفاده نمودند. بآنک و بکاء (Baecke and Bocca, 2017) برای انتخاب بیمه گذاران

1. Self-Organization Map (SOM)
2. Feature Selection
3. Evolutionary Algorithm

مناسب با توجه به سطح ریسک آن‌ها از درخت تصمیم، رگسیون لجستیک و شبکه عصبی استفاده نمودند. همچنين فرپنگ و همکارانش (Frempong, Nicholas and Boateng 2017) از درخت تصمیم به عنوان یک پیش‌بینی کننده برای خسارت اتومبیل استفاده نموده و نتایج قابل توجهی گزارش نموده‌اند.

در یک پژوهش مشابه، ثاکو و ساین (Thakur and Sing, 2013) از درخت تصمیم برای طبقه‌بندی بیمه‌گذاران و تعمیم آن به بیمه‌گذاران جدید بیمه اتومبیل و نیز کاسلان و همکاران (Kaščelan, et al., 2016) از ترکیب روش‌های خوشه‌بندی، رگرسیون، ماشین بردار پشتیبان برای طبقه‌بندی ریسک و میزان خسارت بیمه اتومبیل استفاده نمودند. ضمناً برای پیش‌بینی خسارت بیمه اتومبیل، یونس و همکارانش (Yunos, Ali, Shamsyuddin and Ismail 2016) یک شبکه عصبی ارائه نموده و نتایج قابل قبولی ثبت کرده‌اند.

۲. مروری بر ادبیات تحقیق

در این بخش به بررسی مبانی نظری مرتبط با این مقاله می‌پردازیم:

۱-۲. داده کاوی

در صنعت بیمه، اطلاعات و استفاده از آن بسیار حائز اهمیت است، به طوری که می‌توان موفقیت در بیمه را در گرو توانایی شرکت‌ها در تبدیل داده‌های خام به اطلاعات کاربردی دانست.

داده‌کاوی موجب بهبود روند تصمیم‌گیری در یک سازمان از طریق استخراج اطلاعات مهم از داده‌های موجود و جستجو روابط و الگوهای پنهان و آشکار از مجموعه داده‌های جمع‌آوری شده توسط سازمان خواهد شد و نهایتاً به بهینه‌سازی تصمیمات کسب‌وکار، ارتباطات و بهبود رضایتمندی مشتریان کمک می‌کند (حنفی‌زاده و رستخیز پایدار، ۱۳۹۰).

داده‌کاوی بر طبق تعریف موسسه سیستم تحلیل آماری، فرایند انتخاب، اکتشاف، مدل‌سازی و شفاف‌سازی الگوهای مفید و ناشناخته در حجم زیادی از داده می‌باشد.

۲-۲. روش‌ها و تکنیک‌های داده‌کاوی

بر حسب اینکه در فرایند داده‌کاوی، استنتاج چه نوع دانشی از مجموعه آموزشی مورد نظر است، از روش‌های مختلف داده‌کاوی می‌توان بهره جست. به‌طور کلی الگوریتم‌های یادگیری ماشین از نظر شیوه یادگیری به دو دسته اصلی الگوریتم‌های یادگیری با نظارت و الگوریتم‌های یادگیری بدون نظارت تقسیم می‌شوند که به صورت مختصر اینجا معرفی می‌شوند.

- یادگیری با نظارت

در این نوع از الگوریتم‌ها، با دو نوع از متغیرها سروکار داریم. نوع اول که متغیرهای مستقل نامیده می‌شوند، یک یا چند متغیر هستند که براساس مقادیر آنها، متغیر دیگری را پیش‌بینی خواهیم نمود. نوع دوم هم متغیرهای وابسته یا هدف یا خروجی هستند که مقادیر آنها را به کمک این الگوریتم‌های یادگیری با نظارت پیش‌بینی خواهیم نمود. برای این منظور باید تابعی ایجاد کنیم که ورودی‌ها (متغیرهای مستقل) را گرفته و خروجی مورد نظر (متغیر وابسته یا هدف) را تولید کند.

نمونه‌هایی از این الگوریتم‌ها عبارتند از رگرسیون، درخت‌های تصمیم، جنگل‌های تصادفی، N نزدیک‌ترین همسایه، نایو بیس، ماشین بردار پشتیبان، شبکه عصبی و ... مسائل یادگیری با نظارت، به دو گروه طبقه‌بندی^۱ (برای پیش‌بینی پاسخ‌های گسسته) و رگرسیون^۲ (برای پیش‌بینی پاسخ‌های پیوسته) تقسیم می‌شوند.

1. Classification
2. Regression

- یادگیری بدون نظارت

در این نوع از الگوریتم‌ها، متغیر هدف نداریم و خروجی الگوریتم براساس الگوی درون داده‌ها مشخص می‌شود. بهترین مثال برای این نوع از الگوریتم‌ها، خوشه‌بندی یک جمعیت با داشتن اطلاعات شخصی و خریدهای مشتریان می‌باشد که به صورت خودکار آن‌ها را به گروه‌های همسان و هم‌ارز تقسیم کنیم.

در این دسته از یادگیری، تنها ورودی (x) را داریم و خروجی از پیش تعیین شده نیست. در واقع اینجا ناظری وجود ندارد تا به الگوریتم در یادگیری کمک کند. هدف اصلی یادگیری بدون نظارت، مدل کردن توزیع داده می‌باشد تا بتوان اطلاعات بیشتری درباره داده را بدست آورد. برعکس یادگیری با نظارت، هیچ ناظری وجود ندارد و مدل مجبور است خودش ساختار مخفی داده بدون برچسب را پیدا کند. الگوریتم K-Means از این دسته هستند.

- یادگیری عمیق

یکی از چالش‌های مهم در یادگیری ماشین، انتخاب بهینه ویژگی‌های مؤثر در فرایند یادگیری است. معمولاً ویژگی‌ها به صورت انتخاب مستقیم متغیرهای ورودی و یا به صورت ترکیب خطی یا غیر خطی آن‌ها حاصل می‌شوند. یادگیری عمیق فرایند انتخاب یا استخراج ویژگی را همزمان در طول فرایند یادگیری با هدف رسیدن به حداکثر کارآمدی و حداقل خطا آموزشی انجام می‌دهد. این نوع یادگیری از مباحث نوین و پرکاربرد در علوم کامپیوتر می‌باشد که قابلیت یادگیری الگوهای پیچیده را نیز دارد و به خاطر قدرت و دقت بالایشان در بسیاری از مسائل دنیای واقعی به کار گرفته شده‌است.

۲-۲-۱. درخت تصمیم

درخت تصمیم، نوعی روش یادگیری با نظارت است که با کمک یک ساختار درختی نتایج دسته‌بندی را ارائه می‌دهد. در این درخت هر گره نشانگر یک آزمون برای یک تصمیم بر روی یک متغیر مستقل است و هر شاخه، خروجی آزمون را نمایش می‌دهد.

برگ‌های درخت نیز نمایانگر تصمیم نهائی و کلاس‌ها است. به‌طور عادی، پیچیدگی یک درخت تصمیم با افزایش تعداد ویژگی‌ها افزایش می‌یابد. اگر چه در بعضی از شرایط، تنها تعداد کمی از ویژگی‌ها می‌توانند یک کلاس را تعیین کند و بقیه ویژگی‌ها کم‌تأثیر یا بی‌تأثیر می‌باشد (ایزدپرست، ۱۳۹۰).

۲-۲-۲. ماشین بردار پشتیبان

ماشین بردار پشتیبان، روش به نسبت جدیدی در حوزه داده‌کاوی می‌باشد که در بسیاری از مسائل طبقه‌بندی به‌طور موفقیت‌آمیزی عمل کرده است. ماشین بردار پشتیبان، یک طبقه‌بندی‌کننده دوتایی است که با استفاده از نگاشت داده‌ها از فضای ورودی اصلی به فضایی با بعد بالاتر برای جداسازی آن‌ها عمل می‌کند. این مدل ابر صفحه‌ای را جستجو می‌کند که فاصله‌اش با داده‌های دو طبقه حداکثری است. با امکان تعریف مرزهای انعطاف‌پذیر قدرت تعمیم‌پذیری نسبت به داده‌های جدید را افزایش می‌دهد. ماشین بردار پشتیبان می‌تواند با استفاده از داده‌های آموزشی کمتر نسبت به روش‌های رقیب، مرزهای سیستم را با دقت مناسبی تخمین بزند، بدون آنکه تعمیم‌پذیری سیستم را مخدوش کند. (حاجی‌حیدری و همکاران، ۱۳۹۰)

۲-۲-۳. نایو بیز

اغلب به عنوان یک راه‌کار ساده آماری برای دسته‌بندی و تشخیص برچسب اشیا یا نقاط از روش نایو بیز استفاده می‌شود. الگوریتم نایو بیز، مبتنی بر مشاهدات آماری و احتمال با فرض استقلال ویژگی‌ها نسبت به یکدیگر عمل می‌کنند. در بیشتر مدل‌های نایو بیز از روش حداکثرسازی تابع درست‌نمایی استفاده می‌شود. هر چند تکنیک نایو بیز دارای فرضیات محدود و قابل دسترس است؛ ولی به خوبی می‌تواند از عهده حل مسائل واقعی برآید. یکی از مزایای قابل توجه در الگوریتم نایو بیز، امکان برآورد پارامترهای مدل با اندازه نمونه کوچک به عنوان مجموعه «داده آموزشی» می‌باشد (Salma, et al. 2019).

۲-۲-۴. شبکه‌های عصبی

شبکه‌های عصبی، ساختارهای شبکه‌ای بسیار سازمان‌یافته‌ای هستند و دارای سه نوع لایه می‌باشند؛ لایه‌های ورودی، لایه‌های خروجی و لایه‌های میانی (یا لایه‌های پنهان). هر کدام از گره‌ها (که به نام نرون شناخته می‌شود)، در لایه‌های پنهان و لایه‌های خروجی دارای یک کلاسه‌بند^۱ هستند. نرون‌های ورودی، ابتدا اطلاعات ویژگی‌های شیء را دریافت و سپس به نرون‌های لایه پنهان ارسال می‌کنند. لایه پنهان این اطلاعات را پردازش کرده و نتایج را به لایه پنهان بعدی می‌فرستد. این پروسه ادامه می‌یابد تا اطلاعات به نرون‌های لایه خروجی برسد. در آنجا مقدار^۲ به‌دست‌آمده، تعیین‌کننده احتمال دسته‌بندی قرارگیری شیء می‌باشد. مجموعه این پروسه به عنوان پیش‌انتشار^۳ شناخته می‌شود. نمره به‌دست‌آمده در خروجی، نشانگر دسته‌ای است که مجموعه ورودی‌ها به آن تعلق دارند. به این نوع شبکه عصبی، پرسپترون چندلایه^۴ یا MLP گفته می‌شود. خروجی یک نرون، از جمع وزن‌دار ورودی‌ها و اعمال تابع فعالیت نرون بر آن حاصل می‌شود. ضرایب وزنی که به عنوان وزن^۵ و بایاس^۶ شناخته می‌شوند، در طول فرایند یادگیری شبکه، بر اساس مجموعه نمونه‌های آموزشی و الگوریتم پس‌انتشار خطا تنظیم می‌شوند. قبل از یادگیری، انتخاب پیکره‌بندی مناسب شبکه از جمله تعداد لایه‌ها، تعداد نرون در هر لایه و تابع فعالیت نرون‌ها که بستگی به نوع مسئله و پیچیدگی آن دارد، باید انجام شود.

از کاربردهای شبکه‌های عصبی می‌توان طبقه‌بندی، شناسایی و تشخیص الگو، پیش‌بینی سری‌های زمانی، بهینه‌سازی، سیستم‌های خبره و فازی، مسائل مالی، بیمه،

1. Classifier
2. Value
3. Forward Propagation
4. Multi-Layer Perceptron (MLP)
5. Weight
6. Bias

امنیتی، بازار بورس و وسایل سرگرم‌کننده و ساخت وسایل صنعتی، پزشکی و امور حمل و نقل را نام برد. (عمرانی نوش‌آبادی، ۱۳۹۰)

۳. روش تحقیق

این پژوهش یک پژوهش کاربردی می‌باشد و تلاش شده است تا با استفاده از الگوریتم‌های درخت تصمیم، ماشین بردار پشتیبان، نایو بیز و شبکه عصبی؛ به کشف الگوهای پنهان داده‌ها، در راستای طبقه‌بندی بیمه‌گذاران بیمه شخص ثالث پرداخته شود. مراحل اجرایی در این تحقیق به صورت زیر می‌باشد:

- جمع‌آوری داده‌ها از پایگاه داده بیمه‌گذاران بیمه شخص ثالث یکی از شرکت‌های بیمه؛

- پیش‌پردازش^۱ و پالایش^۲ داده‌ها و تعیین شاخص‌هایی برای تعریف طبقات ریسک بیمه‌گذاران؛

- بررسی آماری و تقسیم داده‌ها به زیر مجموعه‌های متعادل و تصادفی در دو دسته داده‌های آزمایشی و داده‌های آموزشی؛

- استخراج الگوها با استفاده از الگوریتم‌های درخت تصمیم، ماشین بردار پشتیبان، نایو بیز و شبکه عصبی و مقایسه نتایج حاصله از این الگوریتم‌ها؛

- ارائه الگوی کشف‌شده از طبقه‌بندی بیمه‌گذاران و شناسایی ویژگی‌های تعیین‌کننده؛

- ارزیابی نتایج طبقه‌بندی و اعتبارسنجی مدل.

برای ساخت مدل لازم است که ابتدا تکنیک مدل‌سازی انتخاب شود که در این مقاله چهار روش (درخت تصمیم، ماشین بردار پشتیبان، نایو بیز و شبکه عصبی) بررسی شده است و ابزار مورد استفاده در این مقاله کتابخانه سای کیت لرن^۳ در زبان برنامه‌نویسی پایتون می‌باشد.

1. Pre-processing
2. Cleaning
3. Sci-Kit Learn Library in Python

۴. داده‌های پژوهش

جامعه آماری تحقیق شامل کلیه بیمه‌گذاران یکی از شرکت‌های بیمه در رشته بیمه شخص ثالث در بازه زمانی ابتدای سال ۱۳۹۲ تا انتهای سال ۱۳۹۶ می‌باشد. از این میان، تعدادی از بیمه‌گذاران خسارت دریافت کرده و تعدادی دیگر، خسارتی از این شرکت بیمه دریافت ننموده‌اند. در مجموع بیش از چهارصد هزار رکورد در بانک اطلاعاتی بیمه‌گذاران بیمه شخص ثالث ثبت شده بود که پس از مورد ارزیابی و تحلیل قرار گرفتن، متغیرهای موجود در این مجموعه در جدول (۱) ارائه شده است.

شایان ذکر است که در متغیرهای اشاره شده در جدول فوق، دسته‌بندی‌های مربوط به گروه خودرو و نوع پلاک توسط شرکت بیمه صورت گرفته است. در مورد نوع خودرو، طبقه‌بندی ارائه شده با ایده تیم نویسندگان انجام شده که براساس طبقه‌بندی‌های موجود در بازار و یا میزان فراوانی نوع خودروی مدنظر در مجموعه داده‌ها می‌باشد. به عنوان مثال، به طور عرف، دسته‌بندی موتورسیکلت‌ها بر اساس میزان قدرت موتورسیکلت بر حسب سی‌سی صورت می‌گیرد. همچنین برای گروه اتوکار، اتوبوس و ون بیشترین فراوانی را در بین سایر موارد به خود اختصاص داده بودند.

جدول ۱. متغیرهای به‌کار رفته در مدل

مقادیر متغیر	نام متغیر / نوع متغیر	متغیر
خودروهای ژاپنی، کره‌ای، آمریکایی و اروپایی / High	سواری / Car	مشخصات اتومبیل نوع خودرو
خودروهای ایرانی و چینی / LM		
کمتر از ۱۲۵ سی‌سی / M125	موتورسیکلت‌ها / Motor	
بین ۱۲۵ تا ۲۰۰ سی‌سی / M200		
بین ۲۰۰ تا ۳۰۰ سی‌سی / M300		
سایر / MotOther	CarType / متغیر مستقل	
وانت / Truck		
کامیون / Lorry		
تریلر / Trailer		
سایر / ConOther	اتوکار / Autocar	
اتوبوس / Bus		
ون / Van		

متغیر	نام متغیر / نوع متغیر	مقادیر متغیر		
گروه خودرو	/Car Group متغیر مستقل	سایر / AutOther		
		سواری / Car		
		موتورسیکلت‌ها / Motor		
		بارکش / Container		
		اتوکار / Autocar		
		ماشین‌آلات کشاورزی / Farming		
نوع پلاک	/Plate Type متغیر مستقل	شخصی / Private		
		فاقد پلاک / Lack		
		عمومی / Public		
		دولتی / Gov		
		ترانزیت / Transit		
		گذر موقت / Temp		
		معلولان / Disabled		
		منطقه آزاد / Free		
		نظامی / Military		
		سیاسی / Politics		
		از ۰ تا ۴۵ سال		
		سن وسیله نقلیه	/Car Age متغیر مستقل	بیمه‌گذار خسارت مالی ندیده است: ۰
				بیمه‌گذار خسارت مالی دیده است: ۱
خسارت مالی	Property /Damage متغیر وابسته			

مأخذ: مجموعه داده‌های پژوهش

۴-۱. پیش‌پردازش داده‌ها

پیش‌پردازش داده‌ها از گام‌های مهم فرایند داده‌کاوی است که میزان دقت نتایج به دست آمده، تا حد زیادی به اجرای درست آن بستگی دارد. یک تعریف ساده می‌تواند این باشد که پیش‌پردازش داده‌ها مجموعه عملیات و روش‌هایی است برای تبدیل داده‌های خام جمع‌آوری شده از منابع متنوع به اطلاعات پاک‌سازی شده‌ای که برای انجام تحلیل‌ها مناسب باشد.

در واقع کیفیت داده‌ها و اطلاعات مفیدی که از آن حاصل می‌شود؛ به‌طور مستقیم بر توانایی مدل برای یادگیری تأثیر می‌گذارد. بنابراین بسیار مهم است که ما داده‌های

خود را قبل از ارائه به مدل، مورد پیش‌پردازش قرار دهیم. در همین راستا برای پیش‌پردازش داده‌ها، اقدامات به شرح ذیل انجام شد:

۴-۱-۱. حذف متغیرهای نامناسب

برخی از متغیرهای موجود در پایگاه داده مانند شماره بیمه‌نامه، تاریخ صدور بیمه‌نامه و... به دلیل بی‌ارتباط بودن با هدف پژوهش از مجموعه متغیرها حذف شدند.

۴-۱-۲. تبدیل کلیه مقادیر به مقادیر عددی

از بین متغیرهای مستقل نوع خودرو، گروه خودرو، نوع پلاک و سن وسیله نقلیه؛ تنها متغیر سن وسیله نقلیه، متغیر عددی می‌باشد و سایر متغیرها، متغیر اسمی می‌باشند و لذا برای این متغیرهای اسمی، کدگذاری انجام شد.

۴-۱-۳. بررسی صلاحیت داده‌ها جهت ورود به مدل نهایی

در بین حجم انبوهی از داده‌ها، تمام داده‌ها از کیفیت لازم برخوردار نبودند. لذا معیوب‌بودن داده‌ها از لحاظ خطاهای اندازه‌گیری بررسی شد؛ به عنوان مثال در متغیر سن خودرو، اعداد نامتعارفی (اعداد سه رقمی) وارد شده بود که رکوردها از لحاظ وجود داده‌های نامرتبب مورد بررسی قرار گرفت و رکوردهایی که قابلیت اصلاح را داشتند، اصلاح شده و در صورتی که این قابلیت را نداشتند، حذف شدند.

۴-۱-۴. حذف رکوردهای ناقص

پیش‌پردازش در خصوص رکوردهای ناقص نیز صورت گرفت، به این معنا که چنانچه اکثر متغیرهای یک رکورد، گم‌شده باشند، آن رکورد حذف شده است که این موضوع به کاهش داده‌های نهایی منجر شد.

۴-۱-۵. حذف داده‌های دارای نوفه^۱ (نویز)

در برخی از متغیرها، یک داده به دو دسته متفاوت تعلق داشت؛ به‌عنوان مثال در برخی موارد خودرو و وانت یک‌بار به عنوان سواری و یک‌بار به عنوان بارکش ثبت شده بود که تا جای امکان، این‌گونه اطلاعات دارای نویز از داده‌ها حذف گردید.

۵. آمار توصیفی

از آمار توصیفی به منظور سازمان‌دهی، خلاصه‌کردن و توصیف اطلاعات استفاده می‌شود و معمولاً قبل از آنالیز داده‌ها، سازمان‌دهی داده‌ها، می‌تواند منجر به آشکارشدن نکات پنهان داده‌ها شود. لذا به‌همین منظور در این بخش، به بررسی آمار توصیفی مربوط به متغیرهای مستقل (نوع پلاک، گروه خودرو، نوع خودرو و سن وسیله نقلیه) و متغیر وابسته (خسارت مالی داشتن یا نداشتن) به شرح جدول (۲) پرداخته شده است.

جدول ۲. نوع پلاک

نوع پلاک	تعداد	درصد فراوانی	درصد خسارت دیده مالی
خصوصی	۲۹۲۱۴۹	٪۷۰	٪۸۷/۶۰
فاقد پلاک	۱۶۶۷۶	٪۴	٪۷۰
عمومی	۶۷۰۷	٪۲	٪۴/۱۴
دولتی	۳۷۵۰	٪۱	٪۱/۶۶
ترانزیت	۹۵۵۷۵	٪۲۳	٪۲/۳۵
عبور موقت	۹۴	---	٪۰/۰۲
معلولان	۸۰	---	٪۰/۰۳
منطقه آزاد	۶۵۱	---	٪۰/۲
نظامی	۱	---	---
سیاسی	۴	---	---

مأخذ: یافته‌های پژوهش

1. Noise

همانطور که در جدول (۲) مشاهده می‌شود متغیر «نوع پلاک» به ۱۰ زیر گروه تقسیم می‌شود و بیشترین و کمترین فراوانی به ترتیب متعلق به پلاک‌های خصوصی و نظامی می‌باشد. خودروهای با پلاک ترانزیت با وجود درصد فراوانی بالا پس از پلاک‌های خصوصی، لیکن با توجه به مقادیر ستون چهارم خسارت بالایی به بار نیاورده‌اند. به طور مشابه جداول (۳) و (۴) فهرست ویژگی‌ها و دسته‌بندی آن‌ها و آمار توصیفی متغیرهای «گروه خودرو» و «نوع خودرو» را نمایش می‌دهد.

جدول ۳. گروه خودرو

گروه خودرو	تعداد	درصد فراوانی	درصد خسارت دیده مالی
موتور سیکلت	۱۳۴۷۶	۳	۰/۸۶
اتوکار	۱۳۰۷	---	۰/۲۷
بارکش	۱۳۱۹۹۸	۳۲	۱۸/۲۵
سواری	۲۶۷۴۱۵	۶۵	۸۰/۴۴
کشاورزی	۱۴۹۱	---	۰/۱۸

مأخذ: یافته‌های پژوهش

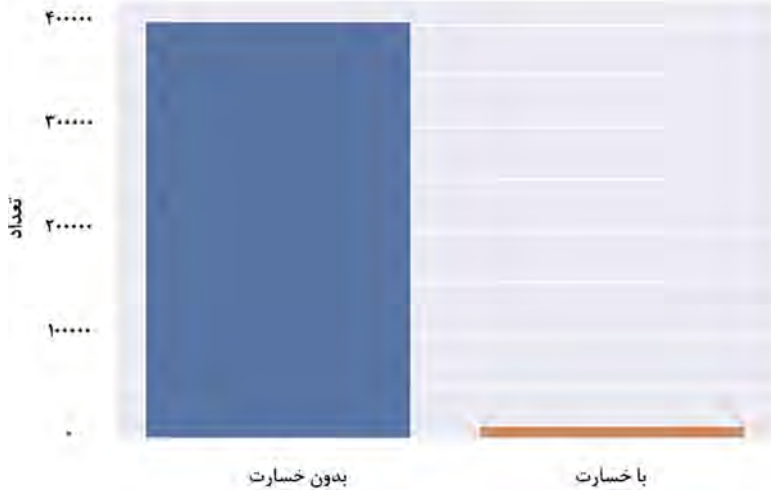
جدول ۴. نوع خودرو

نوع خودرو	تعداد	درصد فراوانی	درصد خسارت دیده مالی
موتور سیکلت‌ها	کمتر از ۱۲۵ سی‌سی	۸۳۲۰	۰/۵۳
	بین ۱۲۵ تا ۲۰۰ سی‌سی	۳۴۴۰	۰/۲۵
	بین ۲۰۰ تا ۳۰۰ سی‌سی	۲۳۶	۰/۰۲
	سایر	۱۴۷۹	۰/۰۶
اتوکار	اتوبوس	۹۴۸	۰/۱۰
	ون	۲۳۰	۰/۱۷
	سایر	۱۱۳	---
بارکش	وانت	۲۶۰۰۱	۹/۰۵
	کامیون	۱۹۹۱۰	۶/۸۲
	تریلر	۸۵۰۷۰	۲/۲۳

درصد خسارت دیده مالی	درصد فراوانی	تعداد	نوع خودرو	
۰/۱۳	۱	۹۸۷	سایر	
۱۱/۱۸	۱۶	۴۳۱۸۴	ژاپنی، کره‌ای، آمریکایی و اروپایی	سواری
۶۹/۲۸	٪۸۴	۲۲۴۲۷۸	خودروهای ایرانی و چینی	
۰/۱۸	۱۰۰	۱۴۹۱	کشاورزی	کشاورزی

مأخذ: مجموعه داده‌های پژوهش

در ادامه توزیع فراوانی متغیرهای مستقل مورد بررسی قرار می‌گیرد. شایان ذکر است که متغیر «سن خودرو»، تنها متغیر عددی داده‌های مورد بررسی می‌باشد. سن خودرو، بر بازه صفر تا چهل سال (با افزایش سن خودرو به سمت چهل سال و بالای چهل سال فراوانی بسیار کم می‌باشد) است. در بین متغیرهای گروه خودرو، نوع خودرو و نوع پلاک، طبقه‌های سواری (Car)، خودروهای ایرانی و چینی (LM) و پلاک شخصی (Private) از بیشترین مقدار فراوانی برخوردار هستند. در مورد متغیر سن خودرو نیز هر چقدر به سمت خودروهای قدیمی‌تر حرکت می‌کنیم؛ از فراوانی آنها کاسته می‌شود و خودروهای بالای ۳۵ سال، بیشتر شامل گروه خودروهای بارکش مانند کامیون و تریلر می‌باشند. متغیر خسارت شامل دو کلاس خسارت دیده (Damaged) و خسارت ندیده (Not Damaged) می‌باشد که به ترتیب با کدهای ۱ و ۰ نامگذاری شده‌اند. شکل (۱) توزیع داده‌ها در دو کلاس فوق را نمایش می‌دهد.



شکل ۱. توزیع فراوانی متغیر وابسته

از مجموع ۴۱۵۶۸۷ نمونه، کلاس خسارت دیده‌ها ($D=1$) نمونه مثبت (Positive)، شامل ۱۱۶۴۱ نمونه و کلاس خسارت ندیده‌ها ($D=0$) که نمونه منفی (Negative) هستند، شامل ۴۰۴۰۴۶ نمونه است. به عبارتی نمونه‌های خسارت دیده ۲/۸ درصد و نمونه‌های خسارت ندیده ۹۷/۲ درصد از داده‌ها را شامل می‌شود. با توجه به اینکه نسبت توزیع دو کلاس ۳۵:۱ است، این مسئله یک مسئله طبقه‌بندی باینری نامتعادل^۱ محسوب می‌شود.

این عدم تعادل در داده‌ها در رشته بیمه شخص ثالث، امری متداول و معمول می‌باشد. از طرفی تعداد خودروهایی که دچار خسارت نمی‌شوند به مراتب بیشتر از خودروهای خسارت دیده می‌باشند و از طرف دیگر بیمه‌گذاران برای خسارت‌های جزئی به شرکت مراجعه نمی‌کنند تا بتوانند از تخفیف عدم خسارت در سال‌های آتی استفاده نمایند.

1. Imbalanced Binary Classification Problem

در ادامه نمودار پراکندگی^۱ نمونه‌ها براساس گروه خودرو و دیگر ویژگی‌ها به تفکیک دو کلاس مثبت و منفی مورد بررسی قرار گرفت. این بررسی مشخص می‌کند در کدام گروه خودرو، نوع خودرو یا نوع پلاک یا سن خودرو؛ چه تعداد نمونه‌ای با خسارت و بدون خسارت وجود دارد و اطلاعاتی در مورد پراکندگی داده‌ها و عدم وجود نمونه در برخی گروه‌ها را مشخص می‌کند. در انتهای این بخش، پراکندگی داده‌ها در فضای ویژگی مورد توجه قرار گرفته است. بررسی پراکندگی نمونه‌ها در فضای ویژگی چهاربعدی نشان می‌دهد که از مجموع ۴۱۵۶۸۷ نمونه داده، تنها ۱۵۱۱ نقطه یکتا^۲ وجود دارد. منظور از نقاط یکتا سطرهایی از مجموعه داده می‌باشند که حداقل به وسیله ارزش مقداری یک متغیر مستقل از یکدیگر متمایز می‌شوند. دلیل اصلی این مشکل، تعداد کم ویژگی‌ها و نوع کدگذاری آن‌ها است. در واقع می‌توان نتیجه گرفت، برای آنکه به سطرهای یکتای بیشتری دست پیدا کنیم می‌بایست متغیرها و ویژگی‌های متنوع‌تری جمع‌آوری گردد. همچنین با توجه به متغیرهایی که در پایگاه داده‌ها در دسترس می‌باشد، این نتیجه حاصل می‌شود که مشخصات فردی راننده در تعیین سطح ریسک مشتریان و تمایز سطرهای داده بسیار تأثیرگذار می‌باشد که امید است با توجه به رویکرد قانون جدید بیمه شخص ثالث مبنی بر صدور بیمه‌نامه شخص ثالث براساس ویژگی‌های راننده (بهادر، استاد رمضان و خانی‌زاده، ۱۳۹۶)، در تحقیقات آتی برطرف گردد. در بخش بعد به چگونگی برخورد با این شرایط و مدل‌های استفاده‌شده در این وضعیت خواهیم پرداخت.

1. Scatter plot
2. Unique points

۶. روش کار و مدل‌ها

همان‌طور که در بخش قبل اشاره شد، در مجموعه داده‌های بیمه شخص ثالث، نسبت تعداد بیمه‌نامه‌هایی که منجر به پرداخت خسارت شده به موارد بدون خسارت در حدود ۱ به ۳۵ است. این نسبت آماری در صنعت بیمه طبیعی می‌باشد و به معنی توزیع نامتعادل داده‌ها^۱ در دو کلاس مثبت (خسارت دیده مالی) و منفی (خسارت ندیده مالی) است. این پدیده در تحلیل داده سبب ایجاد خطای آریبی طبقه‌بندی می‌شود. به عبارتی، با توجه به اینکه تعداد نمونه مثبت، $2/8$ درصد و تعداد نمونه منفی، $97/2$ درصد داده‌ها است، چنانچه خروجی طبقه‌بندی به ازای تمام نمونه‌ها ثابت و منفی باشد، دقتی معادل $97/2$ درصد خواهیم داشت که علی‌رغم دقت قابل توجه فاقد هرگونه ارزش عملیاتی است. در واقع در این حالت مدل، یادگیری خود را تنها براساس خروجی‌های مربوط به داده‌های خسارت ندیده انجام می‌دهد و داده‌های خسارت، نادیده گرفته می‌شوند.

یک روش شناخته‌شده برای تحلیل داده‌ها با توزیع نامتعادل، ارزیابی طبقه‌بندی داده‌ها در زیر مجموعه‌های متعادل^۲ است که به صورت تصادفی از مجموعه اصلی نمونه‌برداری^۳ شده است (Doucette and Heywood, 2008). در این روش با توجه عدم توازن بین نمونه‌های موجود در بین دو کلاس برچسب‌گذاری شده، از بین کلاسی که تعداد مشاهدات آن به میزان چشمگیری از کلاس دیگر بیشتر می‌باشد با روش نمونه‌گیری تصادفی، نمونه‌ای به تعداد رکوردهای موجود در کلاس دیگر به دست خواهیم آورد که با توجه به برابری مشاهدات در هر دو مجموعه نتایج تحلیل‌ها قابل اتکاء می‌باشند. در ادامه برای ساخت مجموعه داده متعادل، نمونه‌های مثبت به صورت ثابت و نمونه‌های منفی با تعداد مساوی به صورت تصادفی از مجموعه اصلی

1. Imbalanced Data
2. Balanced subsamples
3. Subsampling

نمونه برداری می‌شود. این کار ۵۰ بار تکرار شده و نتایج طبقه‌بندی به ازای هر تکرار ثبت می‌شود. نتایج، شامل ماتریس آشفتگی، دقت^۱ (ACC)، نرخ تشخیص صحیح نمونه‌های مثبت^۲ (TPR)، نرخ تشخیص صحیح نمونه‌های منفی^۳ (TNR) و معیار اف وان (F1) است.

در تحقیق حاضر، از الگوریتم‌های طبقه‌بندی یادگیری ماشین و داده‌کاوی بهره گرفته‌ایم. برای طبقه‌بندی از درخت تصمیم، نایو بیس، شبکه عصبی و ماشین بردار پشتیبان استفاده شده و پارامترهای هریک به صورت جداگانه محاسبه و نتایج مقایسه شده است. در هر فرایند آموزش، مجموعه داده به صورت تصادفی به دو بخش داده‌های آموزشی و داده‌های آزمون با نسبت ۷۰ به ۳۰ تقسیم گردید. در حالت نامتعادل تعداد نمونه‌های آزمون ۱۲۴۷۰۷ نمونه و در حالت متعادل تعداد نمونه‌های آزمون ۶۹۸۵ نمونه می‌باشد.

لازم به ذکر است در استفاده از مدل درخت تصمیم و جهت تعیین عمق بهینه درخت تصمیم برای هر بار نمونه‌گیری نمودار ROC-AUC (Gajowniczek, et al., 2014),... (Bowers and Zhou, 2019) نسبت به عمق درخت محاسبه و ترسیم می‌شود و میانگین مقادیر بهینه در ۵۰ نمونه‌گیری تصادفی محاسبه شده و به عنوان نتیجه نهایی در نظر گرفته می‌شود. در ادامه به ارزیابی و بررسی مدل‌های استفاده شده می‌پردازیم و در نهایت نتایج به دست آمده از این تحقیق ارائه می‌گردد.

۷. ارزیابی مدل‌ها

در ابتدا دو مدل درخت تصمیم و شبکه عصبی را بر روی حالتی که داده‌ها نامتعادل هستند بررسی می‌کنیم. برای شبکه عصبی از ۱۰ نرون برای لایه پنهان

1. Accuracy
2. True Positive Rate (TPR)
3. True Negative Rate (TNR)

استفاده شده است. جدول ۵ و ۶ به ترتیب نتایج مربوط به درخت تصمیم و شبکه عصبی را نمایش می‌دهند. همانطور که در بخش قبل اشاره شد؛ ۷۰ درصد از داده‌ها به عنوان مجموعه آموزشی در نظر گرفته می‌شود. همان‌طور که قابل پیش‌بینی نیز بود؛ در حالی که محاسبات بر روی داده‌های نامتعادل انجام گیرد، مدل از دقت بسیار بالایی (۹۷٪) برخوردار می‌باشد. لیکن در این شرایط معیار دقت (ACC) ابزار مناسبی برای تشخیص کارایی مدل محسوب نمی‌شود. در واقع همان‌طور که در جداول ۵ و ۶ مشاهده می‌شود؛ نرخ تشخیص صحیح نمونه‌های منفی ۱۰۰٪ می‌باشد ولی در مقابل، نرخ تشخیص صحیح نمونه‌های مثبت صفر درصد است. نتیجه فوق بیانگر این موضوع است که مدل استفاده شده از نمونه‌های مثبت صرف نظر کرده و یادگیری را تنها براساس نمونه‌های منفی انجام داده است. معیار مناسبی که برای ارزیابی مدل بر روی داده‌های نامتعادل و یا مدل‌هایی که FP^۱ و FN^۲ آن‌ها هم‌ارزش نیستند؛ قابل استناد می‌باشد، معیار F1 است که در جداول ۵ و ۶، صفر می‌باشد و بیانگر ضعف و قابلیت پایین مدل در پیش‌بینی داده‌های جدید می‌باشد. در واقع معیار F1 تابعی از صحت^۳ و فراخوانی^۴ می‌باشد. پس نیاز است برای درک بهتر معیار F1 ابتدا روابط مربوط به صحت و فراخوانی را مشاهده کنیم:

$$Precision = \frac{True\ Positive\ (TP)}{True\ Positive + False\ Positive\ (FP)} \quad (1)$$

$$Recall = \frac{True\ Positive\ (TP)}{True\ Positive + False\ Negative\ (FN)} \quad (2)$$

1. False Positive
2. False Negative
3. Precision
4. Recall

اگر به منجر کسرها توجه کنیم، برای رابطه صحت تمام مقادیری که مثبت پیش‌بینی شده‌اند را خواهیم داشت و در مورد رابطه بازخوانی، تمام مقادیر مثبت واقعی به دست می‌آید. بنابراین روابط بالا می‌توان به صورت ذیل بازنویسی کرد:

$$Precision = \frac{True\ Positive\ (TP)}{Total\ Predicted\ Positive} \quad (3)$$

$$Recall = \frac{True\ Positive\ (TP)}{Total\ Actual\ Positive} \quad (4)$$

همان‌طور که از روابط بالا می‌توان برداشت کرد، معیار صحت زمانی استفاده می‌شود که برای محقق اهمیت دارد مشخص شود که چه تعداد از موارد مثبت پیش‌بینی شده، واقعا مثبت هستند و به درستی پیش‌بینی شده‌اند. کاربرد اصلی این معیار زمانی است که هزینه نتایج مثبت‌های کاذب (FP) بالا باشد. از طرفی دیگر معیار بازخوانی بیانگر تعداد نمونه‌هایی است که از بین کل نمونه‌های واقعی به درستی مثبت پیش‌بینی شده‌اند. به طور مشابه این معیار زمانی استفاده می‌شود که هزینه نتایج منفی‌های کاذب بالا باشد. در نهایت معیار F1 ابزاری است که بین معیارهای صحت و بازخوانی نوعی تعادل برقرار می‌کند و از رابطه زیر به دست می‌آید:

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (5)$$

جدول ۵. طبقه‌بندی با درخت تصمیم (داده‌های نامتعادل)

TestSet=۱۲۴۷۰ TrainSet=۲۹۰۹۸۰	تعداد داده‌های آموزشی و آزمون
$CM = \begin{bmatrix} TP = 0 & FP = 2 \\ FN = 3526 & TN = 121179 \end{bmatrix}$	ماتریس آشفتگی
$ACC = \frac{TP + TN}{TP + FP + TN + FN} = 0.97$	معیار دقت
$TPR = \frac{TP}{TP + FN} = 0.00$	نرخ تشخیص صحیح نمونه‌های مثبت
$TNR = \frac{TN}{FP + TN} = 1.00$	نرخ تشخیص صحیح نمونه‌های منفی
$F_1 = 0.00$	معیار F1

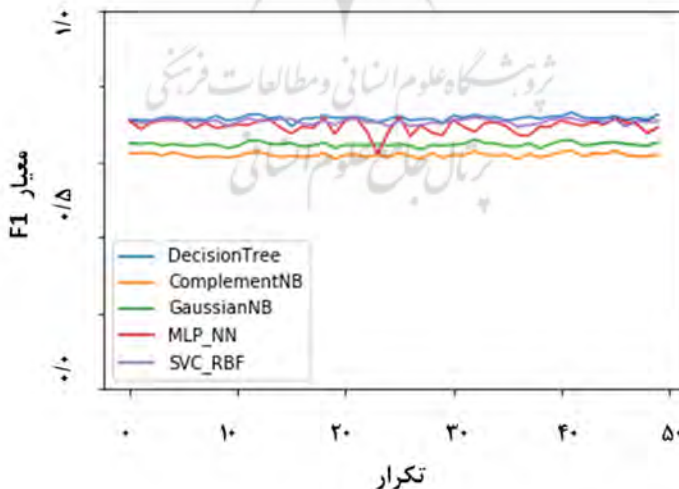
مأخذ: یافته‌های پژوهش

جدول ۶. طبقه‌بندی با شبکه عصبی (داده‌های نامتعادل)

تعداد داده‌های آموزشی و آزمون	۱۲۴۷۰ TestSet= ۲۹۰۹۸ TrainSet=
ماتریس آشفتگی	$CM = \begin{bmatrix} TP = 0 & FP = 0 \\ FN = 3526 & TN = 121181 \end{bmatrix}$
معیار دقت	ACC = 0.97
نرخ تشخیص صحیح نمونه‌های مثبت	TPR = 0.00
نرخ تشخیص صحیح نمونه‌های منفی	TNR = 1.00
معیار F1	F ₁ = 0.00

مأخذ: یافته‌های پژوهش

در ادامه برای رفع مشکل ناشی از نامتعادل بودن داده‌ها و بهبود نتایج، روش نمونه‌گیری از مجموعه داده‌ها استفاده شده است. لذا ابتدا با نمونه‌گیری تصادفی، تعداد ۵۰ مجموعه داده متعادل تولید کرده و مدل‌ها را روی هر مجموعه جداگانه محاسبه کرده و برآیند نتایج (منظور میانگین‌گیری روی اندازه‌ها برای ۵۰ بار تکرار انجام دادید؟) را به عنوان خروجی نهایی در نظر می‌گیریم. در شکل (۲) معیار F1 برای مدل‌های درخت تصمیم، نایو بیز، ماشین بردار پشتیبان و شبکه عصبی بازنه تکرار قابل مشاهده می‌باشد.



شکل ۲. معیار F1 برای ۵۰ تکرار مدل‌ها روی داده‌های متعادل (نمونه‌گیری تصادفی)

چنانچه از نتایج شکل (۲) مشخص می‌شود، با کاربرد داده‌های متعادل عملکرد مدل‌ها بهبود چشمگیری داشته و تقریباً تمام مدل‌ها نتایج نزدیک به نرخ ۷۰ درصد دارند. همان‌طور که در شکل (۲) مشاهده می‌شود درخت تصمیم نسبت به سایر مدل‌های طبقه‌بندی، کارایی بهتری را از خود نشان می‌دهد. در جداول (۷) تا (۱۰) نتایج مربوط به هر مدل (به ترتیب شبکه عصبی، ماشین بردار پشتیبان، درخت تصمیم، نایو بیز) به‌طور جداگانه ارائه گردیده است.

جدول ۷. ارزیابی داده‌های متعادل شده (۵۰ تکرار با نمونه‌های تصادفی) - شبکه عصبی

TestSet=۶۹۸۵ TrainSet=۱۶۲۹۷	تعداد داده‌های آموزشی و آزمون
$CM = \begin{bmatrix} TP = 3200 & FP = 2382 \\ FN = 311 & TN = 1092 \end{bmatrix}$	ماتریس آشفتگی
$ACC = 0.61 \pm 0.01$	معیار دقت
$TPR = 0.91 \pm 0.03$	نرخ تشخیص صحیح نمونه‌های مثبت
$TNR = 0.32 \pm 0.03$	نرخ تشخیص صحیح نمونه‌های منفی
$F_1 = 0.70 \pm 0.01$	معیار F1

مأخذ: یافته‌های پژوهش

جدول ۸. ارزیابی داده‌های متعادل شده (۵۰ تکرار با نمونه‌های تصادفی) - ماشین بردار پشتیبان

TestSet=۶۹۸۵ TrainSet=۱۶۲۹۷	تعداد داده‌های آموزشی و آزمون
$CM = \begin{bmatrix} TP = 3351 & FP = 2543 \\ FN = 119 & TN = 972 \end{bmatrix}$	ماتریس آشفتگی
$ACC = 0.62 \pm 0.01$	معیار دقت
$TPR = 0.95 \pm 0.02$	نرخ تشخیص صحیح نمونه‌های مثبت
$TNR = 0.29 \pm 0.02$	نرخ تشخیص صحیح نمونه‌های منفی
$F_1 = 0.71 \pm 0.01$	معیار F1

مأخذ: یافته‌های پژوهش

TestSet=۶۹۸۵ TrainSet=۱۶۲۹۷	تعداد داده‌های آموزشی و آزمون
$CM = \begin{bmatrix} TP = 3335 & FP = 2551 \\ FN = 143 & TN = 956 \end{bmatrix}$	ماتریس آشفتگی
ACC = 0.62 ± 0.01	معیار دقت
TPR = 0.95 ± 0.01	نرخ تشخیص صحیح نمونه‌های مثبت
TNR = 0.29 ± 0.02	نرخ تشخیص صحیح نمونه‌های منفی
$F_1 = 0.72 \pm 0.01$	معیار F1

مأخذ: یافته‌های پژوهش

جدول ۱۰. ارزیابی داده‌های متعادل شده (۵۰ تکرار با نمونه‌های تصادفی) - نایو بیز گوسی

TestSet=۶۹۸۵ TrainSet=۱۶۲۹۷	تعداد داده‌های آموزشی و آزمون
$CM = \begin{bmatrix} TP = 2689 & FP = 2055 \\ FN = 830 & TN = 1411 \end{bmatrix}$	ماتریس آشفتگی
ACC = 0.58 ± 0.00	معیار دقت
TPR = 0.76 ± 0.01	نرخ تشخیص صحیح نمونه‌های مثبت
TNR = 0.41 ± 0.01	نرخ تشخیص صحیح نمونه‌های منفی
$F_1 = 0.65 \pm 0.00$	معیار F1

مأخذ: یافته‌های پژوهش

باتوجه به جداول فوق ملاحظه می‌گردد که معیار دقت در هر چهار مدل کاهش پیدا کرده؛ لیکن نرخ تشخیص صحیح نمونه‌های مثبت در حال بهبود و افزایش است که نشانه بهبود مدل نیز می‌باشد. به طور کلی در مسائل طبقه‌بندی، پیش‌بینی درست نمونه‌های مثبت از ارزش زیادی برخوردار می‌باشد. همچنین همان‌طور که در بخش‌های قبل اشاره شد، معیار F1 یکی از ابزارهای مناسب جهت ارزیابی مدل‌های مورد استفاده می‌باشد. همان‌طور که در جداول (۷) تا (۱۰) دیده می‌شود در تمامی مدل‌ها، معیار F1 افزایش یافته و باعث بهبود مدل شده است. در این میان مدل درخت

تصمیم از مقدار F1 بالاتری نسبت به سایر مدل‌ها برخوردار است. جدول (۱۱) عملکرد چهار مدل را براساس معیار F1 نمایش می‌دهد.

جدول ۱۱. مقایسه دقت مدل‌ها

معیار F1	مدل
$0/65 \pm 0/00$	نایو بیز
$0/70 \pm 0/01$	شبکه عصبی
$0/71 \pm 0/01$	ماشین بردار پشتیبان
$0/72 \pm 0/01$	درخت تصمیم

مأخذ: یافته‌های پژوهش

در انتها برای مطالعه تاثیر هر یک از ویژگی‌ها اعم از نوع خودرو، گروه خودرو، نوع پلاک و سن خودرو روی نتایج طبقه‌بندی از مدل درخت تصمیم استفاده شده است. مسیر درخت تصمیم از گره اول (ریشه) تا یکی از گره‌های برگ که نمایانگر یکی از قوانین موجود در درخت تصمیم می‌باشد مورد بررسی قرار می‌گیرد. با توجه به نمودار درخت تصمیم اگر نوع خودرو غیر از خودروهای سواری، اتوبوس، ون، کامیون، وانت و سایر اتوکارها باشد؛ آنگاه برای ارزیابی دقیق‌تر ریسک می‌بایست نوع پلاک خودرو نیز مورد بررسی قرار گیرد. اگر نوع پلاک شامل یکی از دسته‌های «خصوصی»، «فاقد پلاک»، «عمومی» و یا «منطقه آزاد» باشد؛ آنگاه به سراغ متغیر گروه خودرو خواهیم رفت. در این حالت اگر گروه خودرو «سواری»، «اتوکار» و یا «بارکش» باشد، سن خودرو را مورد بررسی قرار می‌دهیم و در صورتی که سن خودرو کمتر از ۴۳ سال باشد، پیش‌بینی می‌شود که خودرو خسارت نمی‌بیند و بنابراین بیمه‌گذار و خودرو مورد نظر به طبقه کم‌ریسک تعلق داشته و باید حق بیمه کمتری از وی دریافت گردد. شایان ذکر است برای تعیین دقیق حق بیمه نیاز است تحقیق مشابهی بر روی خسارت‌های جانی شرکت بیمه انجام پذیرد و برآیند نتایج آن با خروجی مقاله فعلی می‌تواند نرخ کاربردی‌تری برای شرکت‌های بیمه ارائه دهد. جمع‌بندی میزان تاثیرگذاری متغیرها در وقوع خسارت مالی به ترتیب اولویت عبارتند از: نوع خودرو، نوع پلاک، سن خودرو و گروه خودرو.

۸. جمع‌بندی و پیشنهادها

پژوهش فوق در صدد آن بود که با توجه به اهمیت زیاد رشته بیمه شخص ثالث در صنعت بیمه ایران، نتایجی در ارتباط با تأثیرگذاری متغیرهای موجود در احتمال وقوع خسارت تنها در نوع مالی و نه جانی ارائه دهد. در همین راستا از الگوریتم‌های یادگیری ماشین و داده‌کاوی جهت تحلیل داده‌ها بهره گرفته شد. مجموعه داده‌های موجود شامل ۴۱۵۶۸۷ نمونه به همراه پنج متغیر بود که متغیرهای مستقل عبارتند از نوع پلاک، گروه خودرو، نوع خودرو و سن خودرو و متغیر وابسته شامل یک متغیر دو ارزشی (باینری) با مقدار صفر برای خودروهای خسارت ندیده (نمونه منفی) و مقدار یک برای خودروهای خسارت دیده (نمونه مثبت) می‌باشد. در میان متغیرهای مستقل به غیر از سن خودرو که یک متغیر عددی می‌باشد، از نوع اسمی هستند.

در این مقاله سعی شد از اکثر الگوریتم‌های طبقه‌بندی شناخته‌شده یادگیری ماشین استفاده شود و از این لحاظ، پژوهش حاضر دارای جامعیت قابل قبولی می‌باشد. نتایج، یافته‌ها و پیشنهادات پژوهش را می‌توان به شرح زیر خلاصه نمود:

- با بررسی داده‌ها مشخص گردید که از میان ۴۱۵۶۸۷ نمونه تنها ۱۵۱۱ نمونه یکتا وجود دارد. براساس این نتیجه، می‌توان پیش‌بینی کرد که متغیرهای بیشتری در محاسبه احتمال وقوع خسارت تأثیرگذارند. همچنین با توجه به نوع داده‌های در دسترس که تنها شامل ویژگی‌های خودرو می‌باشد؛ این نتیجه حاصل می‌شود که ویژگی‌های شخصی راننده می‌تواند در افزایش نقاط و نمونه‌های یکتا بسیار تعیین‌کننده باشد. شایان ذکر است از آنجاکه در سال‌های گذشته بیمه‌نامه‌های شخص ثالث تنها براساس ویژگی‌های خودرو صادر می‌گردید؛ لذا دسترسی به اطلاعات فردی رانندگان بسیار دشوار و غیر قابل اتکا می‌باشد. با این حال تیم نویسنده این مقاله، در حال جمع‌آوری این اطلاعات از طریق تجمیع داده‌های برخی شرکت‌های بیمه می‌باشد و امیدوار است در آینده‌ای نزدیک تحقیقی در همین راستا و بر روی متغیرهایی شامل اطلاعات راننده انجام شود که مکمل مقاله فوق نیز خواهد بود.

- با مرور نتایج ارزیابی مدل‌های استفاده شده در تحقیق، مشاهده می‌شود که مدل درخت تصمیم از کارایی بهتر و قدرت پیش‌بینی بالاتری برخوردار می‌باشد ($F1=0.72 \pm 0.01$). این امر می‌تواند نتیجه‌ای مثبت برای فعالان صنعت بیمه باشد؛ چرا

که درخت تصمیم هم از لحاظ ارائه، قابلیت ارائه ساده و تصویری را داشته و هم قادر است بین ویژگی‌های موجود، اولویت‌بندی مناسبی را براساس میزان تأثیرگذاری ویژگی‌ها انجام دهد (و قدرت تفسیرپذیری بالاتری هم دارد). در انتها نتیجه به دست آمده در رابطه با اولویت‌بندی متغیرهای استفاده‌شده در تحقیق اشاره می‌گردد.

• جهت تعیین اولویت متغیرهای مستقل موجود در مجموعه داده‌ها برای هر یک از ۵۰ نمونه تصادفی، درخت‌های تصمیم رسم گردید و بر اساس میزان فراوانی متغیرها در گره‌ها و سطوح مختلف درخت، میزان اهمیت و اولویت متغیرهای مستقل به ترتیب اولویت عبارتند از:

۱. نوع خودرو؛
۲. نوع پلاک؛
۳. سن خودرو؛
۴. گروه خودرو.

در انتها مجدداً شایان ذکر است که برای دستیابی به یک سیستم ارزیابی ریسک دقیق لازم است که شرکت‌های بیمه در صحت گردآوری اطلاعات شخصی رانندگان نیز کوشا باشند.

پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی

منابع

- اصغری اسکوئی، محمدرضا، (۱۳۹۴). کاربرد روش پنجره لغزان برای انتخاب ساختار شبکه عصبی با تأخیر زمانی در پیش‌بینی سری‌های زمانی مالی. فصلنامه پژوهشنامه اقتصادی، ۱۵(۵۷)، صص ۷۵-۱۰۸.
- اصغری اسکوئی، محمدرضا و قاسم‌زاده، محمد، (۱۳۹۵). کاربرد قواعد کشفی و الگوریتم ژنتیک در ساخت مدل *ARMA* برای پیش‌بینی سری‌های زمانی. فصلنامه مدیریت فناوری اطلاعات، دانشگاه تهران، ۸(۱)، صص ۱-۲۶.
- ایزدپرست، محمود، (۱۳۹۰). دسته‌بندی مشتریان بیمه با استفاده از داده‌کاوی. تازه‌های جهان بیمه، شماره ۱۶۱.
- بهدار، آزاده، استادرمضان، آذین و خانی‌زاده، فرید، (۱۳۹۶). بررسی امکان صدور بیمه‌نامه شخص ثالث بر اساس ویژگی‌های راننده (تبصره ۱ ماده ۱۸ قانون جدید بیمه شخص ثالث) و ارائه آیین‌نامه پیشنهادی. پژوهشکده بیمه.
- ترکستانی، محمد صالح؛ ده‌پناه، آرمان؛ تقوی‌فرد، محمدتقی و شفیعی، شهرام، (۱۳۹۵). ارائه چارچوبی برای اصلاح نرخ حق بیمه در رشته بدنه اتومبیل با استفاده از مدل شبکه‌های عصبی (مطالعه موردی: شرکت بیمه آسیا)، مدیریت فناوری اطلاعات، ۸(۴).
- حاجی‌حیدری، نسترن؛ خاله، سامرند و فراهی، احمد، (۱۳۹۰). طبقه‌بندی میزان ریسک بیمه‌گذاران بیمه بدنه خودرو با استفاده از الگوریتم‌های داده‌کاوی (مورد مطالعه: یک شرکت بیمه). پژوهشنامه بیمه، ۲۶(۴).
- حنفی‌زاده، پیام و رستخیز پایدار، ندا، (۱۳۹۰). مدلی جهت دسته‌بندی ریسکی گروه‌های مشتریان بیمه بدنه اتومبیل بر اساس ریسک با استفاده از تکنیک داده‌کاوی (مورد مطالعه: بیمه بدنه اتومبیل در یک شرکت بیمه‌ای). پژوهشنامه بیمه، ۲۶(۲).

عمرانی نوش آبادی، مصطفی، (۱۳۹۰). ارائه مدل اقتصادسنجی، جهت تعیین حق بیمه بیمه گذار در بیمه شخص ثالث خودرو براساس متغیرهای تأثیرگذار بر آن. پایان نامه کارشناسی ارشد، پژوهشکده بیمه.

فتح‌نژاد، فرامرز و ایزدپرست، محمود، (۱۳۹۰)، ارائه چهارچوب برای پیش‌بینی سطح خسارت مشتریان بیمه بدنه اتومبیل با استفاده از راهکار داده‌کاوی. تازه‌های جهان بیمه، ۱۵۶.

کریم‌زادگان مقدم، داود و بهروان، مجید، (۱۳۹۴)، ارائه راهکاری برای تعرفه‌گذاری پویا در صنعت بیمه با استفاده از تکنیک داده‌کاوی (مورد مطالعه: بیمه شخص ثالث). پژوهشنامه بیمه، ۴.

Baecke, P., and Bocca, L., (2017). The value of vehicle telematics data in insurance risk selection processes. *Decision Support Systems*, 98, 69.

Bowers, A.J. and Zhou, X. (2019). Receiver operating characteristic (ROC) area under the curve (AUC): A diagnostic measure for evaluating the accuracy of predictors of education outcomes. *Journal of Education for Students Placed at Risk*. 24(1), 20-46.

Chawla, N.V. (2009). Data mining for imbalanced datasets: An overview. In *Data mining and knowledge discovery handbook* (875-886). Springer, Boston, MA.

David, M., (2015). Auto insurance premium calculation using generalized linear models. *Procedia Economics and Finance*, 20(15), 147-156.

Doucette, J. and Heywood, M.I. (2008). GP classification under imbalanced data sets: Active sub-sampling and AUC approximation. In *European Conference on Genetic Programming* (266-277). Springer, Berlin, Heidelberg.

Frempong, N.K., Nicholas, N. and Boateng, M.A. (2017). Decision tree as a predictive modeling tool for auto insurance claims. *Int. J. Statist. Appl.*, 7(2), 117-120.

- Gajowniczek, K., Ząbkowski, T. and Szupiluk, R. (2014). Estimating the roc curve and its significance for classification models'assessment. *Metody Ilościowe w Badaniach Ekonomicznych*, 15(2), 382-391.
- Guo, H. and Viktor, H.L. (2004). Learning from imbalanced data sets with boosting and data generation: the databoost-im approach. *ACM Sigkdd Explorations Newsletter*, 6(1), 30-39.
- Kaščelan, V., Kaščelan, L. and Novović Burić, M. (2016). A nonparametric data mining approach for risk prediction in car insurance. *Economic Research-Ekonomska Istraživanja*. 29(1), 545-558.
- Menardi, G. and Torelli, N. (2014). Training and assessing classification rules with imbalanced data. *Data Mining and Knowledge Discovery*, 28(1), 92-122.
- Provost, F. (2000). Machine learning from imbalanced data sets 101. In *Proceedings of the AAI'2000 workshop on imbalanced data sets*, 68(2000), 1-3. AAI Press.
- Salma, D.F., Murfi, H. and Sarwinda, D. (2019). July. The Performance of One-Dimensional Naïve Bayes Classifier for Feature Selection in Predicting Prospective Car Insurance Buyers. *In International Conference on Data Mining and Big Data (124-132)*. Springer, Singapore.
- Thakur, S.S. and Sing, J.K. (2013). Mining Customer's Data for Vehicle Insurance Prediction System using k-Means Clustering-An Application. *International Journal of Computer Applications in Engineering sciences*, 3(4), 148.
- Wang, S., Liu, W., Wu, J., Cao, L., Meng, Q. and Kennedy, P.J. (2016). Training deep neural networks on imbalanced data sets. *In international joint conference on neural networks (4368-4374)*. IEEE.
- Wuyu, S. and Cerna, P. (2019). Risk Assessment Predictive Modelling in Insurance Industry Using Data Mining. *Software Engineering*, 6(4), 121.
- Yunos, Z.M., Ali, A., Shamsuddin, S.M. and Ismail, N. (2016). Predictive Modelling for Motor Insurance Claims Using Artificial Neural Networks. *Int. J. Advance Soft Compu. Appl*, 8(3).

ملاحظات اخلاقی

حامی مالی

این مقاله حامی مالی ندارد.

مشارکت نویسندگان

تمام نویسندگان در آماده‌سازی این مقاله مشارکت کرده‌اند.

تعارض منافع

بنابه اظهار نویسندگان، در این مقاله هیچگونه تعارض منافی وجود ندارد.

تعهد کپی‌رایت

طبق تعهد نویسندگان، حق کپی‌رایت (CC) رعایت شده است.

پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی