



ORIGINAL RESEARCH PAPER

Discovering profitable customers by data mining approach

M. Nezhad Afrasiabi¹, A. Esfahanipour², A.M. Kimiagar^{2,*}

¹ Department of Department of Industrial Engineering, Faculty of Industrial Engineering and Management Systems, Amirkabir University of Technology, Tehran, Iran

² Department of Financial Engineering, Faculty of Industrial Engineering and Management Systems, Amirkabir University of Technology, Tehran, Iran

ARTICLE INFO

Article History:

Received 19 January 2019

Revised 15 May 2021

Accepted 6 June 2021

Keywords:

Behavior of the Insurers

Clustering

Decision tree

Discover of Profitable Customers

K-means

ABSTRACT

BACKGROUND AND OBJECTIVES: Today, customers have become a critical factor in directing investors, producers, and even researchers and innovators. For this reason, organizations need to know about their customers and plan for them. Insurance companies and in general the insurance industry in each country, is one of the most important financial institutions active in financial markets, especially the capital market, which in addition to providing security for economic activities, can play a very fundamental role in providing insurance services. In other words, insurance companies play a vital role in the mobility, dynamic of financial markets and the provision of investable funds in economic activities. In this research, it has been attempted to answer one of the most important questions of insurance organizations, namely, predicting the level of customers' losses and investing on profitable customers.

METHODS: Data mining methods were used to discover knowledge to meet business needs and customer relationship management strategies. In addition, an overview of the various applications of data mining in customer relationship management in various insurance companies has been done. In the model implementation stage, a real data-set is used to evaluate the proposed model. To perform the data mining techniques in the insurance industry as data of customers, the vehicle body insurance from 2015 to 2017 has been under investigation. The total number of data used in this study from the beginning was more than 19,356, which during data preparation using Rapidminer 7.1 software became 19,356. After the initial processing, an attempt is made to extract good features from the 15 variables in the data-set that is tangible and help this research in its goal. As a result, by using clustering, drivers are divided into separate clusters based on the amount of loss, and the characteristics of each cluster are expressed. In the clustering section, three algorithms of data mining are examined. First, k-means, k-medoids, and DBSCAN implemented on data-set. Then, the conclusion of three algorithms compared with each other based on the time of calculation and accuracy.

FINDINGS: Data mining was a good tool in this research, owing to the large volume of data, to discover the needs and identify customers. The data mining technique which was the main approach of this study fully covered the information needs by methods such as classification, prediction and clustering. The k-means algorithm was selected as the most optimal one in time and accuracy. In the following, the implementation of the algorithms in the modeling step, the decision tree algorithm was selected and by the decision tree related to the forecasting model, it can be predicted future customers by what characteristics would be in what category. It will be valuable for the insurance companies. Using a decision tree, a forecasting model is proposed to help insurance companies to identify profitable customers which can be used for future plans of organizations.

CONCLUSION: The customer plays an important role in today's industry. Through studying the data obtained from customer behavior, appropriate action can be taken for marketing-related planning and customer acquisition. The use of predictive models and preventive roadmaps has always been one of the goals of the tools that various organizations have been looking for. In this research, the insurance industry as one of the most important pillars of economic in developing countries has been chosen. By reviewing the share of the insurance industry in the economy of a developing country, it can be seen that insurance has a significant role compared to other services. In this study, the role of insurance companies in optimizing the investment process and ways to expand the interaction between insurance examined. Customers can lead to the growth and development of the insurance industry and the capital market and thus the growth and development of the national economy. Therefore, in the implementation of this research, the data of insurance customers have been used and a forecasting model has been presented. As a good prediction model, the decision tree with 86.21% accuracy was the best model that reached in this study. The insurers' income criterion is considered as the node root, which shows the used method can help insurance companies make more profit by focusing on profitable customers.

*Corresponding Author:

Email: kimiagar@aut.ac.ir

Phone: +9821 64545389

ORCID: [0000-0001-7216-0685](https://orcid.org/0000-0001-7216-0685)

DOI: [10.22056/ijir.2021.03.03](https://doi.org/10.22056/ijir.2021.03.03)





مقاله علمی

کشف مشتریان سودآور با رویکرد داده‌محور

مریم نژاد افراسیابی^۱، اکبر اصفهانی‌پور^۲، علی محمد کیمیاگری^{۳*}

^۱ گروه مهندسی صنایع، دانشکده مهندسی صنایع و سیستم‌های مدیریت، دانشگاه صنعتی امیرکبیر، تهران، ایران

^۲ گروه مهندسی مالی، دانشکده مهندسی صنایع و سیستم‌های مدیریت، دانشگاه صنعتی امیرکبیر، تهران، ایران

اطلاعات مقاله

تاریخ‌های مقاله:

تاریخ دریافت: ۲۹ دی ۱۳۹۷

تاریخ داوری: ۲۵ اردیبهشت ۱۴۰۰

تاریخ پذیرش: ۱۶ خرداد ۱۴۰۰

کلمات کلیدی:

خوشه‌بندی

درخت تصمیم

رفتار بیمه‌گذاران

کشف مشتریان سودآور

کی‌مینگ

*نویسنده مسئول:

ایمیل: kimiagar@aut.ac.ir

تلفن: +۹۸۲۱ ۶۴۵۴۵۳۸۹

ORCID: 0000-0001-7216-0685

DOI: 10.22056/ijir.2021.03.03

چکیده:

پیشینه و اهداف: امروزه مشتریان به عامل بسیار مهم و حیاتی در هدایت سرمایه‌گذاران، تولیدکنندگان و حتی محققان و نوآوران مبدل گشته‌اند. به همین دلیل، سازمان‌ها نیاز دارند مشتریان خود را بشناسند و برای آنان برنامه‌ریزی کنند. در این پژوهش، تلاش شده تا به یکی از اساسی‌ترین سؤالات سازمان‌های بیمه‌ای، یعنی پیش‌بینی سطح خسارت مشتریان، پاسخ داده شود.

روش‌شناسی: در پژوهش حاضر از ابزار داده‌کاوی برای داده‌های مشتریان صنعت بیمه، بخش بیمه بدنه خودرو از سال ۱۳۹۴ تا ۱۳۹۶ استفاده شده‌است. تعداد کل داده‌ها که از ابتدا در این پژوهش مورد استفاده قرار می‌گیرد بیش از ۱۹۳۵۶ بوده که در ادامه و در طی آماده‌سازی آن‌ها با استفاده از نرم‌افزار Rapidminer V/1 تعداد داده‌هایی که در نرم‌افزار لحاظ می‌شود ۱۹۳۵۶ است. پس از پردازش اولیه تلاش می‌شود، از بین ۱۵ متغیر موجود در پایگاه داده ویژگی استخراج شود که ملموس باشد و این پژوهش را در هدف خود یاری دهد. بدین منظور با به کارگیری خوشه‌بندی، رانندگان بر اساس میزان مبلغ خسارت به خوشه‌های مجزا تقسیم می‌شوند و ویژگی‌های هر خوشه بیان می‌شود. در قسمت خوشه‌بندی، ابتدا الگوریتم‌های *k-means*، *k-medoids* و *DBSCAN* استفاده شده‌است. سپس الگوریتم‌های بکار رفته به جهت زمان انجام محاسبات و میزان صحت با یکدیگر مقایسه شدند.

یافته‌ها: در نهایت الگوریتم *k-means* به عنوان الگوریتم بهینه برای این مجموعه داده انتخاب شد. در انتها به کمک درخت تصمیم مدلی پیش‌بینی ارائه می‌شود که شرکت‌های بیمه را در جهت سودآوری بیشتر و کشف مشتریان سودآور کمک می‌کند و برای برنامه‌ریزی و تصمیم‌گیری‌های آتی سازمان قابل استفاده است.

نتیجه‌گیری: برای پیش‌بینی، درخت تصمیم، با میزان صحت ۸۶/۲۱٪ بهترین مدلی بود که در این پژوهش به آن رسیدیم و در مدل درخت تصمیم ارائه شده معیار درآمد بیمه‌گذار به عنوان گره ریشه در نظر گرفته می‌شود که همین نکته نشان‌دهنده آن است روش بکار رفته می‌تواند به شرکت‌های بیمه کمک کند تا با تمرکز بر مشتریان سودآور به درآمد بیشتری برسند.

مقدمه

داده‌کاوی جهت کشف مشتریان سودآور بیمه بدنه خودرو انجام شده است. زیرا داده‌کاوی می‌تواند در تصمیم‌گیری‌های حیاتی کسب‌وکار به شرکت‌های بیمه کمک کند و دانش تازه به دست آمده را به نتایج قابل اقدام در کسب‌وکار شامل محصول، بازاریابی، تحلیل توزیع خسارت، مدیریت دارایی- بدهی و تحلیل توانایی بازپرداخت دیون تبدیل کند (Gharakhani and Abolghasemi, 2011).

مروری بر پیشینه پژوهش

(Baecke and Bocca, 2017) به بررسی چگونگی تأثیر پیدایش اینترنت به عنوان یک سنسور داده بر بهبود انتخاب شرکت‌های بیمه و ریسک آن پرداخته‌اند. (Raedel et al., 2017) یک مجموعه داده گسترده از بیمه سلامت ملی آلمان را با هدف بررسی مداخله مجدد در درمان احیا بیماران دندانپزشکی مورد مطالعه قرار داده‌اند. (Saidur Rahman et al., 2017) با تحلیل داده‌های شرکت‌های بیمه عمر و بررسی واکنش‌های مشتریان در مقابل سیاست‌های مختلف بیمه به پیش‌بینی رفتار آینده بیمه گذاران پرداخته‌اند. (Wanke and Barros, 2015) یک مجموعه داده پانل متوازن (Balanced Panel Data) مربوط به شرکت‌های بیمه برزیلی را مورد مطالعه قرار داده‌اند و ناهمسانی (Heterogeneity) در قسمت‌های مختلف را بررسی کرده‌اند. نتایج نشان داد این ناهمسانی در عملکرد تأثیرگذار است. (Sundarkumar and Ravi, 2015) یک روش ترکیبی برای مسائل با داده‌های غیرمتوازن ارائه کردند، که قادر به کشف خودکار کلاهبرداری‌ها در شرکت‌های بیمه می‌باشد. (Oshini and Caldera, 2013) با تمرکز بر اجرای تکنیک‌های حفظ مشتری به موضوع داده‌کاوی در بیمه‌های زندگی و همچنین تحلیل بیمه‌گذاران پرداخته‌اند. آنها به این نتیجه رسیده‌اند که اجرای تکنیک‌های داده‌کاوی در حوزه بیمه‌های عمر، به راحتی می‌تواند از ریزش بیمه‌گذاران جلوگیری کند. (Thakur and Sing, 2013) مطالعه‌ای بر روی مشتریان بیمه خودرو انجام داده و از الگوریتم بهبود یافته‌تری نسبت به روش k-means استفاده کرده‌اند که با توجه به ویژگی‌های خاص مشتریان قدرت پیش‌بینی عملکرد مشتریان را افزایش می‌دهد. (Balaji and Srivatsa, 2012) تکنیک‌های مورد استفاده برای پیش‌بینی داده‌ها برای بیمه‌گذاران بیمه عمر را طبقه‌بندی کرده‌اند. آنها الگوریتم‌های مختلفی چون روش دسته‌بندی نایو بیس و شبکه‌های بیزین را جهت دسته‌بندی داده‌ها مورد ارزیابی قرار داده‌اند. (Ranjan, 2017) با استفاده از روش مورد کاوی، به مطالعه کاربردهای مدیریت ارتباط با مشتری در شرکت‌های بیمه پرداخته است. (Bhowmik, 2011) به بررسی و کشف تخلفات بیمه‌های خودرویی با استفاده از چندین تکنیک کشف تخلف پرداخته‌اند. همچنین با تمرکز بر رفتار بیمه‌گذاران، به دنبال تحلیل مشتریان سودآور برای شرکت‌های بیمه هستند. (Morik and Kopcke, 2004) با استفاده از یک مطالعه موردی موضوع مدیریت حفظ مشتری در صنعت بیمه را مورد بررسی قرار داده و روش‌هایی برای کاهش

فشار فزاینده در زندگی روزمره، موجب رشد تقاضای محصولات بیمه‌ای می‌شود. در این بین داده‌کاوی می‌تواند به شرکت‌های بیمه‌ای در کشف الگوهای مفید نهفته در دل بانک‌های اطلاعاتی مشتریان کمک کند (Umamaheswari and Janakiraman, 2014). هدف از این پژوهش نیز تبیین این موضوع است که داده‌کاوی چگونه در صنعت بیمه می‌تواند مفید باشد، فنون آن چه نتایج در بخش بیمه به دنبال داشته و چگونه تصمیم‌گیری با استفاده از داده‌های بیمه‌ای ممکن می‌شود. از فرآیند داده‌کاوی در صنعت بیمه در مسائلی مانند بهینه‌سازی قیمت‌ها، بهینه‌سازی خدمات، جذب مشتریان جدید، حفظ مشتریان کنونی و کشف کلاهبرداری‌ها در زمینه ادعای خسارت می‌توان استفاده کرد. هر یک از این موارد جای بحث فراوان دارد. برای این که بتوان مشتریان وفادار را در سازمان حفظ کرد و به مشتریان امیدوار و در معرض خطر، خدمات بهتری ارائه داد باید در اولین قدم آنها را شناخت. با شناخت هر دسته از مشتریان می‌توان متناسب با نیاز آنها خدمات شرکت را بهبود بخشید و به این صورت مشتریان را حفظ کرد. پس از تشخیص مشتریان، مشخصات آنها از قبیل خواسته‌ها، نیازها و انتظارات هر طبقه را شناخت و سپس شناخت خدمات شرکت، ایرادات و اتلاف‌های حین خدمات، یافتن فرصت‌ها و تهدیدها می‌تواند سواوری شرکت را بالا برد (Gharanejad, 2010). در ادامه ساختار مقاله به این صورت است که ابتدا پیشینه پژوهش مرور و مبانی نظری تحقیق بیان شده است. پس از آن به روش‌شناسی تحقیق پرداخته شده و توضیحات لازم درباره انواع روش‌های داده‌کاوی ارائه گردیده است. در مرحله بعدی بر روی مجموعه داده‌ها مطالعه انجام شده و روش‌های معرفی شده بر روی مجموعه داده پیاده می‌شود.

مبانی نظری پژوهش

امروزه سیستم‌های بیمه به سرعت در حال پیشرفت هستند و به دلیل افزایش دغدغه در زندگی روزانه، رشد تقاضای شرکت‌های بیمه‌ای افزایش یافته است. صنعت بیمه یکی از حساس‌ترین صنایع به تغییر رفتار و ارتباطات مشتریان است. زیرا در این صنعت، نقشی پررنگی دارد و مشخص‌کننده و جهت‌دهنده سیاست‌های سازمانی و رفتاری شرکت‌های بیمه‌ای است. به همین دلیل، نظارت و کنترل ارتباطات مشتریان در یک شرکت بیمه‌ای امری بسیار مهم است. به نحوی که شرکت‌های بیمه با کنترل کامل چرخه بازاریابی، فروش و خدمات در تمامی رشته‌های بیمه‌ای و رسیدگی دقیق به تمام درخواست‌های مشتریان می‌توانند گامی موثر در جهت حفظ مشتریان بردارند. از این رو، شناسایی مولفه‌های اصلی موثر بر رضایتمندی مشتریان باید در اولویت برنامه‌های شرکت‌های بیمه‌ای قرار گیرد (Motarjem and Niakan, 2020). در این میان داده کاوی به شرکت‌های بیمه کمک می‌کند تا الگوهای مفیدی را از بانک اطلاعات مشتری کشف کنند. بنابراین، پژوهش حاضر با هدف بررسی چگونگی استفاده از

مشتری از دست رفته ارائه داده‌اند.

$$X^* = \frac{X - \text{Min}(X)}{\text{Max}(X) - \text{Min}(X)}$$

۲. نرمال سازی با تابع نرمال استاندارد:

$$X^* = \frac{X - \text{Mean}(X)}{\text{SD}(X)}$$

پس از نرمال سازی داده‌ها با استفاده از الگوریتم‌های نظیر k-means یا Kohonen، کم کردن فاصله‌ی درون خوشه‌ای و زیاد کردن فاصله بین خوشه‌ای مد نظر خواهد بود.

الگوریتم K-means

مهم‌ترین الگوریتم خوشه‌بندی k میانگین یا k-means است که در این پژوهش نیز از آن استفاده شده است. این الگوریتم به‌طور خلاصه از مراحل زیر تشکیل شده است:

✓ تعداد خوشه‌های مدنظر کاربر را از او می‌پرسد (فرض کنیم k خوشه).

✓ به صورت تصادفی k نمونه را به عنوان مراکز خوشه‌های k گانه انتخاب می‌کند.

✓ برای هریک از نمونه‌های دیگر فاصله تا این k مرکز را حساب کرده و آن نمونه را به نزدیک‌ترین مرکز خوشه اختصاص می‌دهد.

✓ برای هریک از خوشه‌ها، مرکز خوشه جدیدی حساب کرده و به قدم ۳ می‌رود.

✓ این کار را آنقدر تکرار می‌کند تا شرط توقف حاصل گردد (Shahrabi, 2012).

الگوریتم K-medoids

در الگوریتم k-means عنصری که به عنوان نماینده خوشه انتخاب می‌شد با توجه به معیار میانگین است. به این صورت که هر بار فاصله تمام اعضا با نماینده محاسبه و میانگین فاصله‌ها به‌عنوان نماینده جدید انتخاب می‌شود.

اما k-medoids الگوریتمی مبتنی بر شی بوده و نماینده خوشه را از میان خود داده‌ها انتخاب می‌کند. معیار انتخاب نماینده عنصر میانه بوده که در این صورت نماینده انتخابی همیشه عضوی از داده‌ها می‌باشد. هدف از انتخاب میانه کم کردن حساسیت خوشه نسبت به مقادیر بسیار بزرگ و یا داده‌های پرت می‌باشد و فاصله اعضای خوشه همیشه با نماینده‌ای که خودش هم عضوی از اعضا می‌باشد سنجیده می‌شود. مراحل اجرای این الگوریتم به صورت زیر است:

✓ مانند k-means ابتدا k نمونه را به عنوان نماینده خوشه‌ها انتخاب می‌شود.

✓ برای هر نمونه نزدیکترین نماینده مشخص و نمونه مربوطه را در آن خوشه قرار داده می‌شود.

روش‌شناسی پژوهش

پژوهش حاضر از نوع داده‌محور است که با استفاده از فرایند استاندارد داده‌کاوی در صنعت (Cross-Industry Standard Process for Data Mining (CRISP-DM) انجام شده است. این مراحل شامل درک مسئله کسب و کار، درک داده‌ها، آماده‌سازی داده‌ها، مدل‌سازی، ارزیابی نتایج و به کارگیری مدل است (Haji Heydari et al., 2011). قدم اول، آماده‌سازی داده‌ها است که در آن انواع داده‌های پرت، ناقص و اشتباه از مجموعه داده حذف و یا اصلاح شدند. در ادامه به معرفی الگوریتم‌های استفاده شده در این تحقیق پرداخته می‌شود. در قسمت مدل‌سازی، از مدل خوشه‌بندی و از الگوریتم k-means استفاده شده است. نتایج این مرحله از آن جهت اهمیت دارند که علاوه بر پی بردن به ویژگی‌های هر خوشه، مقادیر خسارت خوشه‌ها، به‌صورت چند سطحی تعریف شده تا بتوان از آنها در ساخت مدل‌های دسته‌بندی و پیش‌بینی استفاده نمود. لازم به توضیح است که داده‌های مورد استفاده برای مدل‌های دسته‌بندی، به دو دسته آموزشی (Training) و آزمایشی (Testing) تقسیم شده‌اند.

خوشه‌بندی

برای شروع عملیات روی داده‌ها، ابتدا باید آنها خوشه‌بندی شوند. منظور از خوشه‌بندی نیز عملیاتی است که در آن نمونه‌ها و مشاهدات بر اساس ویژگی‌های مشابه بین یکدیگر به دسته‌های گوناگون تقسیم می‌گردند. یک خوشه، شامل یک مجموعه از نمونه‌هاست که بیشترین شباهت را با یکدیگر و بیشترین تفاوت را با خوشه‌های دیگر دارد. نکته قابل توجه این که خوشه‌بندی با دسته‌بندی متفاوت است. زیرا، خوشه‌بندی جزء مدل‌های «غیر هدایت شونده» است. اما دسته‌بندی با توجه به حجم داده‌های موجود، یک هدف را یاد گرفته و به ازای داده‌های جدید آن هدف را پیش‌بینی می‌کند. تفاوت دیگر خوشه‌بندی با دسته‌بندی این است که در خوشه‌بندی هیچ مشخصه هدفی تعریف نمی‌شود. در حالی که در دسته‌بندی همواره بایستی مشخصه هدف موجود باشد و بر اساس مشخصه‌های ورودی، پیش‌بینی شود.

در خوشه‌بندی با کمک الگوریتم‌های آن، کل داده‌ها به زیرگروه‌ها و یا خوشه‌های همگن بخش‌بندی می‌شود. نکته حائز اهمیت این که در خوشه‌بندی تعداد خوشه‌ها به شکل دلخواه انتخاب می‌شود. ولی در مسائل بزرگتر که ابعاد بیشتری دارد، در انتخاب تعداد خوشه‌ها باید به شاخص‌های ارزیابی خوشه‌ها که در ادامه توضیح داده شده، توجه گردد.

برای اجرای خوشه‌بندی و به منظور یکسان کردن اثر داده‌های مختلف، در قدم اول باید مقادیر داده‌ها نرمالیزه (Normalization) شود که عموماً از دو روش زیر برای نرمال کردن داده‌ها استفاده می‌شود (Shahrabi, 2012):

۱. نرمال سازی بیشینه-کمینه:

مدل مدنظر برای تحلیل و تصمیم‌گیری‌های آتی سازمان برگزیده می‌شود. بنابراین قدم‌های صورت گرفته تا حصول نتیجه به صورت زیر است:

- ✓ آماده‌سازی و پالایش داده‌ها
- ✓ خوشه‌بندی به کمک الگوریتم‌های k-means, k-me-
- DBSCAN و doids
- ✓ تقسیم داده‌ها به دو دسته‌ی آموزشی و آزمایشی
- ✓ دسته‌بندی به کمک مدل درخت تصمیم مقایسه‌ی میزان صحت هریک از مدل‌ها جهت تعیین مدل مناسب.

درخت تصمیم

درخت تصمیم، روشی معروف برای دسته‌بندی است که نتایج آن در یک نمودار، شبیه ساختار درخت ارائه می‌شود که هر گره نشانگر یک تست بر روی ارزش مشخصه و هر شاخه، خروجی هر تست را نمایش می‌دهد؛ برگ‌های درخت نیز نمایانگر کلاس‌ها هستند. مفاهیم اصلی در هر درخت تصمیم به شرح زیر می‌باشد:

- ✓ گره (Node): متغیر مستقلی که روی آن آزمون انجام می‌شود.
- ✓ گره ریشه (Root Node): گره‌ای که در بالاترین نقطه‌ی درخت وجود دارد.

- ✓ برگ (Leaf): به برجسب دسته، برگ گفته می‌شود.
- ✓ شاخه (Branch): مقیاسی که خروجی از آن تعیین می‌شود (Ghazanfari et al., 2008).

شکل ۱ ساختار یک درخت تصمیم را نشان می‌دهد که امکان دریافت وام از بانک را برای یک فرد متقاضی نشان می‌دهد. با توجه به این درخت تصمیم بانک برای دادن وام به فرد ابتدا به سطح درآمد او توجه می‌کند و این متغیر به عنوان گره ریشه در نظر گرفته می‌شود. به دنبال آن متغیرهای دیگری نظیر داشتن سوءپیشینه و میزان سابقه فرد در شغل فعلی سایر گره‌های درخت هستند که در تصمیم‌گیری

✓ آن‌گاه در k خوشه ایجاد شده مجدد میانه‌ها به دست آمده و الگوریتم تکرار می‌گردد.

✓ این تکرار را تا زمانی ادامه داده که دیگر میانه‌ها در مرحله جدید با میانه‌های انتخاب شده قبل تفاوت نکنند.

الگوریتم DBSCAN

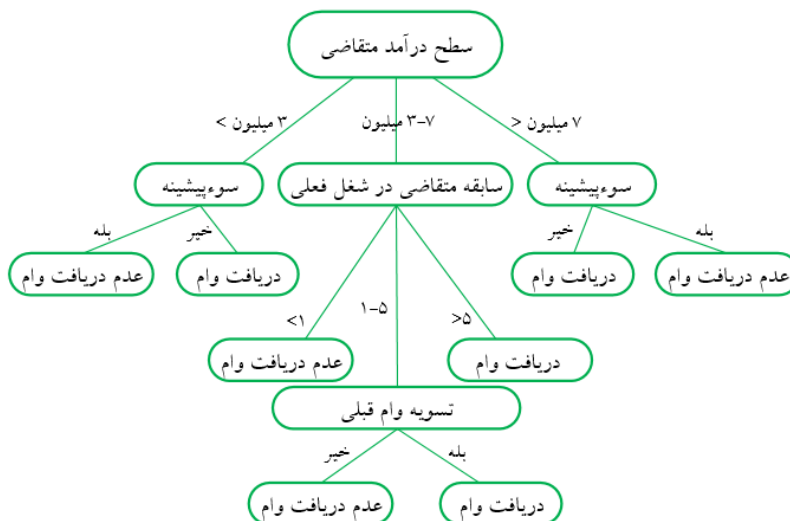
الگوریتم خوشه‌بندی مبتنی بر چگالی یا به اختصار DB-SCAN یکی از الگوریتم‌های خوشه‌بندی مهم و مطرح در داده‌کاوی می‌باشد. این الگوریتم از آن جهت مورد توجه است که برخلاف سایر الگوریتم‌های خوشه‌بندی مانند k-Mean ، k-medoids یا الگوریتم FCM که وابسته به تعداد خوشه می‌باشند و باید از قبل تعداد خوشه‌ها مشخص باشد، خود می‌تواند تعداد خوشه‌ها را مشخص کند و نیازی نیست که تعداد خوشه‌ها به آن اعلام شود. مزیت این روش به نسبت روش‌های دیگری خوشه‌بندی این است که نسبت به شکل داده‌ها حساس نیست و می‌تواند اشکال غیر منظم را نیز در داده‌ها تشخیص دهد.

این الگوریتم نیز مانند دیگر روش‌های خوشه‌بندی نیازمند روشی برای یافتن نزدیکی داده‌ها است. در این الگوریتم می‌توان از فاصله اقلیدسی جهت اندازه‌گیری فاصله (شباهت) استفاده نمود. برای تشریح الگوریتم، نیازمند آشنایی با پارامترهای ϵ و μ می‌باشد که توضیح داده می‌شود:

✓ هر نقطه از داده با نقاط دیگر فاصله‌ای دارد. هر نقطه‌ای که فاصله‌اش با یک نقطه مفروض کمتر از ϵ شد به عنوان همسایه آن نقطه حساب می‌شود.

✓ هر نقطه مفروض که μ همسایه داشته باشد، یک نقطه مرکزی است.

در ادامه از مدل درخت تصمیم استفاده نموده و در نهایت با مقایسه میزان صحت (Accuracy) و اعتبار هر یک از این مدل‌ها،



شکل ۱. نمایش درخت تصمیم برای دریافت وام یک مشتری از بانک

سالی ثبت شده بود که اطلاعات بیمه‌ای جمع‌آوری شده بودند (یکی از سال‌های ۱۳۹۴ تا ۱۳۹۶). بنابراین، در چنین مواردی تنها به خاطر یک ثبت نادرست ناگزیر باید رکورد حذف می‌شد. تعدا قابل توجهی رکورد تکراری نیز در داده‌های این سه سال وجود داشت که اطلاعات یک فرد، اغلب ۲ بار پشت سر هم و در محدود دفعاتی تا ۳ بار پشت سر هم ثبت شده بود که همه این موارد تکراری از پایگاه داده حذف گردیدند. داده‌های پرت از طریق نمودار جعبه‌ای شناسایی شد و به کمک ابزار Subset Worksheet در Minitab و ابزار Filter در Ex-cel، مقادیر نامعتبر از مجموعه داده، شناسایی و حذف گردیدند که در نهایت تعداد ۱۹۳۵۶ رکورد برای انجام مدل‌سازی و تحلیل داده‌ها مورد استفاده قرار گرفتند. اصلاح دیگری که در داده‌ها صورت گرفت، اصلاحاتی بر روی اسامی شهرها بود. از آنجایی که تعداد شهرها در این مجموعه داده بسیار زیاد بود (۲۹۶ نام مجزا)، در صورتی که بدون تغییر باقی می‌ماند، نرم‌افزار، نوع فیلد مورد نظر را تشخیص نمی‌داد و از فرایند مدل‌سازی حذف می‌شدند. بنابراین برای حفظ تأثیر این عامل، نام ۱۶ شهر که بیش از ۷۷ درصد داده‌ها را تشکیل می‌داد، حفظ و مابقی به "Others" تغییر نام یافتند. نام این ۱۶ شهر به شرح جدول ۱ است:

در مورد فیلد «گروه وسیله نقلیه» نیز، به علت تفاوت در فراوانی‌ها، رکوردهای مربوط به آن به ۳ دسته «شخصی»، «بارکش» و «مابقی» تقسیم شدند.

پس از آن تلاش شد، از بین متغیرهای موجود در پایگاه داده ویژگی‌ای استخراج شود که ملموس باشد و پژوهش را در دستیابی به هدف یاری نماید. جمع کل حق بیمه در واقع داده استخراجی از مجموع حق بیمه شکست، حق بیمه سرقت، حق بیمه سیل، حق بیمه پاشیدن رنگ، حق بیمه ماده ۱۰، حق بیمه حوادث، حق بیمه حوادث شخصی، حق بیمه برخورد قطعات، حق بیمه بند ۷ ماده ۹، حق بیمه ایاب و ذهاب، حق بیمه خطر اصلی وسایل اضافی و حق بیمه خطر اصلی می‌باشد. سن نیز داده‌ای است که با توجه به سال تولد مشتریان که در سیستم اطلاعات بیمه ثبت شده و با در نظر گرفتن سال جاری محاسبه شده و به‌عنوان یکی از ویژگی‌های اصلی مورد استفاده قرار گرفت. در ادامه ارتباط بین ویژگی‌های استخراج شده سن، جمع کل حق بیمه و ارزش خودرو مورد توجه قرار گرفت. از آنجا که ارزش خودرو از ملاک‌های مطلوب برای شرکت‌های بیمه است، در ادامه نمودار Scatter plot ارزش خودرو به همراه دو ویژگی استخراج شده سن و جمع کل حق بیمه مطابق شکل ۲ ترسیم شد.

بانک برای دادن وام به فرد متقاضی کمک می‌کند. به‌طور عادی پیچیدگی یک درخت تصمیم با افزایش تعداد مشخصه‌ها افزایش می‌یابد. اگرچه در بعضی شرایط دیده شده است که تنها تعداد کمی از مشخصه‌ها می‌توانند کلاسی را که هر شیء به آن تعلق دارد، تعیین کنند و بقیه مشخصه‌ها کم یا بی‌تأثیرند. اندازه‌گیری کیفیت یک درخت تصمیم، مسئله مهمی است. دقت درخت تعیین شده‌ی دسته‌بندی با استفاده از داده‌های آزمایشی، بطور آشکار یک شاخص مطلوب است.

نتایج و بحث

مجموعه داده‌ها

مجموعه داده‌ها شامل داده‌های مشتریان صنعت بیمه در بخش بیمه بدنه خودرو از سال ۱۳۹۴ تا ۱۳۹۶ است که تعداد ۱۹۳۵۶ داده با استفاده از نرم‌افزار Rapidminer 7.1 آماده‌سازی شد و در نرم‌افزار لحاظ گردیدند. تعداد متغیرهایی که برای این پژوهش در نظر گرفته شد، شامل ۱۵ متغیر سن، جنسیت، شهر، گروه خودرو، نوع خودرو، ارزش خودرو، جمع کل حق بیمه، مقصر، محل حادثه، علت حادثه، نوع حادثه، درصد حادثه، تعداد مصدوم، علت پرداخت و مبلغ خسارت بود. البته، ابتدا بر روی داده‌های مشتریان پردازش صورت گرفت. سپس سه روش مختلف خوشه‌بندی k-means، k-medoids و DBSCAN روی داده‌های آموزشی پیاده شد تا بهترین گزینه از بین آن‌ها انتخاب شود. معیار انتخاب روش بهتر نیز سرعت الگوریتم و خوشه‌بندی منطقی و میزان صحت در داده‌ها بود. برای اعتبارسنجی متقابل (Cross Validation) نیز از 10-fold استفاده گردید.

آماده‌سازی داده‌ها

در این مرحله سعی شد با توجه به داده‌های جمع‌آوری شده تحلیلی آماری صورت گیرد. در این مرحله تحلیل‌های مختلفی می‌توان انجام داد، که در ادامه جزییات آن تشریح می‌شود. رکوردهای موجود در این پایگاه داده، شامل اطلاعات نامعتبر زیادی به شکل داده‌های پرت، داده‌های تکراری، رکوردهای ناقص و رکوردهای اشتباه بود. عملیات آماده‌سازی داده‌ها به کمک نرم‌افزارهای Minitab و Excel صورت گرفت. تعداد بسیار کمی از این رکوردها قابل اصلاح بود و تصحیح شدند. به‌طور مثال سال تولد فرد به‌جای ۱۳۳۳، بصورت ۱۳۳۳۳، نوشته شده بود. البته، تعداد چنین مواردی بسیار کم بود. در بسیاری از رکوردها، در قسمت سال تولد فرد به اشتباه،

جدول ۱. نام شهرها

اهواز	اراک	بندر عباس	بندر انزلی
اصفهان	همدان	کرج	مشهد
قزوین	قم	رشت	شیراز
تبریز	تهران	تنگابن	یزد

مثلا این امر نیز امکان پذیر بود که از تمام فیلدها جهت خوشه بندی استفاده شود که این شکل از خوشه بندی نتایج متفاوت برای هدفی متفاوت را به دنبال خواهد داشت. مثلا در این حالت، شهر بیمه گر، سن، جنسیت و از این حیث ویژگی های فردی که ممکن است در بسیاری از افراد نیز مشترک باشد، ملاکی برای خوشه بندی می شوند و نتایج را تغییر می دهند.

بهبود مدل

تعداد کل داده ها بیمه بدنه که از ابتدا مورد استفاده قرار گرفته بیش از ۱۹۳۵۶ بود. اما پس از آماده سازی داده ها، تعداد داده هایی که برای خوشه بندی در نرم افزار Rapidminer 7.1 لحاظ شد ۱۹۳۵۶ است.

گام یک

ابتدا با استفاده از کلیه ویژگی هایی که از شرکت بیمه استخراج شده بود، درخت تصمیم پیاده سازی شد. نتیجه ای که از اجرای مدل در نرم افزار حاصل شد به این صورت بود که ویژگی سال ساخت خودرو به عنوان ریشه در نظر گرفته شد. از آنجا که به لحاظ منطقی چنین ویژگی نه برای شرکت های بیمه و نه برای مشتریان معیار اصلی برای تصمیم گیری نیست، نباید آن قدر مهم باشد که ریشه قرار بگیرد. در گام بعدی سعی شد ویژگی های مفید نگهداشته شود و بقیه ویژگی ها حذف شوند.

گام دوم

ویژگی های در نظر گرفته شده مطابق جدول ۴ می باشد. نوع ویژگی مطابق گزینه ای است که در نرم افزار انتخاب شده است.

در نهایت متغیرهای سن، جنسیت، شهر، گروه خودرو، نوع خودرو، ارزش خودرو، جمع کل حق بیمه، مقصر، محل حادثه، علت حادثه، نوع حادثه، درصد حادثه، تعداد مصدوم و علت پرداخت به عنوان متغیر مستقل و عوامل تأثیرگذار در سانحه و مبلغ خسارت به عنوان متغیر وابسته مدنظر قرار گرفت (به شرح جدول ۲).

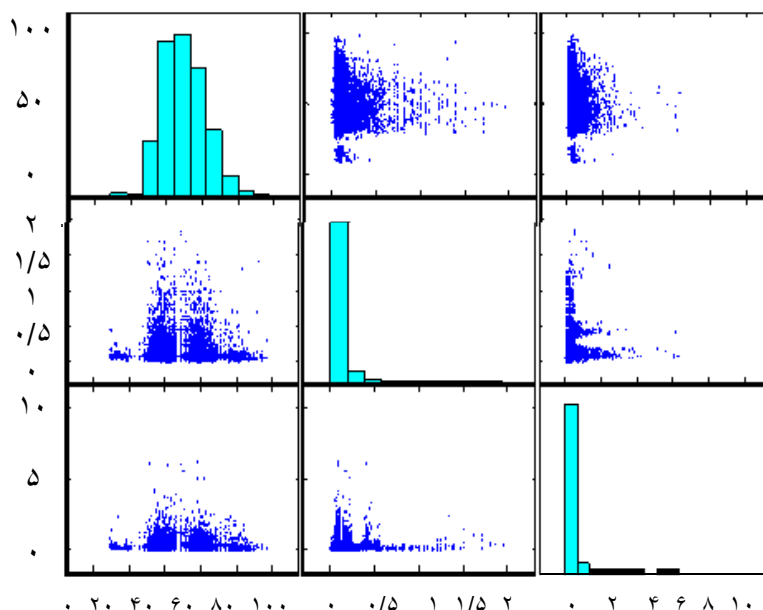
پیاده سازی مدل

قدم اول در خوشه بندی، نرمال کردن داده ها است. اطلاعات مربوط به میزان خسارت در ابتدا از طریق فرمول $x^* = \frac{x - \text{Mean}(X)}{\text{SD}(X)}$ در نرم افزار نرمال شده و سپس از مدل k-means برای خوشه بندی استفاده شده است. در قدم دوم ویژگی مبلغ خسارت به عنوان کلاس هدف در نظر گرفته شد. سپس در مرحله اول روش k-means پیاده سازی شد. مطابق با آن داده ها به سه سطح مطابق جدول ۳ تقسیم شدند. سرعت اجرای این الگوریتم نزدیک به ۱ دقیقه است.

در مرحله دوم روش k-medoids پیاده می شود که در این روش زمان الگوریتم ۳۶ دقیقه است و همه داده ها در یک خوشه قرار می گیرند.

در مرحله سوم الگوریتم DBSCAN با مقدار E برابر ۱ و μ برابر ۵ اجرا می شود. مدت اجرای الگوریتم ۴ دقیقه و ۳۷ ثانیه است و به دنبال آن تعداد خوشه های حاصل شده ۵۵۸ است و بسیار بالاست و هیچ گونه الگوی منطقی برای تحلیل خوشه ها نمی توان یافت.

با مقایسه نتایج سه الگوریتم بکار رفته در این پژوهش، برای خوشه بندی از داده های مربوط به «مبلغ خسارت» استفاده شده است و مطابق با روش k-means به سه سطح تقسیم می شوند که در جدول ۱ قابل مشاهده است. لازم به ذکر است که انجام خوشه بندی، بسته به هدف کار، هدف سازمان و نظر خبره می تواند تغییر کند.



شکل ۲. نمودار اسکتر پلات سن، ارزش خودرو و جمع کل حق بیمه

جدول ۲. داده‌های آماده‌شده و ویژگی آن‌ها

ردیف	ویژگی	نوع ویژگی	سطوح اندازه‌گیری	شرح
1	سن	عددی	فاصله	سن راننده در سال
2	جنسیت	رده‌ای	باینری	زن، مرد
3	شهر	رده‌ای	ترتیبی	محل اقامت راننده
4	گروه خودرو	رده‌ای	ترتیبی	سواری، بارکش، موتورسیکلت، ...
5	نوع خودرو	رده‌ای	ترتیبی	انواع خودرو
6	ارزش خودرو	عددی	فاصله	قیمت خودرو
7	جمع کل حق بیمه	عددی	فاصله	مبلغ حق بیمه
8	مقصر	رده‌ای	ترتیبی	بیمه‌گذار، ناشناخته و ...
9	محل حادثه	رده‌ای	ترتیبی	شهرهای مختلف وقوع حادثه
10	علت حادثه	رده‌ای	ترتیبی	انواع علت حادثه
11	نوع حادثه	رده‌ای	ترتیبی	شکست شیشه، آتش سوزی و ...
12	درصد حادثه	عددی	نسبی	درصدهای مختلف از میزان حادثه
13	تعداد مصدوم	عددی	فاصله	تعداد مصدومین در حادثه
14	علت پرداخت	رده‌ای	ترتیبی	خسارت، کارشناسی و بازیافت
15	مبلغ خسارت	پیوسته	فاصله	انواع هزینه‌ها

جدول ۳. خوشه‌بندی با استفاده از مبلغ خسارت

عنوان خوشه	مبلغ خسارت
سطح 1	[1900350-∞)
سطح 2	[1900350-4427651]
سطح 3	[4427651-∞)

در این پژوهش با سه سطح ایجاد شده ناشی از مرحله خوشه‌بندی، شش حالت ممکن در ماتریس اختلال ایجاد شد که نتایج در **جدول ۶** خلاصه شده است. میزان صحت (Accuracy) مدل به کار رفته مطابق **جدول ۶** نزدیک به ۵۵/۲۸ درصد بدست آمده است. مقدار مربوط به میزان صحت، دقت کلاس (Class Precision) و فراخوانی کلاس (Class Precision) به ترتیب از روابط ۳، ۴ و ۵ پیروی می‌کند [۱۱]:

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN) \quad (۳)$$

$$\text{Precision} = TP / (TP + FP) \quad (۴)$$

$$\text{Recall} = TP / (TP + FN) \quad (۵)$$

نتیجه‌ای که از اجرای مدل در نرم‌افزار حاصل شد به این صورت بود که ویژگی ارزش خودرو به عنوان گره ریشه درخت در نظر گرفته شد.

روش ارائه شده در هر پژوهشی باید از لحاظ اعتبار، مورد ارزیابی قرار گیرد. بنابراین از آنجا که این پژوهش از نوع «داده‌محور» است، میزان صحت هر مدل از طریق ماتریس اختلال (Confusion matrix) سنجیده می‌شود. در حالت کلی اگر دو دسته «درست» و «نادرست» داشته باشیم، چهار حالت ممکن برای نتایج وجود خواهد داشت و ماتریس اختلال مطابق **جدول ۵** می‌باشد. چهار حالت بیان شده در ماتریس اختلال **جدول ۵** شامل به درستی نتیجه آزمون مثبت باشد (TP)، به اشتباه نتیجه آزمون مثبت باشد (FP)، نتیجه آزمون به اشتباه منفی باشد (FN) و نتیجه آزمون به درستی منفی باشد (TN) است

جدول ۴. ویژگی‌های در نظر گرفته شده

ردیف	ویژگی	نوع ویژگی	مقادیر گم شده
1	سن	عدد صحیح	0
2	جنسیت	چندجمله‌ای	0
3	شهر	چندجمله‌ای	0
4	گروه خودرو	چندجمله‌ای	0
5	نوع خودرو	چندجمله‌ای	0
6	ارزش خودرو	عدد صحیح	0
7	جمع کل حق بیمه	عدد صحیح	0
8	مقصر	چندجمله‌ای	0
9	محل حادثه	چندجمله‌ای	0
10	علت حادثه	چندجمله‌ای	0
11	نوع حادثه	چندجمله‌ای	0
12	درصد حادثه	چندجمله‌ای	2
13	تعداد مصدوم	عدد صحیح	9999
14	علت پرداخت	چندجمله‌ای	0
15	مبلغ خسارت	عدد صحیح	0

جدول ۵. ماتریس اختلال

		مشاهده شده	
		درست	نادرست
پیش‌بینی	درست	TP	FP
	نادرست	FN	TN

جدول ۶. ماتریس اختلال حاصل از داده‌های پژوهش

Accuracy	55.28%				
Confusion matrix					
		سطح 1	سطح 2	سطح 3	Class precision
سطح 1	[-∞-1900350]	2810	1942	607	80.36%
سطح 2	[1900350-4427651]	16	28	37	50.47%
سطح 3	[4427651-∞]	533	1337	2689	84.85%
	Class recall	64.32%	96.31%	45.78%	

نهایت خالص درآمد شرکت را به دست آورد. با استخراج ویژگی به نام درآمد یا عایدی شرکت بیمه، تعداد ویژگی‌های در نظر گرفته شده ۱۶ است که در جدول ۷ خلاصه شده‌اند.

با استخراج ویژگی به نام درآمد نرم‌افزار دوباره اجرا شد و نتیجه ماتریس اختلال مطابق جدول ۸ است.

در جدول ۸ مشاهده می‌شود که دقت مدل در هر سه سطح در مقایسه با جدول ۶ بهبود یافته و به حد قابل قبولی رسیده است. علاوه بر آن درخت تصمیم حاصله مطابق شکل ۳ می‌باشد که در آن ویژگی درآمد به‌عنوان گره ریشه در نظر گرفته شده که یک معیار اصلی برای شرکت‌های بیمه محسوب می‌شود. با توجه به این درخت تصمیم مشتریان با جمع کل حق بیمه کمتر و درآمد بالاتر برای

با نتیجه حاصله از نرم‌افزار می‌توان چنین تحلیل کرد که مدل بهبود یافته است. زیرا ریشه درخت معیار ارزش خودرو در نظر گرفته شده که نسبت به حالت مشابه در گام یک معیار بهتری است. اما دقت مدل پایین می‌باشد. به دنبال بهبود مدل در گام بعدی سعی شد در داده‌های مورد استفاده تحلیلی عمیق‌تری صورت گیرد تا ویژگی اضافه مدنظر حذف و ویژگی‌های مفیدتری استخراج گردد و در ادامه مقادیر گم شده در داده‌ها نیز بهبود یابد.

گام سوم

با تحلیل و مطالعه بیشتر در داده‌ها این نتیجه حاصل شد که با استفاده از دو ویژگی مبلغ خسارت و مجموع حق بیمه می‌توان در

جدول ۷. ویژگی‌های نهایی در نظر گرفته شده

ردیف	ویژگی	نوع ویژگی	مقادیر گم‌شده
1	سن	عدد صحیح	0
2	جنسیت	چندجمله‌ای	0
3	شهر	چندجمله‌ای	0
4	گروه خودرو	چندجمله‌ای	0
5	نوع خودرو	چندجمله‌ای	0
6	ارزش خودرو	عدد صحیح	0
7	جمع کل حق بیمه	عدد صحیح	0
8	مقصر	چندجمله‌ای	0
9	محل حادثه	چندجمله‌ای	0
10	علت حادثه	چندجمله‌ای	0
11	نوع حادثه	چندجمله‌ای	0
12	درصد حادثه	چندجمله‌ای	0
13	تعداد مصدوم	عدد صحیح	0
14	علت پرداخت	چندجمله‌ای	0
15	مبلغ خسارت	عدد صحیح	0
16	درآمد	عدد صحیح	0

جدول ۸. ماتریس اختلال مدل بهبود یافته

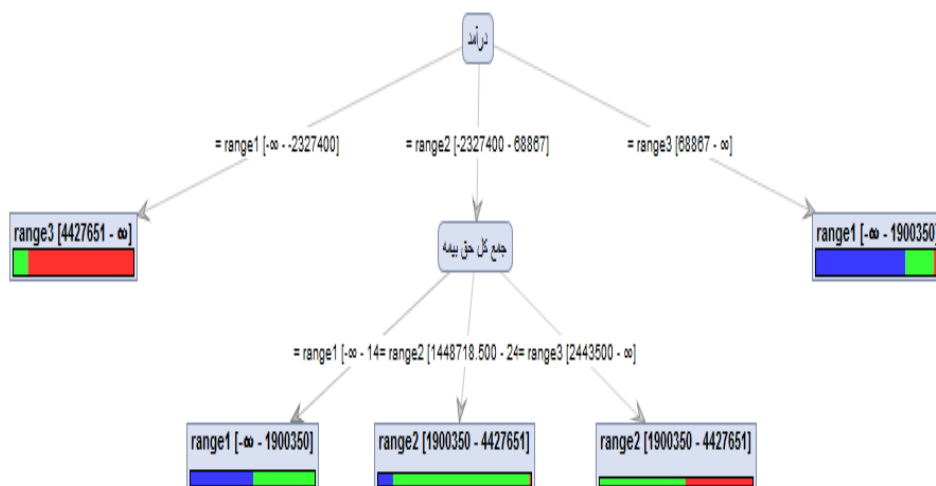
Accuracy	86.21%				
Confusion matrix					
		سطح 1	سطح 2	سطح 3	Class precision
سطح 1	[-∞-1900350]	3241	1040	23	81.67%
سطح 2	[1900350-4427651]	105	2151	82	94.52%
سطح 3	[4427651-∞]	13	116	3228	97.63%
Class recall		97.26%	76.14%	97.69%	

بیمه‌گزاران سود بیشتری به ارمغان می‌آورند.

جمع بندی و پیشنهادها

پژوهش حاضر سعی نمود با تکیه بر فنون مختلف داده‌کاوی، به یکی از نیازهای اساسی صنعت بیمه، یعنی میزان سوآوری حاصل از مشتریان پاسخ دهد. برای این منظور، در ابتدا تعداد بیش از ۱۹۳۵۶ داده از یک شرکت بیمه بدنه خودرو از سال ۱۳۹۴ تا سال ۱۳۹۶ مورد مطالعه قرار گرفت. پس از آن پردازش اولیه روی داده‌ها صورت گرفت و با استفاده از نرم‌افزار Rapidminer 7.1 داده‌ها پاکسازی شدند. در ادامه برای کشف مشتریان سودآور، تنها به یک یا دو الگوریتم اکتفا نشد و الگوریتم‌های k-means، k-medoids و DBSCAN که وظیفه خوشه‌بندی را به عهده دارند، به کار گرفته شد تا ضمن آشنایی با عملکرد آن‌ها و این که

استفاده از روش‌های داده‌کاوی برای تحلیل داده‌های مشتریان در میان شرکت‌های بزرگ و کوچک در سال‌های اخیر رشد فزاینده‌ای داشته است. لذا صاحبان بسیاری از کسب‌وکارها می‌خواهند بدانند چگونه با استفاده از ابزارهایی مانند داده‌کاوی می‌توانند به فروش بیشتری دست یابند و به دنبال آن رویکرد بازاریابی داده‌محور، داده‌کاوی و تحلیل رفتار مشتریان را در کسب‌وکار خود به اجرا در آورند. در این راستا،



شکل ۳. درخت تصمیم (به ترتیب رنگ‌های آبی، سبز و قرمز نشان‌دهنده روند سودآوری بیشتر به کمتر است)

در یک چگونه سازمان را در نیل به هدف یاری می‌کنند، با مقایسه میزان صحت هر یک، الگوریتمی که بیش از سایرین قابل اعتماد است، برای تصمیم‌گیری‌های آتی، انتخاب و مورد بررسی قرار گیرد که نتایج نشان داد از میان سه روش بالا، روش k-means برای خوشه‌بندی بهتر است. برای پیش‌بینی نیز درخت تصمیم با میزان صحت ۸۶/۲۱ درصد بهترین مدلی بود که این پژوهش به آن دست یافت و در مدل درخت تصمیم ارایه شده معیار درآمد بیمه‌گذار به‌عنوان گره ریشه در نظر گرفته شد که این نکته نشان‌دهنده آن است روش به کار رفته می‌تواند به شرکت‌های بیمه کمک کند تا با تمرکز بر مشتریان سودآور به درآمد بیشتری برسند. لازم به توضیح است که ممکن است با کاهش تعداد داده‌ها میزان صحت مدل تغییر کند و حتی افزایش یابد. ولی از آنجایی که از ابتدای انجام این پژوهش سعی بر این بود که یک مساله تجاری حل شود، حجم زیاد داده‌ها بدون کاستن تعداد خاصی از آن‌ها بکار گرفته شد.

تعارض منافع

نویسندگان اعلام می‌دارند که در مورد انتشار این مقاله تضاد منافع وجود ندارد. علاوه بر این، موضوعات اخلاقی شامل سرقت ادبی، رضایت آگاهانه، سوءرفتار، جعل داده‌ها، انتشار و ارسال مجدد و مکرر توسط نویسندگان رعایت شده است.

دسترسی آزاد

کپی‌رایت نویسنده(ها) ©2021: این مقاله تحت مجوز بین‌المللی Creative Commons Attribution 4.0 اجازه استفاده، اشتراک‌گذاری، اقتباس، توزیع و تکثیر را در هر رسانه یا قالبی مشروط به درج نحوه دقیق دسترسی به مجوز CC منوط به ذکر تغییرات احتمالی بر روی مقاله می‌باشد. لذا به استناد مجوز مذکور، درج هرگونه تغییرات در تصاویر، منابع و ارجاعات یا سایر مطالب از اشخاص ثالث در این مقاله باید در این مجوز گنجانده شود، مگر اینکه در راستای اعتبار مقاله به اشکال دیگری مشخص شده باشد. در صورت عدم درج مطالب مذکور و یا استفاده فراتر از مجوز فوق، نویسنده ملزم به دریافت مجوز حق نسخه‌برداری از شخص ثالث می‌باشد.

به‌منظور مشاهده مجوز بین‌المللی Creative Commons Attribution 4.0 به آدرس زیر مراجعه گردد:

<http://creativecommons.org/licenses/by/4.0>

یادداشت ناشر

ناشر نشریه پژوهشنامه بیمه با توجه به مرزهای حقوقی در نقشه‌های منتشرشده بی‌طرف باقی می‌ماند.

انجام خوشه‌بندی، بسته به هدف کار، هدف سازمان و نظر خبره می‌تواند تغییر کند. همان‌طور که عنوان شد در این پژوهش معیار مبلغ خسارت به‌عنوان کلاس هدف برای خوشه‌بندی انتخاب شد که برای پژوهش‌های آتی می‌توان از بین داده‌های موجود معیار دیگری برای خوشه‌بندی انتخاب نمود.

مشارکت نویسندگان

نویسنده اول مسئولیت جمع‌آوری داده‌ها، تدوین روش‌شناسی پژوهش و تدوین مرور ادبیات را داشته است. نویسنده دوم مسئولیت بازنگری نهایی پژوهش را داشته است. نویسنده سوم مسئول بازنگری اولیه پژوهش بوده است.

تشکر و قدردانی

با تشکر از شرکت بیمه بدنه خودرو که اطلاعات مورد نیاز را

منابع

- Baecke, P.; Bocca, L., (2017). The value of vehicle telematics data in insurance risk selection processes. *Decis. Support Syst.*, 98(4): 69-79 (10 Pages).
- Balaji, S.; Srivatsa, S.K., (2012). Naïve bayes classification approach for mining life insurance databases for effective prediction of customer preferences over life insurance products. *Int. J. Comput. Appl.*, 51(3): 22-26 (24 Pages).
- Bhowmik, R., (2011). Detecting auto insurance fraud by data mining techniques. *J. Emerging Trends Comput. Inf. Sci.*, 2(4): 156-162 (6 Pages).
- Gharakhani, M.; Abolghasemi, M., (2011). Applications of data mining in the insurance industry. *Insur. World News*, 14(158): 5-21 (16 Pages). [In Persian]
- Gharanejad, S., (2010). The need to retain insurance customers using data mining tools. *Insur. World News*, 13(150-151): 15-23 (8 Pages). [In Persian]
- Ghazanfari, M.; Alizadeh, S.; Timurpour, B., (2008). Data mining and knowledge discovery. Tehran: University of Science and Technology Press, First edition. [In Persian]
- Haji Heydari, N.; Khale, S.; Farahi, A., (2011). Classifying the risk level of car body insurance policyholders using data mining algorithms (Case study: An insurance company). *Insur. Res. J. (Insur. Ind.)*, 26(4): 107-129 (22 Pages). [In Persian]
- Morik, K.; Kopcke, H., (2004). Analysing customer churn in insurance data—a case study. In *Knowledge Discovery in Databases: PKDD 2004: 8th European Conference on Principles and Practice of Knowledge Discovery in Databases*, Pisa, Italy, September 20-24, 2004. Proceedings 8 (pp. 325-336). Springer Berlin Heidelberg.
- Motarjem, K.; Niakan, L., (2020). Measuring and evaluating the satisfaction of life insurance customers. *Iran. J. Insur. Res.*, 10(1): 87-119 (32 Pages). [In Persian]
- Oshini, G.T.L.; Caldera H.A., (2013). Mining life insurance data for customer attrition analysis. *J. Ind. Intell. Inf.*, 1(1): 52-58 (6 Pages).
- Raedel, M.; Hartmann, A.; Priess, H.W.; Bohm, S.; Samietz, S.; Konstantinidis, I.; Walter, M.H., (2017). Re-interventions after restoring teeth-mining an insurance database. *J. Dent.*, 57(7): 14-19 (5 Pages).
- Ranjan, R., (2011). Self insurance and insurance demand under self-deception. *Asia-Pac. J. Risk Insur.*, 5(2): 1-27 (26 Pages).
- Saidur Rahman, M.; Arefin, K.Z.; Masud, S.; Sultana, S.; Rahman, R.M., (2017). Analyzing life insurance data with different classification techniques for customers' behavior analysis. *Adv. Top. Intell. Inf. and Database Syst.* 710 (24): 15-25 (10 Pages).
- Shahrabi, J., (2012). Data analysis. Tehran: Publications of Amir Kabir University of Technology (Tehran Polytechnic), First edition. [In Persian]
- Sundarkumar, G.G.; Ravi, V., (2015). A novel hybrid under sampling method for mining unbalanced datasets in banking and insurance. *Eng. Appl. Arti. Intell.*, 37(9): 368-377 (9 Pages).
- Thakur, S.S.; Sing, J.K., (2013). Mining customer's data for vehicle insurance prediction system using k-means clustering - An Application. *Int. J. Comput. Appl. Eng. Sci.*, 3(4): 148-153 (5 Pages).
- Umamaheswari, K.; Janakiraman, S., (2014). Role of data mining in insurance industry. *Int. J. Adv. Comput. Technol.*, 3(6): 961-966 (5 Pages).
- Wanke, P.; Barros, C.P., (2015). Efficiency drivers in Brazilian insurance: A two-stage DEA meta frontier-data mining approach. *Econ. Modell.*, 53(C): 8-22 (14 Pages).

AUTHOR(S) BIOSKETCHES	معرفی نویسندگان
<p>مریم نژاد افراسیابی، دانشجوی دکتری گروه مهندسی صنایع، دانشکده مهندسی صنایع و سیستم‌های مدیریت، دانشگاه صنعتی امیرکبیر، تهران، ایران</p> <ul style="list-style-type: none"> ▪ Email: mafrasiabi@aut.ac.ir ▪ ORCID: 0000-0001-6477-2390 ▪ Homepage: https://ie.aut.ac.ir/index.php?sid=9&slc_lang=fa 	
<p>اکبر اصفهانی پور، دانشیار گروه مهندسی مالی، دانشکده مهندسی صنایع و سیستم‌های مدیریت، دانشگاه صنعتی امیرکبیر، تهران، ایران</p> <ul style="list-style-type: none"> ▪ Email: esfahaa@aut.ac.ir ▪ ORCID: 0000-0003-3222-5186 ▪ Homepage: https://aut.ac.ir/cv/2135/%D8%A7%DA%A9%D8%A8%D8%B1-%D8%A7%D8%B5%D9%81%D9%87%D8%A7%D9%86%DB%8C-%D9%BE%D9%88%D8%B1?slc_lang=fa&cv=2135&mod=scv 	
<p>علی محمد کیمیاگری، دانشیار گروه مهندسی مالی، دانشکده مهندسی صنایع و سیستم‌های مدیریت، دانشگاه صنعتی امیرکبیر، تهران، ایران</p> <ul style="list-style-type: none"> ▪ Email: kimiagar@aut.ac.ir ▪ ORCID: 0000-0001-7216-0685 ▪ Homepage: https://aut.ac.ir/cv/2267/%D8%B9%D9%84%DB%8C%20%D9%85%D8%AD%D9%85%D8%AF%20%DA%A9%DB%8C%D9%85%DB%8C%D8%A7%DA%AF%D8%B1%DB%8C 	
<p>HOW TO CITE THIS ARTICLE</p> <p><i>Nezhad Afrasiabi, M.; Esfahanipour, A.; Kimiagari, A.M., (2021). Discovering profitable customers by data mining approach. Iran. J. Insur. Res, 10(3): 185-196.</i></p> <p>DOI: 10.22056/ijir.2021.03.03</p> <p>URL: https://ijir.irc.ac.ir/article_134716.html?lang=en</p>	