



Predicting Bank Customer Churn Using Machine Learning

Mohammad Hossein Mahmoudzadeh 

MSc., Department of Computer Engineering, Faculty of Mathematics and Computer Science, Amirkabir University of Technology (Tehran Polytechnic), Tehran, Iran. E-mail: mhmahmoudzadeh@aut.ac.ir

Mohammad Hassan Shirali Shahreza * 

*Corresponding Author, Assistant Prof., Department of Computer Education, Faculty of Mathematics and Computer Science, Amirkabir University of Technology (Tehran Polytechnic), Tehran, Iran. E-mail: hshirali@aut.ac.ir

Abstract

Objective

In competitive markets, companies focus on establishing long-term relationships with their customers and strengthening their loyalty. Due to the high costs associated with acquiring new customers, businesses tend to focus on retaining existing ones. Predicting which customers are likely to churn in the future plays a crucial role in shaping effective customer retention strategies. To predict customer churn and identify its drivers, companies use customer information and historical data recorded from them. This paper investigates customer churn prediction in the banking industry using real customer data from one of the largest banks in Iran.

Methods

To analyze and predict customer behavior, two equal and consecutive periods were considered. Customer behavior in the first period was used to predict a target variable in the second period. A significant drop in the average effective balance during the second period, compared to the first, was defined as the indicator of customer churn. By processing a large volume of banking transactions in the first period and aggregating them at different levels, various behavioral features for customers. We selected. One fixed validation set and three

Citation: Mahmoudzadeh, Mohammad Hossein & Shirali Shahreza, Mohammad Hassan (2025). Predicting Bank Customer Churn Using Machine Learning. *Financial Research Journal*, 27(2), 218-245. <https://doi.org/10.22059/FRJ.2024.357770.1007453> (in Persian)



training sets of different sizes were selected. To address the issue of dataset imbalance, class weights were determined based on the ratio of class sizes, ensuring that the minority class received greater weight during the training process. To predict customer churn, widely used machine learning algorithms—including Naive Bayes, k-Nearest Neighbors, Support Vector Machine, Logistic Regression, and Decision Tree—were applied, along with ensemble learning methods such as Random Forest, Adaptive Boosting, and Gradient Boosting. Subsequently, deep learning methods were applied, and a model incorporating modern modules such as residual connections and layer normalization—similar to state-of-the-art architectures—was proposed. Exhaustive experiments were conducted to evaluate the performance of the aforementioned methods.

Results

The results showed that ensemble learning algorithms and the proposed deep learning models outperformed the baseline models. Additionally, increasing the size of the training set contributed to improved model performance. Among the traditional machine learning classification algorithms, the decision tree trained on two training sets obtained the highest AUC ROC on the validation set with 0.8531 and 0.8597. The gradient boosting model obtained the overall highest AUC ROC on the validation set with 0.8984 and 0.9010. Deep learning-based single models achieved AUC-ROC values of 0.8825, 0.8909, and 0.8958, outperforming all traditional methods and two ensemble learning approaches while performing competitively with the gradient boosting algorithm.

Conclusion

Extracting behavioral features from customers' banking transactions and applying ensemble methods, along with the proposed deep learning-based models, proves effective in predicting banking customer churn, particularly in cases of a significant decrease in the average effective balance.

Keywords: Banking, Customer churn prediction, Deep learning, Machine learning.

پیش‌بینی ریزش مشتریان بانک با استفاده از یادگیری ماشین

محمدحسین محمودزاده

کارشناسی ارشد، گروه مهندسی کامپیوتر، دانشکده ریاضی و علوم کامپیوتر، دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران)، تهران، ایران.
رایانame: mhmahmoudzadeh@aut.ac.ir

محمدحسن شیرعلی شهرضا *

* نویسنده مسئول، استادیار، گروه آموزشی علوم کامپیوتر، دانشکده ریاضی و علوم کامپیوتر، دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران)،
تهران، ایران. رایانame: hshirali@aut.ac.ir

چکیده

هدف: در بازارهای رقابتی، شرکت‌ها روی برقراری روابط بلندمدت با مشتریان و تقویت وفاداری تمرکز دارند. به علت هزینه‌های سنگین جذب مشتری جدید، کسب‌وکارها روی نگهداری مشتریان موجود تمرکز می‌کنند. پیش‌بینی مشتریانی که احتمال روی‌گردانی آن‌ها در آینده وجود دارد، بخش مهمی از راهبرد حفظ مشتری است. در این مقاله پیش‌بینی ریزش مشتری در صنعت بانکداری، روی دادگان واقعی مشتریان یکی از بانک‌های بزرگ ایران انجام شده است.

روش: در صنعت بانکداری، کاهش شدید میانگین مانده مؤثر یک مشتری در یک بازه زمانی نسبت به بازه زمانی قبلی، به عنوان ریزش مشتری در نظر گرفته می‌شود. در این مقاله، ابتدا با پردازش حجم زیادی از تراکنش‌های بانکی در یک بازه زمانی مشخص، ویژگی‌های رفتاری متفاوت در سطوح مختلف برای مشتریان به دست آمد؛ سپس برای پیش‌بینی ریزش، از الگوریتم‌های پُراستفاده در یادگیری ماشین و روش‌های یادگیری جمعی استفاده شد. در ادامه، با استفاده از روش‌های یادگیری عمیق و واحدهای نوین آن، معماری مدل مدنظر ارائه شد. در نهایت، با انجام آزمایش‌های جامع، عملکرد روش‌های نام برده بررسی شد.

یافته‌ها: این پژوهش در یکی از بانک‌های بزرگ ایران اجرا شد و آزمایش‌ها روی دادگان واقعی مشتریان بانک صورت پذیرفت. در این آزمایش‌ها، از اطلاعات جمعیت‌شناختی و رفتار گذشته مشتریان بهره گرفته شد؛ اما از اطلاعات شخصی افراد استفاده نشد تا حریم خصوصی مشتریان حفظ شود. در آزمایش‌های صورت‌گرفته، پیش‌بینی ریزش مشتری روی بازه زمانی یک ماهه انجام گرفت. بدین ترتیب دو بازه یک‌ماهه متواالی مدنظر قرار گرفت و ویژگی‌های رفتاری مشتریان، از بازه زمانی اول استخراج شد. متغیر هدف نیز از مقایسه میانگین مانده مؤثر در بازه‌های زمانی اول و دوم به دست آمد. در صورتی که میانگین مانده مؤثر یک مشتری، در بازه دوم نسبت به بازه اول با بیش از ۷۰ درصد کاهش همراه بود، به عنوان ریزش در نظر گرفته شد. در نتایج به دست آمده الگوریتم‌های یادگیری جمعی و

استناد: محمودزاده، محمدحسین و شیرعلی شهرضا، محمدحسن (۱۴۰۴). پیش‌بینی ریزش مشتریان بانک با استفاده از یادگیری ماشین. تحقیقات مالی، ۲(۲)، ۲۱۸-۲۴۵.

تاریخ دریافت: ۱۴۰۲-۰۱-۲۷

تاریخ پذیرش: ۱۴۰۳-۰۷-۱۵

تاریخ ویرایش: ۱۴۰۳-۰۷-۲۵

نasher: دانشکده مدیریت دانشگاه تهران

تاریخ انتشار: ۱۴۰۴-۰۲-۳۱

نوع مقاله: علمی پژوهشی

© نویسنده‌گان

doi: <https://doi.org/10.22059/FRJ.2024.357770.1007453>

همچنین مدل‌های عمیق ارائه شده، عملکرد بهتری را نسبت به مدل‌های مبتنا نشان دادند. افزایش اندازه مجموعه آموزش در عملکرد بهتر مدل‌ها مؤثر بود. مدل تقویت گرادیان با ۸۹۸۴٪ بیشترین مساحت زیر منحنی مشخصه عامل گیرنده نسبت به مجموعه اعتبارسنجی را به دست آورد.

نتیجه‌گیری: استخراج ویژگی‌های رفتاری از تراکنش‌های بانکی مشتریان و استفاده از روش‌های یادگیری جمعی و همچنین مدل‌های ارائه شده مبتنی بر یادگیری عمیق، در پیش‌بینی ریزش مشتری مؤثرند. پس از تحلیل رفتار و شناسایی مشتریان در شرف روی گردانی، پیشنهادهایی برای جلوگیری از ریزش و حفظ مشتری ارائه شد. برای نمونه، تفکیک مشتریان بر اساس سن، شغل، تحصیلات و غیره به منظور ارائه خدمات و تولید محصولات بانکی بر این مبنای ایجاد تنوع در خدمات موجود، ارائه خدمات مورد نیاز مشتریان از طریق بسترها مجازی و در صورت نیاز در محل فعالیت و زندگی مشتریان، تسهیل در ارائه خدمات مطلوب به مشتریان، افزایش اعتماد مشتریان از طریق ارائه کاربردی امن و همچنین، حفاظت از اطلاعات مشتریان، به حفظ مشتری کمک می‌کند.

کلیدواژه‌ها: پیش‌بینی ریزش مشتری، بانکداری، یادگیری عمیق، یادگیری ماشین.

مقدمه

برای هر کسبوکاری مشتریان اساس درآمد و موفقیت‌اند و یکی از دارایی‌های مهم در نظر گرفته می‌شوند. در بازارهای رقابتی که خدمات عرضه شده توسط ارائه‌دهندگان مختلف، مشابه یکدیگرند، مشتریان قدرت انتخاب دارند و به راحتی می‌توانند خدمت‌رسان خود را تغییر دهنند. در چنین شرایطی، سازمان‌ها روی برقراری روابط بلندمدت با مشتریان خود تمرکز دارند و در طول زمان رفتار مشتری را رصد می‌کنند (Amin و Hemkaran, ۲۰۱۷). مدیریت ارتباط با مشتری، راهبردی جامع برای ایجاد و تقویت روابط وفادارانه و طولانی‌مدت با مشتریان است. مدیریت ارتباط با مشتری در کسبوکارهای مختلف مانند مخابرات، بانک، بیمه، بازار خرده‌فروشی و غیره به رسمیت شناخته شده و به طور گسترده استفاده می‌شود (Vafeiadis, Diamantaras, Sarigiannidis & Chatzisavvas, ۲۰۱۵). امروزه رویکرد مشتری محور برای پیش‌بینی رفتار مشتریان بر اساس داده‌های تاریخی ذخیره شده در بانک‌های اطلاعاتی و سامانه‌های مدیریت ارتباط با مشتری بسیار رایج است (Amin و Hemkaran, ۲۰۱۷).

نگهداری مشتری و پیش‌بینی ریزش

یکی از اهداف اصلی در مدیریت ارتباط با مشتری، حفظ مشتریان موجود است. طبق بررسی‌های صورت گرفته، جذب مشتریان جدید بسیار پُرهزینه‌تر از حفظ مشتریان موجود است. در برخی موارد هزینه جذب مشتری تا ۲۰ برابر بیشتر از حفظ مشتری برآورد شده است (Vafeiadis و Hemkaran, ۲۰۱۵). این واقعیت اهمیت نگهداری مشتریان موجود را آشکار ساخته و تمرکز کسبوکارها نیز به سمت آن تغییر کرده است. همچنین تقویت وفاداری و حفظ مشتری به بهبود فروش و کاهش هزینه‌های بازاریابی منجر می‌شود (Amin و Hemkaran, ۲۰۱۷).

بخشی از راهبرد حفظ مشتری که از اهمیت زیادی برخوردار است، پیش‌بینی ریزش مشتری است. فرایند جستجو و شناسایی مشتریانی که تمایل زیادی به ترک شرکت نشان می‌دهند، به عنوان پیش‌بینی ریزش مشتری شناخته می‌شود. برای مواجهه با مسئله ریزش، می‌توان از داده‌های نگهداری شده از مشتریان استفاده و آن را به دانش معنادار تبدیل کرد. در پیش‌بینی ریزش مشتری، یک مدل امتیازدهنده بر اساس رفتار گذشته مشتری، برآورد احتمال روی‌گردانی داوطلبانه در آینده را فراهم می‌آورد (دی کاینی، کوسمن و دی باک, ۲۰۱۸). هدف مدل‌های پیش‌بینی ریزش، تشخیص نشانه‌های اولیه روی‌گردانی و شناسایی مشتریانی است که احتمال ترک آن‌ها افزایش یافته است. بدین ترتیب می‌توان قبل از روی‌گردانی، مشتریان هدف را شناسایی کرد و به منظور ارائه پیشنهادها و برگزاری پویش‌های بازاریابی، از آن‌ها بهره برد. ریزش مشتری می‌تواند نتیجه سطح پایین رضایت، راهبردهای رقابتی تهاجمی، محصولات جدید، مقررات یا موارد دیگر باشد (Vafeiadis و Hemkaran, ۲۰۱۵). علاوه‌بر پیش‌بینی مشتریان محتمل، محققان در پی این هستند تا به دلایل ریزش مشتری و محرك‌های مهم آن، مانند رضایت مشتری پی ببرند (دی کاینی و Hemkaran, ۲۰۱۸). از این‌رو، توانایی عملکردی مدل‌های پیش‌بینی ریزش مشتری حائز اهمیت است. استفاده از این مدل‌های پیش‌بینی کننده، باید به

1. Amin et al.

2. Vafeiadis, Diamantaras, Sarigiannidis & Chatzisavvas

3. De Caigny, Coussement & De Bock

بینش‌های قابل اجرا و درک بهتر محرك‌های ریزش مشتری، به منظور اتخاذ تصمیم‌های آگانه مدیریتی منجر شود (دی کاینی و همکاران، ۲۰۱۸). این حقایق در نهایت باعث می‌شود که فعالیت پیش‌بینی ریزش مشتری، به عنوان بخشی ضروری از تصمیم‌گیری راهبردی و فرایند برنامه‌ریزی شرکت‌ها و سازمان‌ها باشد (امین و همکاران، ۲۰۱۷). بدین ترتیب اهمیت اطلاعات تاریخی مشتریان که می‌تواند برای ساختن مدل‌های پیش‌بینی کننده استفاده شود، روشن است و به عنوان دارای مهم برای مقابله با ریزش مشتری شمرده می‌شود. از این رو شرکت‌های زیادی برای سرمایه‌گذاری در مدیریت ارتباط با مشتری و حفظ اطلاعات مشتریان اقدام کرده‌اند.

شرح مسئله و چالش‌های آن

پیش‌بینی ریزش مشتری، مسئله‌ای برای دسته‌بندی مشتریان به گروه‌های ریزش و عدم ریزش است. بدین منظور با در اختیار داشتن اطلاعات مشتریان و داده‌های تاریخی ثبت‌شده برای آن‌ها، باید به مهندسی ویژگی پرداخت. ویژگی‌های استخراج شده برای مشتری، به عنوان ورودی به مدل پیش‌بینی کننده داده می‌شود. مدل تخمین خود را از دسته‌های مشتری به آن مربوط است، به عنوان خروجی ارائه می‌دهد. خروجی مدل می‌تواند به صورت برآورده از احتمال تعلق مشتری به دستهٔ ریزش باشد.

از پیش‌بینی این مدل‌ها، برای انتخاب مشتریان هدف به منظور انجام بررسی‌های بیشتر و تلاش برای جلوگیری از ترک مشتری استفاده می‌شود. علاوه‌بر شناسایی مشتریان در شرف روی گردانی، یافتن علل و محرك‌های مؤثر بر ریزش مشتری به منظور اتخاذ تصمیم‌های مدیریتی نیز ارزش‌آفرین است. از این رو توانایی عملکردی این مدل‌ها در پیش‌بینی احتمال ریزش حائز اهمیت است. بدین منظور، استفاده از الگوریتم‌های دسته‌بندی کارا و وجود ویژگی‌های مرتبط و مفید برای مشتریان لازم است. این ویژگی‌ها باید از پردازش و تجمعیح حجم زیادی از دادگان ذخیره شده از فعالیت مشتریان در گذر زمان مهندسی و استخراج شوند.

از دیگر چالش‌هایی که در مسئلهٔ پیش‌بینی ریزش با آن مواجه می‌شویم، مشکل عدم تعادل مجموعه داده است. مشتریان مربوط به دستهٔ ریزش که مورد توجه قرار دارند، درصد کمتری از مجموعه داده را تشکیل می‌دهند. بدین ترتیب برای مقابله با عدم تعادل مجموعه داده نیز، باید راه کاری را اتخاذ کرد.

پیش‌بینی ریزش مشتری در زمینه‌های مختلف مخابرات، خدمات مالی، بیمه و غیره مطالعه و بررسی شده و مقاله‌های زیادی نیز با محوریت آن انتشار یافته است. در بیشتر مقاله‌ها با استفاده از الگوریتم‌های دسته‌بندی یادگیری ماشین سنتی، فرایند پیش‌بینی ریزش مشتری صورت پذیرفته است (احمد، جعفر و الجمیع^۱؛ کاروانا، یزید، سیالیم و مورسانتو^۲؛ لالوانی، میشرا، چادا و ستی^۳؛ ۲۰۲۲). همچنین روش‌هایی از یادگیری عمیق نیز برای پیش‌بینی ریزش مشتری، به طور محدود استفاده شده است (خان و همکاران^۴؛ ۲۰۱۹؛ اسپانودس و نگوین^۵؛ ۲۰۱۷).

1. Ahmad, Jafar & Aljoumaa
2. Karvana, Yazid, Syalim & Mursanto
3. Lalwani, Mishra, Chadha & Sethi
4. Khan et al.
5. Spanoudes & Nguyen

در این مقاله، مسئلهٔ پیش‌بینی ریزش مشتری در صنعت بانکداری مطالعه و بررسی می‌شود. این پژوهش در یکی از بزرگ‌ترین بانک‌های ایران انجام شده و از داده‌های واقعی مشتریان بانک بهره گرفته شده است. برای حفظ حریم خصوصی مشتریان، از ذکر نام بانک در این مقاله خودداری شده است. با فرض اینکه حجم زیاد داده، از تعامل مشتریان بانک در طول زمان در اختیار است، نحوه استفاده از دادگان به‌طور مؤثر و استخراج ویژگی‌های مفید برای پیش‌بینی ریزش مشتری، مورد توجه قرار گرفته است. با استفاده از معیار مهم میانگین مانده حساب مشتری و مقایسه آن در بازه‌های زمانی مختلف، متغیر هدف را مشخص می‌کنیم. با بررسی تراکنش‌های انجام‌شده مشتری در یک بازه زمانی و فراداده‌های مربوط به آن، ویژگی‌هایی از نحوه رفتار مالی مشتری به‌دست می‌آوریم.

همان‌طور که اشاره شد، در اغلب کارهای پیشین در زمینهٔ پیش‌بینی ریزش مشتری، از روش‌های یادگیری ماشین سنتی استفاده شده است. با توجه به اهمیت عملکرد مدل پیش‌بینی کننده مورد استفاده و ظرفیت یادگیری آن، این سؤال مطرح می‌شود که آیا استفاده از یادگیری عمیق نوین، می‌تواند به بهبود عملکرد کمک کند؟ در این راستا آزمایش‌های جامعی برای بررسی استفاده از روش‌های مختلف یادگیری ماشین و تأثیر اندازه مجموعه دادگان آموزش انجام می‌دهیم. بدین منظور از الگوریتم‌های پُرکاربرد و مرسوم دسته‌بندی در یادگیری ماشین بهره می‌بریم. از روش‌های یادگیری عمیق استفاده می‌کنیم و با به کارگیری برخی از واحدهای نوین آن، شبکه‌هایی را تشکیل می‌دهیم که معماری آن‌ها مشابه با آخرین دستاوردهای این حوزه است.

ساختار این مقاله در ادامه بدین شرح است: در ادامه، به پیشینهٔ تحقیق در حوزهٔ پیش‌بینی ریزش مشتری و روش‌های پُراستفاده پرداخته و چند نمونه از کارهای صورت‌گرفته در سال‌های اخیر مرور می‌شود. سپس مسئلهٔ پژوهش و روش تحقیق تشریح می‌شود. در بخش بعد، نتایج حاصل از آزمایش‌های انجام‌شده ارائه می‌شود. در پایان مقاله نیز نتیجه‌گیری و جمع‌بندی نهایی ارائه خواهد شد.

پیشینهٔ تجربی پژوهش

تحلیل رفتار مشتریان بانک و دسته‌بندی آن‌ها از مباحثی است که در سال‌های اخیر، در کانون توجه پژوهشگران حوزه تحقیقات مالی بوده است. به‌طور نمونه، در مقالهٔ رحیمی، روستا و آسایش (۱۴۰۳) به ارائه الگوی ارتقای توان رقابت‌پذیری خدمات ارزی مشتریان در صنعت بانکداری پرداخته شده است. برای تعیین ابعاد، مؤلفه‌ها و شاخص‌ها، از آزمون دلفی و برای تعیین وضع موجود و همچنین، تعیین عوامل اثرگذار، از آزمون تحلیل معادلات ساختاری و تحلیل عاملی اکتشافی استفاده شده است. بر اساس نتایج آزمون معادلات ساختاری و آماره‌های به‌دست آمده، دریافت‌های اند که ضرایب تمامی عوامل در نظر گرفته شده، بر مدل نهایی تأثیر معناداری می‌گذارند و اثرگذاری تمامی عوامل، بر الگوی ارتقای توان رقابت‌پذیری خدمات ارزی مشتریان در صنعت بانکداری تأیید شده است. در این مقاله اشاره شده است که بانک‌ها با ایستی با ارزیابی وضعیت در حوزه بانکی، به ایجاد مزیت رقابتی برای جذب، حفظ و نگهداری مشتریان اقدام کنند. در مقالهٔ احمدی کوشان، احمدی، رنجبر و کردلوئی (۱۴۰۳) شناسایی شاخص‌های اعتبارسنجی و رتبه‌بندی مشتریان

در تسهیلات خُرد در بانک خاورمیانه صورت گرفته است. در این پژوهش روش‌های آماری در دو بخش آمار توصیفی و استنباطی انجام گرفته است. در بخش آمار توصیفی، برخی عوامل شخصیتی همچون سن، جنسیت، تحصیلات، کسب‌وکار... تجزیه و تحلیل و از طریق جداول و نمودارها بررسی شده است. الگوریتم‌هایی بیز ساده اجرا و برای رده‌بندی معیارها و ایجاد الگو استفاده شده است. بهمنظور ارزیابی نهایی مدل اعتبارسنجی و رتبه‌بندی مشتریان تسهیلات خُرد بدون پشتوانه، از آزمون تی استفاده شده است. یافته‌های این مقاله نشان می‌دهد که بانک باید بر اساس شاخص‌هایی در خصوص پرداخت یا عدم پرداخت وام تصمیم گیری کند و درباره اشخاص متقاضی تسهیلات، شناخت کامل داشته باشد. بدین منظور باید داده‌های لازم را از آن‌ها جمع‌آوری کند و از حسن اعتبار و شهرت آن‌ها اطمینان یابد.

در مقاله باجلان، فلاچپور و رئیسی (۱۴۰۳) مدلی برای اندازه‌گیری ریسک و بهینه‌سازی پرتفوی اعتباری بانک‌ها ارائه شده است. ابتدا به بررسی ریسک اعتباری پرتفوی تسهیلات بانکی با استفاده از رویکرد اکچوئری پرداخته و سپس با استفاده از شبکه عصبی پرسپترون و با توجه به محدودیت‌های بانک در ارائه تسهیلات، ترکیب بهینه پرتفوی اعتباری تعیین شده است. نتایج بهدست آمده در این مقاله نشان می‌دهد که بازده پرتفوی بهینه از پرتفوی فعلی بانک بالاتر است. براساس یافته‌ها و تأیید فرضیه‌های پژوهشی نتیجه گرفته شده است که استفاده از مدل اکچوئری برای تعیین ریسک اعتباری و سپس بهینه‌سازی با استفاده از شبکه عصبی مصنوعی، به بهبود فرایند بهینه‌سازی پرتفوی اعتباری بانک‌ها منجر می‌شود.

در مقاله احمدی سرتختی، هژبر کیانی، حسینی و معمان‌زاد (۱۴۰۲) انتخاب بهترین روش برای اعتبارسنجی مشتریان ضمانت‌نامه‌های اعتباری و تفکیک اشخاص خوش حساب از بدحساب، برای کاهش نکول اعتبارات اعطا شده صندوق ضمانت صادرات ایران صورت گرفته است. به‌کمک مدل شبکه عصبی مصنوعی، مدلی برای ارزیابی ریسک متقاضیان تسهیلات و ضمانت‌نامه‌ها از این صندوق طراحی شده که بیشترین قدرت پیش‌بینی احتمال نکول تسهیلات اعطایی را داشته باشد. بر اساس نتایج حاصل از این پژوهش، نتیجه گیری شده است که می‌توان ضمن تفکیک مشتریان خوش حساب از بدحساب و بر اساس رتبه اعتباری مشتریان، میزان وثایق اخذشده از مشتریان را متناسب با وضعیت اعتباری گروه‌های اعتباری تنظیم کرد.

با بررسی این مقاله‌ها می‌توان به اهمیت دسته‌بندی مشتریان در نظام بانکی پی برد. همچنین در تعدادی از کارها، به اعتبارسنجی مشتریان بهمنظور ارائه تسهیلات پرداخته شده است. اما متفاوت با پژوهش‌های ذکر شده، این مقاله به پیش‌بینی ریزش مشتریان بانک با تحلیل رفتار مشتری بر اساس تراکنش‌های بانکی انجام شده در بازه زمانی مشخص می‌پردازد.

در گذشته مسئله پیش‌بینی ریزش مشتری در حوزه‌های کسب‌وکاری مختلف مانند مخابرات، بانکداری و خدمات مالی، بیمه و غیره مورد مطالعه قرار گرفته است. در بسیاری از مطالعات ابتدا روی دادگان ذخیره شده از سابقه فعالیت مشتریان مهندسی ویژگی صورت گرفته است. با تجمعیg دادگان و استخراج ویژگی‌های رفتاری و همچنین اضافه کردن اطلاعات جمعیت‌شناختی، مجموعه داده‌های ساختاری‌افته بهدست آمده است. در این موارد مطالعات روی داده‌های واقعی

شرکت‌ها یا مؤسسه‌هایی که پژوهش در آن‌ها صورت گرفته، انجام شده است (احمد و همکاران، ۲۰۱۹؛ کاروانا و همکاران، ۲۰۱۹). در تعدادی از کارها نیز استفاده از محدود مجموعه داده‌های عمومی مرتبط با پیش‌بینی ریزش مشتری در دستور کار قرار گرفته است (هالیاس و همکاران^۱، ۲۰۱۹؛ پوستوخین، نگوین، الهوسنی و شانکار^۲، ۲۰۲۱). در اغلب پژوهش‌ها عدم تعادل مجموعه داده، از چالش‌های موجود بوده که برای مواجهه با آن، از روش‌های کمونونه‌برداری و بیش‌نمونه‌برداری استفاده شده است (رحمان و کومار^۳، ۲۰۲۰). برای پیش‌بینی ریزش مشتری، از مدل‌هایی مانند درخت تصمیم^۴، بیز ساده^۵، k نزدیک‌ترین همسایه^۶، رگرسیون لجستیک^۷، ماشین بردار پشتیبان^۸ استفاده بیشتری شده است. همچنین مدل‌های یادگیری جمعی مانند جنگل تصادفی^۹، تقویت تطبیقی^{۱۰}، تقویت گرادیان^{۱۱} و غیره نیز استفاده شده‌اند. بدین ترتیب، در اغلب مقاله‌های منتشرشده، استفاده از الگوریتم‌های دسته‌بندی یادگیری ماشین سنتی روی دادگان ساختاریافته بررسی شده و الگوریتم‌های یادگیری جمعی عملکرد خوبی داشته است.

استفاده از روش‌های یادگیری ماشین

در ادامه به چند نمونه از مطالعات انجام شده در سال‌های اخیر می‌پردازیم. در مقاله احمد و همکاران (۲۰۱۹) پیش‌بینی ریزش مشتری در صنعت مخابرات بررسی شده است. بدین منظور از روش‌های یادگیری ماشین استفاده شده و برای سنجش عملکرد مدل نیز معیار مساحت زیر منحنی انتخاب شده است. مجموعه دادگان شامل، اطلاعات مشتریان در طول ۹ ماه فعالیت بوده است. علاوه بر ویژگی‌های آماری، نوع دیگری از ویژگی‌ها بر اساس فعالیت اجتماعی مشتریان از طریق پیامک و تماس محاسبه شده است. برای تأیید اعتبار و بهینه‌سازی ابرپارامترها، اعتبارسنجی متقابل ده برابری اجرا شده است. برای رفع نامتعادل بودن داده‌ها از زیرنمونه‌برداری یا اجرای الگوریتم‌های مبتنی بر درخت که نامتعادل بودن دادگان تأثیری روی آن‌ها ندارد، استفاده شده است. چهار الگوریتم مبتنی بر درخت شامل درخت تصمیم، جنگل تصادفی، ماشین تقویت گرادیان و تقویت گرادیان شدید آزمایش شده‌اند. استفاده از ویژگی‌های شبکه اجتماعی، بیشترین تأثیر را روی بهبود عملکرد مدل داشته و معیار مساحت زیر منحنی را از $\frac{۹۳}{۳}$ به ۸۴ افزایش داده است. بهترین نتایج توسط الگوریتم تقویت گرادیان شدید به دست آمد و پس از آن نیز به ترتیب الگوریتم‌های تقویت گرادیان، جنگل تصادفی و درخت تصمیم قرار گرفتند. در مقاله لالوانی و همکاران (۲۰۲۲) نیز پیش‌بینی ریزش مشتری در صنعت مخابرات انجام شده است. سپس با استفاده از الگوریتم جست‌وجوی گرانشی، انتخاب ویژگی اجرا شده است. مدل‌هایی مانند رگرسیون

1. Halibas et al.

2. Pustokhina, Pustokhin, Nguyen, Elhoseny & Shankar

3. Rahman & Kumar

4. Decision Tree

5. Naive Bayes

6. k-Nearest Neighbors

7. Logistic Regression

8. Support Vector Machine

9. Random Forest

10. Adaptive Boosting (AdaBoost)

11. Gradient Boosting

لجستیک، بیز ساده، ماشین بردار پشتیبان، درخت تصمیم و غیره، روی دادگان آموزش اعمال شده است. همچنین تأثیر شیوه‌های تقویتی و جمعی مانند تقویت تطبیقی، تقویت گرادیان شدید، جنگل تصادفی و غیره نیز بر دقت مدل بررسی شده است. برای تنظیم ابرپارامترها و جلوگیری از بیش برازش مدل‌ها از اعتبارسنجی k برابری روی مجموعه آموزش استفاده شده است. عملکرد مدل‌های مختلف روی مجموعه آزمایش با استفاده از ماتریس درهم‌ریختگی و مساحت زیر منحنی ارزیابی شده‌اند. نتایج به دست آمده نشان می‌دهد که دسته‌بندهای تقویت تطبیقی و تقویت گرادیان شدید، به ترتیب با درصد ۸۱/۷۱ و ۸۰/۸ درصد بیشترین دقت‌ها را به دست آورده‌اند. همچنین بیشترین مساحت زیر منحنی نیز ۸۴ درصد توسط دسته‌بندهای تقویت تطبیقی و تقویت گرادیان شدید به دست آمد.

در مقالهٔ پوستوخینا و همکاران (۲۰۲۱) مدلی برای پیش‌بینی ریزش مشتری در بخش مخابرات ایجاد شده است. بدین منظور شیوه‌ای برای بیش نمونه برداری ساختگی از دستهٔ اقلیت همراه با مدل یادگیری ماشین شدید با وزن‌های بهینه ارائه شده است. برای تعیین نرخ بهینه نمونه برداری و تنظیم پارامترهای مدل از الگوریتم بهینه‌سازی باران با چند هدف استفاده شده است. برای مدیریت عدم تعادل مجموعه داده، بیش نمونه برداری به کار گرفته شده و در نهایت، مدل دسته‌بندی روی داده‌های آماده‌سازی شده اعمال شده است. آزمایش‌های گسترده‌ای برای مدل ارائه شده روی چند مجموعه داده مخابراتی برای پیش‌بینی ریزش انجام شده است. نتایج نشان می‌دهد که مدل ارائه شده در سه مجموعه داده نسبت به مدل‌های دیگر برتری دارد.

در مقالهٔ ویجايا و سیواسانکار^۱ (۲۰۱۹) چند الگوریتم یادگیری ماشین برای پیش‌بینی ریزش مشتری در صنعت مخابرات بررسی شده است. علاوه بر این استفاده از الگوریتم‌های فرآبتكاری برای پیش‌بینی ریزش نیز به بحث گذاشته شده است. بدین منظور، تحلیلی روی استفاده از الگوریتم بهینه‌سازی ازدحام ذرات، به عنوان یک پیش‌بینی کننده ریزش انجام شده است. دسته‌بندهای پیشنهادی برای بررسی سطوح پیش‌بینی و جنبه‌های عملکردی با مدل‌هایی مانند درخت تصمیم، بیز ساده، k نزدیک ترین همسایه، ماشین بردار پشتیبان، جنگل تصادفی و سه مدل ترکیبی مقایسه شده‌اند. برای بررسی عملکرد مدل‌ها از معیارهایی مانند دقت، درستی، امتیاز F1، نمودار مشخصه عامل گیرنده و غیره استفاده شده است. آزمایش‌های صورت گرفته نشان می‌دهد که بهینه‌سازی ازدحام ذرات و انواع دیگر آن روی مجموعه داده‌های بزرگ و نامتعادل عملکرد خوبی داشته‌اند. همچنین مشاهده شد که استفاده از بهینه‌سازی ازدحام ذرات برای انتخاب ویژگی و بازپخت شبیه‌سازی شده به عنوان الگوریتم جستجوی محلی در شرایط نامتعادل به خوبی عمل می‌کنند.

در مقالهٔ کاروانا و همکاران (۲۰۱۹) پیش‌بینی ریزش مشتری در صنعت بانکداری مورد بررسی قرار گرفته است. روش‌های دسته‌بندی درخت تصمیم، شبکهٔ عصبی، ماشین بردار پشتیبان، بیز ساده و رگرسیون لجستیک روی مجموعه داده‌ای در یک بانک خصوصی در اندونزی آزمایش شده است. داده‌های به دست آمده شامل ویژگی‌های جمعیت‌شناسنامه مشتریان از قبیل سن، حرفه، محصولات مورد استفاده، ویژگی‌هایی مربوط به تراکنش‌های ورودی یا خروجی، تراکنش‌های داخلی یا خارجی و همچنین ویژگی‌هایی مانند میانگین موجودی در بازه‌های زمانی مختلف، تعداد

حساب‌های تحت مالکیت و غیره بوده است. آزمایش‌ها با نسبت‌های نمونه‌برداری مختلف از دسته‌های ریزش و عدم ریزش انجام شده است. نتایج به دست آمده نشان می‌دهند که تعداد نمونه‌های مورد استفاده برای آموزش تأثیر زیادی بر نتایج مدل دارد. همچنین نسبت اندازه دسته‌ها تأثیر زیادی بر نتایج معیار یادآوری دارد و با نسبت برابر دسته‌ها به معیار یادآوری بیشتری دست یافتند. ماشین بردار پشتیبان با نسبت نمونه‌گیری برابر از دسته‌ها بیشترین سود حاصل از پیش‌بینی ریزش مشتری را به دست آورد.

در مقالهٔ چلیک و عثمان اوغلو^۱ (۲۰۱۹) الگوریتم‌های یادگیری ماشین مانند شبکهٔ عصبی، درخت تصمیم، ماشین بردار پشتیبان، بیز ساده، k نزدیک‌ترین همسایه، تقویت گرادیان شدید و غیره برای تحلیل و پیش‌بینی ریزش مشتریان مقایسه شده‌اند.

در مقالهٔ هالیاس و همکاران (۲۰۱۹) تحلیل اکتشافی دادگان و مهندسی ویژگی برای پیش‌بینی ریزش مشتری در صنعت مخابرات مورد مطالعه قرار گرفته است. همچنین چندین روش دسته‌بندی شامل بیز ساده، رگرسیون لجستیک، درخت تصمیم، جنگل تصادفی، تقویت گرادیان و... روی مجموعه داده اعمال شده است. برای بررسی عملکرد مدل‌ها از معیارهای مختلفی مانند دقت، خطای دسته بندی، درستی، یادآوری، امتیاز F1 و مساحت زیر منحنی استفاده شده است. با توجه به نتایج به دست آمده در آزمایش، الگوریتم تقویت گرادیان بهترین عملکرد را داشته و تقریباً در تمام معیارها از سایر دسته‌بندها بهتر عمل کرده است. علاوه بر این، عملکرد همهٔ دسته‌بندها پس از انجام پیش‌نمونه‌برداری از دسته اقلیت، با بهبود همراه بوده است.

در مقالهٔ جین، یاداو و مانوف^۲ (۲۰۲۱) پیش‌بینی ریزش مشتری در سه حوزهٔ مختلف بانکداری، مخابرات و فناوری اطلاعات با استفاده از روش‌های یادگیری ماشین مورد مطالعه قرار گرفته است. برای پیش‌بینی ریزش مشتری در حوزه بانکی، از مجموعه داده‌ای عمومی، شامل ۱۰۰۰۰ نمونه و ۱۴ ویژگی استفاده شده است. ویژگی‌های مشتریان شامل امتیاز اعتباری، موقعیت جغرافیایی، جنسیت، سن، موجودی حساب، تعداد محصولات مورد استفاده، داشتن کارت اعتباری و غیره بوده است. برای پیش‌بینی ریزش از چهار مدل الگوریتمی رگرسیون لجستیک، جنگل تصادفی، ماشین بردار پشتیبان و تقویت گرادیان شدید استفاده شده و عملکرد الگوریتم‌های مختلف در حوزه‌های مختلف مقایسه شده است. همچنین با تحلیل اکتشافی داده‌ها، راهبردهایی برای حفظ مشتری به دست آمده است. طبق نتایج به دست آمده، جنگل تصادفی در بخش بانکی با دقت ۸۶/۳۱۲ درصد، رگرسیون لجستیک در بخش فناوری اطلاعات با دقت ۹۰/۱۳۶ درصد و تقویت گرادیان شدید در بخش مخابرات با دقت ۸۲/۹۴۲ درصد بهترین عملکردها را به دست آوردند.

در مقالهٔ رحمان و کومار (۲۰۲۰) پیش‌بینی ریزش مشتری با استفاده از یادگیری ماشین در صنعت بانکداری انجام شده است. در این مطالعه دسته بندهای k نزدیک‌ترین همسایه، ماشین بردار پشتیبان، درخت تصمیم و جنگل تصادفی تحت شرایط مختلف آزمایش شده‌اند. همچنین برخی از روش‌های انتخاب ویژگی برای یافتن ویژگی‌های مرتبط‌تر و بررسی عملکرد سیستم مورد استفاده قرار گرفته‌اند. این آزمایش‌ها روی مجموعه داده ذکر شده در مورد قبلی انجام شده

1. Çelik & Osmanoglu

2. Jain, Yadav & Manoov

است. در نتایج به دست آمده، مدل جنگل تصادفی پس از بیش نمونه برداری با دقت پیش‌بینی ۹۵/۷۴ درصد عملکرد بهتری را نسبت به سایر مدل‌های آزمایش شده به دست آورده است.

در مقاله سید، عبدالفتاح و خلیف^۱ (۲۰۱۸) نیز پیش‌بینی ریزش مشتری در صنعت بانکداری مطالعه شده است. به منظور تولید مدلی برای پیش‌بینی احتمال روی‌گردانی مشتریان بانک از الگوریتم درخت تصمیم استفاده شده است. در مقاله لی و همکاران^۲ (۲۰۲۱) پیش‌بینی ریزش مشتری در صنعت پخش رسانه‌های جمعی رادیو و تلویزیون انجام شده است. بدین منظور عوامل تأثیرگذار بر ریزش مشتری در صنعت شبکه‌های کابلی جمع‌آوری و بررسی شده است. برای کاهش ابعاد متغیرهای موجود از روش تحلیل مؤلفه اصلی و برای پیش‌بینی ریزش مشتری از الگوریتم دسته‌بندی رگرسیون لجستیک استفاده شده است.

جدول ۱ حوزه کسب‌وکاری و مدل‌های استفاده شده در چند مقاله بیان شده آورده شده است.

جدول ۱. حوزه کسب‌وکار و مدل‌های استفاده شده برای پیش‌بینی ریزش مشتری در تعدادی از مقاله‌های سال‌های اخیر

مقاله	حوزه	مدل‌ها
(احمد و همکاران، ۲۰۱۹)	مخابرات	تقویت گرادیان، تقویت گرادیان شدید، جنگل تصادفی، درخت تصمیم
(لالوani و همکاران، ۲۰۲۲)	مخابرات	بیز ساده، تقویت طبیقی، تقویت گرادیان شدید، جنگل تصادفی، درخت تصمیم، رگرسیون لجستیک، ماشین بردار پشتیبان
(ویجاia و سیواسانکار، ۲۰۱۹)	مخابرات	k نزدیک‌ترین همسایه، بیز ساده، جنگل تصادفی، درخت تصمیم، ماشین بردار پشتیبان
(کاروانا و همکاران، ۲۰۱۹)	بانکداری	بیز ساده، درخت تصمیم، رگرسیون لجستیک، شبکه عصبی، ماشین بردار پشتیبان
(هالیپاس و همکاران، ۲۰۱۹)	مخابرات	بیز ساده، تقویت گرادیان، جنگل تصادفی، درخت تصمیم، رگرسیون لجستیک
(جین و همکاران، ۲۰۲۱)	بانکداری، مخابرات، فناوری اطلاعات	تقویت گرادیان شدید، جنگل تصادفی، رگرسیون لجستیک، ماشین بردار پشتیبان
(رحمان و کومار، ۲۰۲۰)	بانکداری	k-نزدیک‌ترین همسایه، جنگل تصادفی، درخت تصمیم، ماشین بردار پشتیبان
(سید و همکاران، ۲۰۱۸)	بانکداری	درخت تصمیم
(لی و همکاران، ۲۰۲۱)	پخش رسانه‌ای	رگرسیون لجستیک

استفاده از یادگیری عمیق

در میان روش‌های یادگیری ماشین، از یادگیری عمیق نیز در برخی مقالات پیش‌بینی ریزش مشتری استفاده شده است. در اکثر این کارها از چندین لایه پیش‌خور با تابع فعال‌سازی یک‌سوساز^۳ به عنوان لایه‌های پنهان در ساختار شبکه عصبی

1. Sayed, Abdel-Fattah & Kholief

2. Li et al.

3. Rectifier (ReLU)

استفاده شده است. در ادامه به چند نمونه از این موارد اشاره می‌کیم. در مقاله اسپانودس و نگوین (۲۰۱۷) یک شبکه عصبی پیش‌خور عمیق برای پیش‌بینی ریزش مشتری پیاده‌سازی شده است. لایه‌های پنهان با تابع فعال‌سازی یک‌سوساز و لایه خروجی با دو نورون و تابع بیشینه هموار و به کارگیری روش‌های منظم‌سازی L1 و L2 پیاده‌سازی شده است. همچنین از عملیات حذف تصادفی^۱ برای قدرت تعیین بهتر استفاده شده است. برای آموزش مدل از الگوریتم بهینه‌سازی نزول گرادیان تصادفی همراه با تکانه و روش توقف زودهنگام استفاده شده است. در مقاله خان و همکاران (۲۰۱۹) از شبکه عصبی به صورت پرسپترون چندلایه برای پیش‌بینی ریزش مشتری در صنعت مخابرات استفاده شده است. در مقاله اومایاپارواتی و ایاکوتی^۲ (۲۰۱۷) از شبکه عصبی پیش‌خور و شبکه عصبی پیچشی یک‌بعدی برای یادگیری خودکار ویژگی از دادگان ورودی برای پیش‌بینی ریزش مشتری در صنعت مخابرات استفاده شده است. در مقاله دومینگوس، اوجمه و دارامولا^۳ (۲۰۲۱) تحلیل تجربی تأثیر ابرپارامترهای مختلف در استفاده از شبکه‌های عصبی عمیق برای پیش‌بینی ریزش مشتری در صنعت بانکی انجام شده است. از آزمایش‌های انجام‌شده نتیجه‌گیری شده است که شبکه‌های عصبی عمیق با استفاده از تابع فعال‌سازی یک‌سوساز در لایه‌های پنهان، نسبت به مدل‌های پرسپترون چندلایه عملکرد بهتری دارند. همچنین استفاده از الگوریتم بهینه سازی انتشار جذر میانگین مربعات در مقایسه با سایر الگوریتم‌ها، به یادگیری مدل با عملکرد بهتری منجر شده است.

روش‌شناسی پژوهش

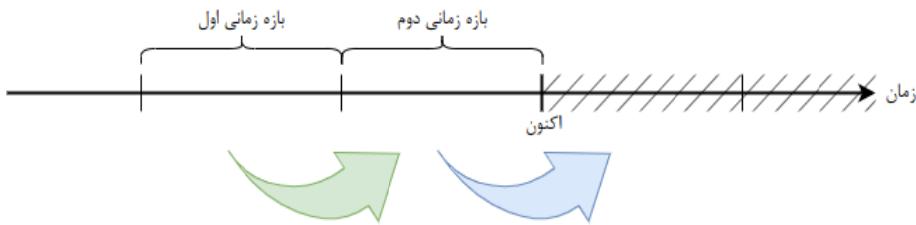
پیش‌بینی ریزش مشتری یک مسئله دسته‌بندی دودویی است. برای مشتری نام بدار ویژگی $x_i \in R^n$ تشکیل می‌شود و $\{y_i \in \{0, 1\}$ نشان‌دهنده متغیر هدف به معنای ریزش یا عدم ریزش است. بدین ترتیب با در اختیار داشتن مجموعه داده $\{(x_1, y_1), \dots, (x_m, y_m)\}$ ، هدف یادگیری تابع f است؛ به طوری که (x_i, f) تخمین مناسبی از مقدار y_i باشد. برای این منظور باید با استفاده از داده‌هایی که از مشتریان در اختیار است، ویژگی‌هایی را برای اشخاص جمع‌آوری کرد. با استفاده از این ویژگی‌ها، مشتریانی که احتمال روی‌گردانی آن‌ها در آینده وجود دارد، شناسایی می‌شوند.

برای پیش‌بینی ریزش مشتریان، دو بازه زمانی هماندازه و متواالی را در نظر می‌گیریم. انتهای بازه زمانی دوم باید قبل از زمان حال باشد تا دادگان ثبت شده در اختیار باشند. رفتار مشتریان در بازه زمانی اول را بررسی و ویژگی‌های رفتاری را برای مشتریان استخراج می‌کنیم. برای مشتریان کمیت معیار را در بازه‌های زمانی اول و دوم مقایسه کرده و متغیر هدف را مشخص می‌کنیم. بدین ترتیب با توجه به رفتار مشتری در بازه زمانی اول، سعی در پیش‌بینی متغیر هدف نسبت به بازه زمانی دوم داریم. پس از آموزش مدل دسته بند روی بازه زمانی اول، در نهایت از رفتار مشتری در بازه زمانی دوم استفاده می‌شود تا پیش‌بینی روی بازه زمانی متواالی در آینده صورت پذیرد. در شکل ۱ شما بی از نحوه انتخاب بازه‌های زمانی برای پیش‌بینی رفتار مشتری به تصویر کشیده شده است.

1. Dropout

2. Umayaparvathi & Iyakutti

3. Domingos, Ojeme & Daramola



شکل ۱. نحوه انتخاب بازه‌های زمانی برای پیش‌بینی رفتار مشتری

مهندسی ویژگی

در اجرای این طرح، دادگانی شامل اطلاعات مشتریان، اطلاعات حساب‌ها و کارت‌ها، مانده روزانه حساب‌ها و اطلاعات تراکنش‌های بانکی برای استفاده در دسترس بودند. با استفاده از اطلاعات ذکر شده به مهندسی ویژگی برای مشتریان بانک می‌پردازیم. مشتریان متشکل از اشخاص حقیقی و حقوقی هستند. برای اشخاص اطلاعات جمعیت‌شناسنامه مانند نوع مشتری، جنسیت، تاریخ تولد، وضعیت اقامت و غیره را در اختیار داریم. علاوه‌بر اطلاعات جمعیت‌شناسنامه، نیاز داریم تا ویژگی‌هایی از نحوه رفتار مشتریان نیز به دست آوریم. مشتریان بانک می‌توانند دارای حساب‌های مختلفی باشند و با آن‌ها به انجام تراکنش‌های مالی پردازند. در این پژوهش به منظور تحلیل رفتار مشتریان، اطلاعاتی را که برای حساب‌های مختلف به دست می‌آیند، تجمعی و ویژگی‌ها را به ازای اشخاص محاسبه می‌کنیم. در این پژوهش حساب‌های از نوع مشترک را در نظر نمی‌گیریم.

برای اشخاص مختلف، میانگین مانده مؤثر در طول بازه‌های زمانی اول و دوم را به طور جداگانه محاسبه می‌کنیم. طبق گفتوگو با متخصصان بانکی، مانده مؤثر معیار مهمی بوده و به معنای بخشی از سرمایه بانک است که متعلق به هر مشتری است و بر مبنای آن اهمیت مشتریان سنجیده می‌شود. میانگین مانده مؤثر در بازه زمانی مورد مطالعه را به عنوان ویژگی برای مشتری در نظر می‌گیریم. همچنین از مقایسه میانگین مانده مؤثر در بازه‌های زمانی اول و دوم، متغیر هدف را به دست می‌آوریم. اگر میانگین مانده مؤثر شخصی در بازه زمانی دوم نسبت به خود شخص در بازه زمانی اول کاهش شدیدی داشته باشد، می‌تواند به معنای خروج سرمایه از بانک باشد. بدین ترتیب درصد زیاد کاهش در میانگین مانده مؤثر مشتری را به عنوان ریزش مشتری در نظر می‌گیریم.

برای جمع‌آوری سایر ویژگی‌های رفتاری، تاریخچه تراکنش‌های مشتریان در بازه زمانی مورد مطالعه را بررسی می‌کنیم. تراکنش‌های صورت حساب مشتریان می‌توانند به صورت واریز یا برداشت باشند. همچنین تراکنش‌ها می‌توانند به صورت جزئی تراز انواع مختلفی مانند برداشت از پایانه فروش، پرداخت قبض، خرید شارژ، خرید اینترنتی و غیره باشند. از انواع دیگر تراکنش‌ها می‌توان به انتقال از/به کارت شتاب یا داخلی، انتقالی داخلی، پایا و ساتنا نیز اشاره کرد.

برای مشتریان از تراکنش‌های آن‌ها در سطوح مختلف آمارهایی از قبیل تعداد، مجموع و میانگین مبالغ، مبالغ کمینه و بیشینه را محاسبه می‌کنیم. ابتدا آمارهای ذکر شده را برای تراکنش‌ها به تفکیک واریز و برداشت محاسبه می‌کنیم. یعنی برای یک شخص در بازه زمانی مشخص به صورت کلی تعداد تراکنش‌های واریز و برداشت، مجموع و

میانگین مبالغ واردشده به حساب‌ها و مبالغ خارج شده از حساب‌ها، کمترین و بیشترین مبالغ واردشده و خارج شده مشخص می‌شوند. سپس همین آماره‌ها را برای تراکنش‌های واریز و برداشت به تفکیک نوع تراکنش نیز به دست می‌آوریم. برای مثال برای یک شخص در بازه زمانی مشخص تعداد برداشت از پایانه فروش و مجموع، میانگین، کمینه و بیشینه مبالغ برداشت شده مشخص می‌شوند. همچنین این اعداد برای نوع تراکنش‌های دیگر مانند انتقالی کارت داخلی و خارجی و انتقالی حساب داخلی و خارجی به دست می‌آیند. بدین ترتیب با استفاده از تراکنش‌های بانکی به ویژگی‌هایی برای مشتریان که نمایانگر رفتار مالی آن‌ها باشد می‌رسیم.

آماده‌سازی و پیش‌پردازش دادگان

دادگان را به مجموعه‌های آموزش، اعتبارسنجی و آزمایش تقسیم می‌کنیم. مدل‌های مختلف را روی دادگان آموزش یادگیری کرده و برای مقایسه، معیار مساحت زیر منحنی مشخصه عامل گیرنده نسبت به دادگان اعتبارسنجی را مدنظر قرار می‌دهیم. هر ویژگی اسمی را به تعداد مقادیر متمایزش به متغیرهای دودویی تبدیل می‌کنیم. هر ویژگی عددی را نیز به منظور داشتن توزیع نرمال با میانگین $0 = \mu$ و انحراف معیار $1 = \sigma$ مقیاس می‌کنیم. بدین ترتیب با الحاق کردن اعداد به دست آمده، برای هر مشتری به یک بردار $x \in \mathbb{R}^n$ می‌رسیم.

از چالش‌هایی که برای مسئله پیش‌بینی ریزش مشتری وجود دارد، نامتعادل بودن مجموعه دادگان است. بدین صورت که مشتریان مربوط به دسته ریزش، درصد کمتری از مجموعه دادگان را شامل می‌شوند. برای مقابله با این چالش، عکس نسبت اندازه دسته‌ها را به عنوان وزن دسته‌ها در نظر می‌گیریم. در این صورت طی مراحل آموزش مدل، وزن بیشتری به نمونه‌های دسته ریزش اختصاص داده می‌شود تا تعادل برقرار شود.

الگوریتم‌های دسته‌بندی یادگیری ماشین

ابتدا به عنوان دسته‌بندهای اولیه، مدل‌های k نزدیک‌ترین همسایه، درخت تصمیم، بیز ساده، رگرسیون لجستیک و ماشین بردار پشتیبان را آزمایش می‌کنیم. همچنین از روش‌های یادگیری جمعی جنگل تصادفی، تقویت تطبیقی و تقویت گرادیان نیز استفاده می‌کنیم.

استفاده از یادگیری عمیق

بعد پنهان مدل را برابر d در نظر گرفته و $x \in \mathbb{R}^n$ را فرض می‌کنیم. با انجام یک تبدیل خطی روی $h^{(0)}$ ، بعد آن را به d تغییر داده و سپستابع غیرخطی یک سوساز را اعمال می‌کنیم تا به $h^{(1)} \in \mathbb{R}^d$ برسیم. به این تبدیل صورت گرفته، لایه پیش‌پردازش می‌گوییم. در ادامه بازنمایی پنهان مدل در لایه $l-1$ ($l > 1$) برابر خواهد بود با:

$$h^{(l)} = \max(0, W_l h^{(l-1)} + b_l) \quad (1)$$

که $b_l \in \mathbb{R}^d$ و $W_l \in \mathbb{R}^{d \times d}$ خواهد بود. با تکرار لایه ذکر شده می‌توان شبکه‌هایی عمیق را تشکیل داد.

در ادامه به منظور سهولت آموزش شبکه‌های عمیق‌تر، اتصال باقی‌ماندهای (هی، زنگ، رن و سان^۱، ۲۰۱۶) به صورت $F(h) + h$ را به کار می‌گیریم. در حالت اول، بهارای هر لایه و در حالت دوم، برای هر دو لایه اتصال باقی‌ماندهای را برقرار می‌کنیم. در حالت تک لایه، تابع F مشابه رابطه ۱ خواهد بود. در حالت اتصال باقی‌ماندهای برای دو لایه، تابع را به صورت زیر در نظر می‌گیریم:

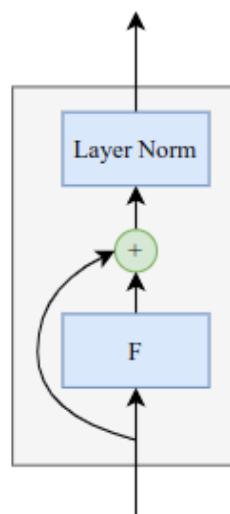
$$F(h) = W_2 \max(0, W_1 h + b_1) + b_2 \quad (\text{رابطه } 2)$$

که عبارت است از یک تبدیل خطی، اعمال تابع غیرخطی یکسوساز و سپس یک تبدیل خطی دیگر. برای انجام عمل جمع، ابعاد ورودی و خروجی تابع F را برابر در نظر می‌گیریم. در صورت قرارگیری دو لایه پشت سر هم، بعد بردار میانی حاصل از تبدیل خطی اول می‌تواند با بعد ورودی و خروجی برابر نباشد. بنابراین $W_2 \in \mathbb{R}^{d' \times d}$ و $W_1 \in \mathbb{R}^{d' \times d}$ که d' و d ممکن است نابرابر باشند.

همچنین برای کاهش زمان آموزش، نرمال‌سازی لایه (با، کیروس و هیتون^۲، ۲۰۱۶) را پس از اتصال باقی‌ماندهای قرار می‌دهیم. در این صورت بازنمایی پنهان مدل در لایه $l-1$ ($l > 1$) برابر خواهد بود با:

$$h^{(1)} = \text{Layer Norm}(F(h^{(1-1)}) + h^{(1-1)}) \quad (\text{رابطه } 3)$$

با تکرار بلوک‌های فوق می‌توان شبکه‌هایی عمیق را تشکیل داد. بلوک حاصل مشابه زیرلایه‌های استفاده شده در معماری تبدیل کننده (وسانی و همکاران^۳، ۲۰۱۷) است. نمایشی از این بلوک در شکل ۲ نشان داده شده است. همچنین نرمال‌سازی لایه را پس از لایه پیش‌پردازش قرار می‌دهیم.



شکل ۲. بلوک شامل اتصال باقی‌ماندهای و نرمال‌سازی لایه

1. He, Zhang, Ren & Sun

2. Ba, Kiros & Hinton

3. Vaswani et al.

سرانجام برای اجرای دسته بندی نهایی، روی بازنمایی پنهان حاصل از آخرين بلوک $h^{(L)}$ ، یک رگرسیون لجستیک اعمال می‌کنیم. در این رابطه، که $\hat{y} = \sigma(w^T h^{(L)} + b)$ تابع سیگموید و $P(y = 1 | x)$ است.

$$\hat{y} = \sigma(w^T h^{(L)} + b) \quad \text{رابطه ۴}$$

یافته‌های پژوهش

این پژوهش در یکی از بزرگ‌ترین بانک‌های ایران اجرا شده و آزمایش‌های انجام‌شده روی دادگان واقعی مشتریان بانک صورت پذیرفته است. در این آزمایش‌ها از اطلاعات جمعیت شناختی و رفتار گذشته مشتریان بهره گرفته شده و اطلاعات قابل شناسایی شخصی افراد استفاده‌ای نداشته است. در آزمایش‌های صورت گرفته، پیش‌بینی ریزش مشتری روی بازه زمانی یک ماهه انجام شده است. بدین ترتیب، دو بازه یک ماهه متولی در نظر گرفته شده و ویژگی‌های رفتاری مشتریان از بازه زمانی اول استخراج شده است. متغیر هدف نیز از مقایسه میانگین مانده مؤثر در بازه‌های زمانی اول و دوم به دست می‌آید. در صورتی که میانگین مانده مؤثر یک مشتری در بازه دوم نسبت به بازه اول با بیش از ۷۰ درصد کاهش همراه باشد، به عنوان ریزش در نظر گرفته می‌شود.

چنانچه گفته شد، تراکنش‌های بانکی شامل انواع مختلفی هستند. تعداد نوع‌های متمایز تراکنش‌ها زیاد بوده و برخی از آن‌ها برای مصارف خاص در نظر گرفته شده‌اند و به صورت مکرر استفاده نمی‌شوند. از این جهت برای استخراج ویژگی‌های رفتاری، نوع تراکنش‌های پُرتکرار مورد استفاده قرار گرفتند. بدین ترتیب تنها انواعی که در مجموع بیش از ۹۰ درصد تراکنش‌های انجام‌شده در بازه زمانی اول را شامل می‌شوند، برای محاسبه آماره‌ها انتخاب شدند.

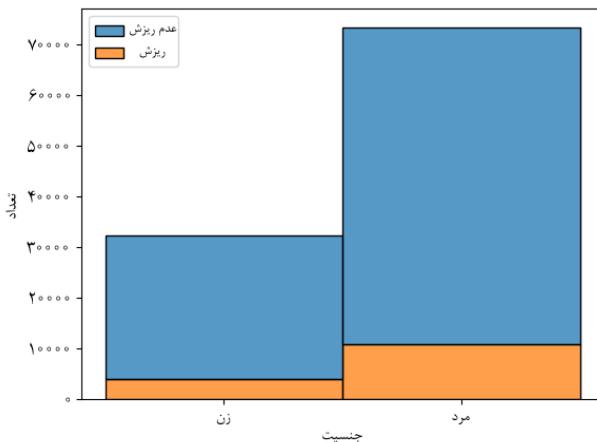
در مراحل پالایش دادگان، مشتریانی که در اطلاعات ثبت‌شده برای آن‌ها نقص یا ناسازگاری وجود داشت، کنار گذاشته شدند. به علت تعداد زیاد مشتریان بانک، برای انجام پیش‌بینی ریزش زیرمجموعه‌ای از اشخاص انتخاب و جداسازی شدند. در اطلاعات موجود برای تراکنش‌های انتقالی کارت داخلی، مبدأ و مقصد جابه‌جایی پول مشخص است. بدین ترتیب برای امکان بررسی تأثیر ارتباطات مالی میان اشخاص، جامعه هدف به مشتریانی که در بازه زمانی مورد مطالعه، تراکنش انتقالی کارت داخلی داشتند، محدود شد. مشتریان به دست آمده بر اساس میانگین مانده مؤثر در بازه زمانی مورد مطالعه، به صورت نزولی مرتب‌سازی شدند. برای آزمایش‌های مختلف با اندازه مجموعه داده آموزش متفاوت، تعدادی از مشتریان با مانده مؤثر بیشتر به صورت تصادفی انتخاب شدند.

دادگانی به اندازه تقریبی ۱۱۰۰۰ نمونه به عنوان مجموعه اعتبارسنجی جداسازی شدند. سه مجموعه با اندازه‌های مختلف به عنوان دادگان آموزش در نظر گرفته شدند. اندازه تقریبی هر مجموعه و درصد دسته‌های ریزش و عدم ریزش در جدول ۲ نشان داده شده است. عدم تعادل مجموعه دادگان نیز در جدول ۲ مشاهده می‌شود. به ازای هر مجموعه آموزش، مدل‌ها را روی آن مجموعه یادگیری بررسی و عملکرد مدل روی مجموعه اعتبارسنجی را گزارش می‌کنیم. در این جدول K به معنای هزار است.

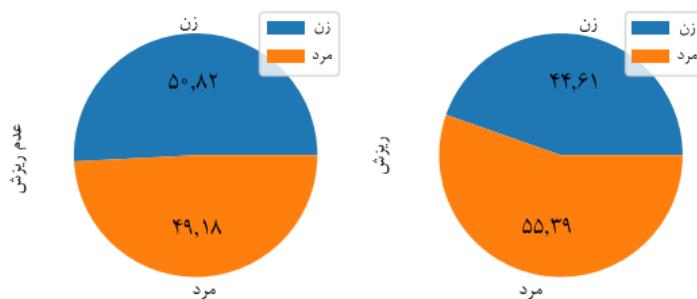
جدول ۲. اندازهٔ مجموعه‌های آموزش و درصد دسته‌های ریزش و عدم ریزش در هر مجموعه

	ریزش	ریزش عدم
۱۰۶K	۱۳/۷۱ درصد	۸۶/۲۹ درصد
۲۳۶K	۱۳/۷۵ درصد	۸۶/۲۵ درصد
۱۰۲۷K	۱۳/۷۷ درصد	۸۶/۲۳ درصد

در ادامه، به برخی از ویژگی‌های موجود در مجموعه آموزش ۱۰۶K نگاه می‌کنیم. در شکل ۳، نمودار ستونی فراوانی مقادیر متمایز ویژگی جنسیت در میان مشتریان حقیقی موجود در مجموعه داده ۱۰۶K نمایش داده شده است. همچنین در هر ستون، تعداد مشتریان بر حسب متغیر هدف با رنگ متفاوت مشخص شده است. مشاهده می‌شود که تعداد بیشتری از اشخاص حقیقی موجود در مجموعه داده را مردان تشکیل می‌دهند. در شکل ۴، درصد مردان و زنان به تفکیک دسته‌های ریزش و عدم ریزش در میان مشتریان حقیقی مجموعه داده ۱۰۶K به صورت نمودار دایره‌ای نمایش داده شده است. مشاهده می‌شود که ریزش در میان مردان بیشتر از زنان اتفاق افتاده است.

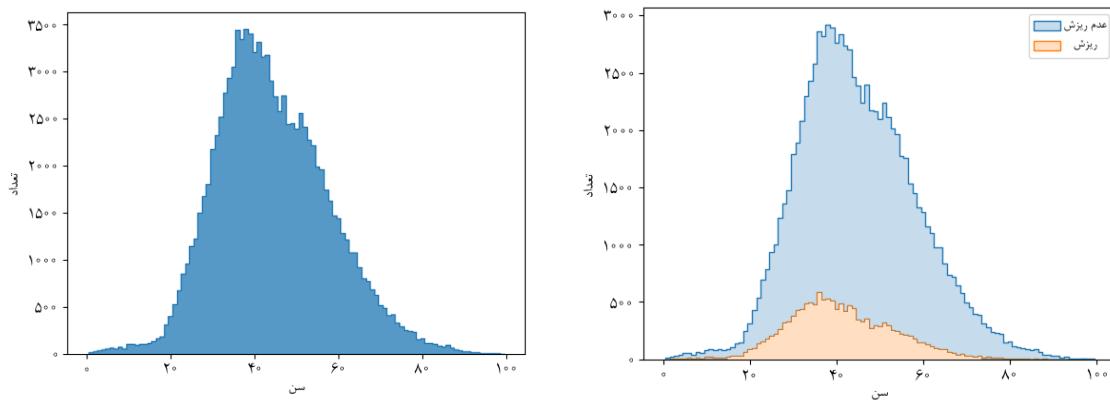


شکل ۳. نمودار فراوانی مقادیر متمایز ویژگی جنسیت در میان مشتریان حقیقی مجموعه داده ۱۰۶K

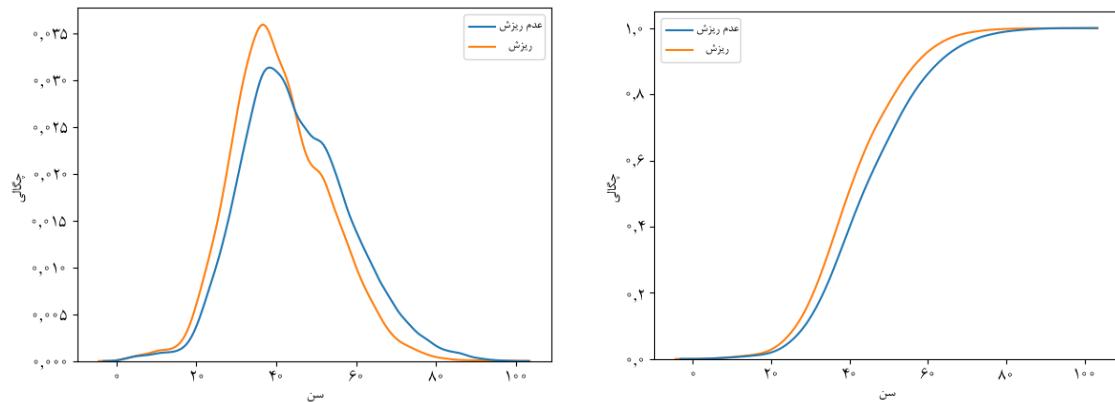


شکل ۴. نمودار دایره‌ای درصد شمول مردان و زنان، به تفکیک تعلق به دسته ریزش در مشتریان حقیقی مجموعه داده ۱۰۶K

در شکل ۵ توزیع ویژگی سن در میان مشتریان حقیقی مجموعه داده $106K$ به صورت کلی و به تفکیک تعلق به دسته ریزش نمایش داده شده است. در شکل ۶ نیز نمودار تخمین چگالی کرنلی ویژگی سن مشتریان حقیقی، به تفکیک تعلق به دسته ریزش، به صورت غیرتجمعی و تجمعی نمایش داده شده است. مشاهده می‌شود که ریزش در سنین پایین‌تر، بیشتر اتفاق افتاده است.



شکل ۵. نمودارهای توزیع ویژگی سن به صورت کلی و به تفکیک تعلق به دسته ریزش برای اشخاص حقیقی مجموعه داده $106K$



شکل ۶. نمودارهای تخمین چگالی کرنلی ویژگی سن به تفکیک تعلق به دسته ریزش به صورت غیر تجمعی و تجمعی برای اشخاص حقیقی مجموعه داده $106K$

در جدول ۳ آماره‌هایی از ویژگی میانگین مانده مؤثر و ویژگی‌های مربوط به تراکنش‌های واریز و برداشت نمایش داده شده است. اعداد ارائه شده برای مبالغ، به ریال هستند. همان طور که مشاهده می‌شود، ویژگی میانگین مانده مؤثر دارای انحراف زیادی از میانگین است. با توجه به محل تجمع دادگان، توزیع این ویژگی چولگی مثبت دارد. این اتفاق به علت وجود اشخاصی با میانگین مانده مؤثر بیشتر نسبت به اکثریت مشتریان موجود در مجموعه داده رخ داده است.

اشخاص حقیقی موجود در مجموعه داده $106K$ در بازه زمانی یک ماهه مورد مطالعه، به طور میانگین تعداد ۷۸ تراکنش برداشت داشته‌اند و حدود نیمی از اشخاص، بیش از ۶۱ برداشت از حساب انجام داده‌اند. حدود نیمی از مشتریان در تراکنش‌های برداشت خود، در مجموع بیش از $406/8$ میلیون ریال و به طور میانگین بیش از $6/9$ میلیون ریال، مبالغی را پرداخت کرده‌اند. همچنین حدود نیمی از مشتریان در تراکنش‌های واریز به حساب، در مجموع کمتر از 464 میلیون ریال و به طور میانگین کمتر از $30/3$ میلیون ریال، مبالغی را دریافت کرده‌اند. بیشترین مبلغ برداشت شده برای حدود ۲۵ درصد از اشخاص بیش از 400 میلیون ریال بوده است. همچنین، بیشترین مبلغ واریز شده برای حدود ۲۵ درصد از مشتریان کمتر از 50 میلیون ریال بوده است. سایر اطلاعات آماری مشابه نیز در جدول ۳ مشاهده می‌شود.

جدول ۳. آماره‌هایی از ویژگی‌های میانگین مانده مؤثر و ویژگی‌های مربوط به تراکنش‌های واریز و برداشت در مجموعه داده $106K$

ویژگی	میانگین	معیار انحراف	۲۵ صد	۵۰ صد	۷۵ صد
مؤثر مانده میانگین	۸۱۴/۷ میلیون	۳/۰ میلیارد	۱۶۸/۹ میلیون	۲۹۰/۳ میلیون	۶۳۸/۱ میلیون
برداشت تعداد	۷۸	۷۰	۲۷	۶۱	۱۱۰
مبالغ مجموع	۱/۶ میلیارد	۶/۵ میلیارد	۱۲۵/۷ میلیون	۴۰۶/۸ میلیون	۱/۳ میلیارد
مبالغ میانگین	۲۶/۱ میلیون	۱۱۳/۶ میلیون	۲/۶ میلیون	۶/۹ میلیون	۲۰/۶ میلیون
مبالغ کمینه	۴۶۱/۰ هزار	۱۳/۹ میلیون	۱/۲ هزار	۳/۰ هزار	۲۰/۰ هزار
مبالغ بیشینه	۴۵۵/۴ میلیون	۱/۶ میلیارد	۳۰/۰ میلیون	۱۰۰/۰ میلیون	۴۰۰/۰ میلیون
واریز تعداد	۲۴	۷۵	۷	۱۴	۲۷
مبالغ مجموع	۱/۸ میلیارد	۸/۶ میلیارد	۱۴۰/۴ میلیون	۴۶۴/۰ میلیون	۱/۴ میلیارد
مبالغ میانگین	۸۴/۳ میلیون	۲۲۱/۹ میلیون	۱۲/۵ میلیون	۳۰/۳ میلیون	۷۹/۱ میلیون
مبالغ کمینه	۳/۷ میلیون	۵۶/۱ میلیون	۴/۳ هزار	۱۰۰/۰ هزار	۷۳۷/۵ هزار
مبالغ بیشینه	۵۴۸/۴ میلیون	۱/۷ میلیارد	۵۰/۰ میلیون	۱۰۴/۷ میلیون	۵۰۰/۰ میلیون

در جدول ۴ آماره‌هایی از ویژگی‌های مربوط به تراکنش‌های برداشت از پایانه فروش، به عنوان نمونه‌ای از نوع تراکنش‌های مختلف، نشان داده شده است. مشاهده می‌شود که اشخاص حقیقی موجود در مجموعه داده $106K$ در بازه زمانی یک ماهه مورد مطالعه، به طور میانگین ۵۰ بار تراکنش برداشت از پایانه فروش انجام داده‌اند. همچنین در حدود نیمی از مشتریان حقیقی بیش از ۳۷ بار برداشت از پایانه فروش داشته‌اند. حدود نیمی از اشخاص در تراکنش‌های برداشت شده از پایانه فروش، در مجموع بیش از $74/۹$ میلیون ریال و به طور میانگین بیش از $1/۸$ میلیون ریال، مبالغی را پرداخت کرده‌اند. کمترین مبلغ برداشت شده در یک تراکنش توسط یک شخص، برای حدود ۲۵ درصد از مشتریان بیش از 85 هزار ریال بوده است. به همین ترتیب، بیشترین مبلغ برداشت شده در یک تراکنش برای حدود ۲۵ درصد از مشتریان بیشتر از 64 میلیون ریال بوده است. سایر اطلاعات مشابه، در جدول ۴ درج شده است.

جدول ۴. آمارهایی از ویژگی‌های مربوط به تراکنش‌های برداشت از پایانه فروش در مجموعه داده ۱۰۶K

ویژگی	میانگین	انحراف معیار	صد ک	۵۰	صد ک	۷۵
تعداد تراکنش‌ها	۵۰	۴۹	۱۴	۳۷	۷۴	۷۴
مجموع مبالغ	۲۰۷/۱	۴۶۴/۸	۲۵/۶	۷۴/۹	۲۰۰/۷	۴ میلیون
میانگین مبالغ	۷/۷	۳۴/۴	۹۳۹/۸	۱/۸	۴/۶	۴ میلیون
کمترین مبلغ	۱/۵	۲۳/۰	۱۸/۰	۳۰/۰	۸۵/۰	هزار ۸۵/۰
بیشترین مبلغ	۶۹/۹	۱۷۱/۶	۶/۰	۱۶/۲	۶۴/۰	۶۴ میلیون

جزئیات آموزش

مدل‌های مبتنی بر شبکه عصبی را با استفاده ازتابع ضرر آنتروپی متقطع دودویی و الگوریتم بهینه‌سازی برآورد گشتاور تطبیقی (کینگما و با^۱، ۲۰۱۴) با نرخ یادگیری اولیه ۰/۰۱ آموزش می‌دهیم. در محاسبه کردن تابع ضرر، وزن بیشتری به نمونه‌های دسته ریزش اختصاص داده می‌شود تا تعادل برقرار شود. برای بررسی و مقایسه عملکرد مدل‌ها در دسته‌بندی دودویی، از معیار مساحت زیرمنحنی مشخصه عامل گیرنده استفاده می‌کنیم. مدل‌هایی با پیکربندی‌های مختلف را روی مجموعه‌های آموزش یادگیری کرده و نتایج عملکردی آن‌ها روی دادگان اعتبارسنجی را به دست می‌آوریم. یادگیری مدل در صورت سپری شدن سه دوره متوالی و عدم افزایش مساحت زیرمنحنی مشخصه عامل گیرنده نسبت به مجموعه اعتبارسنجی، متوقف می‌شود. در تمامی مدل‌های شرح‌داده شده پس از اعمال تابع F در هر بلوک و همچنین پس از لایه پیش‌پردازش، عملیات حذف تصادفی با نرخ ۰/۱۰ اعمال شده است.

آزمایش‌ها

به عنوان مبنایی برای مقایسه‌ها، از چندین مدل یادگیری ماشین سنتی پُراستفاده، شامل درخت تصمیم، بیز ساده، کنزدیک‌ترین همسایه، رگرسیون لجستیک و ماشین بردار پشتیبان استفاده می‌کنیم. همچنین الگوریتم‌هایی از یادگیری جمعی شامل جنگل تصادفی، تقویت تطبیقی و تقویت گرادیان مورد استفاده قرار گرفته‌اند. به‌منظور یادگیری مدل‌های ذکر شده، در موارد ممکن برای برقرار تعادل میان دسته‌های ریزش و عدم ریزش با توجه به نسبت اندازه دسته‌ها وزن بیشتری به دسته اقلیت داده شده است. نتایج به دست آمده از این آزمایش‌ها در جدول ۵ آمده است. هر کدام از مدل‌ها به‌ازای پیکربندی‌های متفاوت روی مجموعه‌های آموزش ۱۰۶K و ۲۳۶K یادگیری شده و تنها بهترین نتیجه به دست آمده روی مجموعه اعتبارسنجی، بر حسب معیار مساحت زیرمنحنی مشخصه عامل گیرنده گزارش شده است. مشاهده می‌شود که در میان روش‌های سنتی یادگیری ماشین، درخت تصمیم بهترین عملکرد را با توجه به معیار مساحت زیرمنحنی نتیجه داده است. همچنین روش‌های جمعی نسبت به مدل‌های تکی نتایج بهتری را کسب کرده‌اند که از میان آن‌ها روش تقویت گرادیان منجر به به دست آمدن بیشترین مساحت زیر منحنی مشخصه عامل گیرنده شده است.

**جدول ۵. عملکرد الگوریتم‌های یادگیری ماشین سنتی به همراه روش‌های جمعی روی مجموعه‌های آموزش مختلف
(برحسب معیار مساحت زیر منحنی مشخصه عامل گیرنده)**

۲۳۶K	۱۰۶K	مدل
۰/۷۴۵۴	۰/۷۴۵۴	بیز ساده
۰/۷۴۸۹	۰/۷۵۰۳	K نزدیک‌ترین همسایه
۰/۷۸۵۱	۰/۷۷۷۰	ماشین بردار پشتیبان
۰/۸۳۹۶	۰/۸۳۸۲	رگرسیون لجستیک
۰/۸۵۹۷	۰/۸۵۳۱	درخت تصمیم
۰/۸۷۹۱	۰/۸۷۴۱	جنگل تصادفی
۰/۸۷۲۵	۰/۸۷۱۶	تقویت تطبیقی
۰/۹۰۱۰	۰/۸۹۸۴	تقویت گرادیان

در جدول ۶ نتایج بدست‌آمده از مدل‌های حاصل از تکرار لایه‌های کاملاً متصل به صورت متوالی آمده است. حدودی از تعداد پارامترهای قابل یادگیری مدل‌ها به ازای مقادیر مختلف d و L ذکر شده است. مقدار $L = 0$ به معنای قرارگیری لایه دسته‌بندی کننده نهایی بلا فاصله بعد از لایه پیش‌پردازش است. اعداد ذکر شده، نشان‌دهنده معیار مساحت زیرمنحنی مشخصه عامل گیرنده نسبت به مجموعه اعتبارسنجی به ازای مجموعه‌های آموزش با اندازه متفاوت است. مشاهده می‌شود که بهترین نتیجه بدست‌آمده به‌ازای هر مجموعه آموزش با افزایش اندازه بهبود یافته است؛ اما افزایش تعداد لایه‌های پشت‌سرهم قرارگرفته، فرایند یادگیری مدل‌ها را دشوار کرده است و اغلب به کاهش معیار مساحت زیر منحنی منجر شده است.

جدول ۶. نتایج مدل‌های حاصل از انباسته کردن چندین لایه کاملاً متصل پس از لایه پیش‌پردازش

d	L	۱۰۶K	۲۳۶K	۱۰۲۷K	تعداد پارامتر ($\times 10^{4}$)
۱۲۸	۰	۰/۸۶۲۱	۰/۸۷۴۳	۰/۸۷۵۰	۱۸
	۱	۰/۸۶۳۴	۰/۸۶۵۱	۰/۸۶۶۸	۳۴
	۲	۰/۸۳۶۵	۰/۸۴۰۹	۰/۸۲۰۲	۵۱
	۳	۰/۷۴۰۷	۰/۵۰۰۰	۰/۵۰۰۰	۶۷
	۴	۰/۵۰۰۰	۰/۷۷۶۸	۰/۵۰۰۰	۸۴
۲۵۶	۰	۰/۸۶۴۸	۰/۸۷۸۹	۰/۸۷۱۹	۳۶
	۱	۰/۸۵۳۷	۰/۸۶۵۶	۰/۸۵۹۲	۱۰۱
	۲	۰/۸۰۴۰	۰/۸۱۰۹	۰/۵۰۰۰	۱۶۷
	۳	۰/۵۰۰۰	۰/۵۰۰۰	۰/۵۰۰۰	۲۳۳
	۴	۰/۶۶۴۰	۰/۵۰۰۰	۰/۵۰۰۰	۲۹۹

در جدول ۷ نتایج مدل‌هایی که برای هر لایه کاملاً متصل، اتصال باقی‌ماندهای و نرمال‌سازی لایه قرار داده شده، آمده است. مقدار $L = 0$ به معنای اعمال نرمال‌سازی لایه پس از لایه پیش‌پردازش و سپس قرارگیری لایه دسته‌بندی کننده نهایی است. در مقایسه با آزمایش قبل، مشاهده می‌شود که استفاده از اتصال‌های پرشی به صورت باقی‌ماندهای و واحدهای نرمال‌سازی لایه، فرایند یادگیری مدل‌های عمیق‌تر را راحت‌تر و پایدار‌تر می‌سازند. افزایش اندازه مجموعه آموزش نیز، به بهبود عملکرد مدل روی اعتبارسنجی منجر شده است.

جدول ۷. نتایج حاصل از مدل‌های شامل اتصال باقی‌ماندهای برای هر یک لایه کاملاً متصل و نرمال‌سازی لایه

d	L	۱۰۶K	۲۳۶K	۱۰۲۷K	تعداد پارامتر ($\times 10^{+0}$)
۱۲۸	۰	۰/۸۸۰۶	۰/۸۸۷۷	۰/۸۸۹۷	۱۸
	۱	۰/۸۸۲۵	۰/۸۹۰۹	۰/۸۹۵۵	۳۵
	۲	۰/۸۸۰۵	۰/۸۸۵۹	۰/۸۹۴۸	۵۱
	۳	۰/۸۸۱۷	۰/۸۹۰۶	۰/۸۹۳۶	۶۸
	۴	۰/۸۸۱۰	۰/۸۸۹۷	۰/۸۹۵۸	۸۵
۲۵۶	۰	۰/۸۸۱۲	۰/۸۸۶۰	۰/۸۹۱۵	۳۶
	۱	۰/۸۸۰۷	۰/۸۸۵۸	۰/۸۹۴۳	۱۰۲
	۲	۰/۸۸۰۷	۰/۸۸۸۳	۰/۸۹۴۸	۱۶۹
	۳	۰/۸۸۰۱	۰/۸۸۷۹	۰/۸۹۳۴	۲۳۵
	۴	۰/۸۷۹۳	۰/۸۸۶۷	۰/۸۹۳۰	۳۰۱

در جدول ۸، نتایج به دست آمده از مدل‌های متشكل از اتصال باقی‌ماندهای برای هر دو لایه کاملاً متصل و نرمال‌سازی لایه آمده است. مشابه با نتایج قبل، در این آزمایش افزایش اندازه مجموعه دادگان آموزش، به بهبود عملکرد مدل‌ها منجر شده است. با مقایسه نتایج جداول ۷ و ۸، مشاهده می‌کنیم که بهترین عملکردهای به دست آمده به ازای هر مجموعه آموزش، بسیار نزدیک به هم هستند؛ اما تعداد پارامترهای قابل یادگیری بهترین مدل‌ها در جدول ۷، بسیار کمتر از بهترین مدل‌ها در جدول ۸ است؛ بنابراین استفاده از مدل‌های مشتمل بر اتصال باقی‌ماندهای برای هر یک لایه کاملاً متصل، به لحاظ هزینه محاسباتی مناسب‌تر به نظر می‌آید.

بدین ترتیب مدل‌های پُراستفاده در پیشینه پژوهش پیاده‌سازی شده و به‌منظور مقایسه نتایج با روش ارائه شده، بر روی مجموعه داده به دست آمده در این تحقیق اعمال شدند. به طور کلی مشاهده می‌کنیم که با استفاده از الگوریتم‌های یادگیری جمعی و همچنین روش‌های یادگیری عمیق، عملکرد مدل‌ها نسبت به مدل‌های مبنا بهبود می‌یابند. همچنین استفاده از مجموعه آموزش بزرگ‌تر، به بهبود عملکرد مدل روی مجموعه اعتبارسنجی منجر می‌شود؛ بنابراین می‌توان از مجموعه داده‌های بزرگ‌تر و روش‌های یادگیری عمیق بهره گرفت و پیش‌بینی ریزش مشتری را بهبود بخشد.

جدول ۸. نتایج حاصل از مدل‌های شامل اتصال باقی‌مانده‌ای برای هر دو لایه متوالی و نرمال‌سازی لایه

L	d	d'	۱۰۶K	۲۳۶K	۱۰۲۷K	تعداد پارامتر ($1000 \times$)
۱	۱۲۸	۱۲۸	۰/۸۸۱۱	۰/۸۸۹۰	۰/۸۹۴۹	۵۱
		۲۵۶	۰/۸۸۰۴	۰/۸۸۷۹	۰/۸۹۴۲	۸۴
		۵۱۲	۰/۸۸۲۳	۰/۸۹۰۲	۰/۸۹۴۹	۱۵۰
۲	۲۵۶	۲۵۶	۰/۸۸۲۱	۰/۸۹۰۸	۰/۸۹۵۰	۱۶۸
		۵۱۲	۰/۸۸۰۴	۰/۸۸۷۰	۰/۸۹۵۷	۳۰۰
		۱۰۲۴	۰/۸۷۹۵	۰/۸۸۸۱	۰/۸۹۵۸	۵۶۲
۲	۱۲۸	۱۲۸	۰/۸۷۹۱	۰/۸۸۸۳	۰/۸۹۳۹	۸۴
		۲۵۶	۰/۸۷۹۶	۰/۸۸۸۷	۰/۸۹۳۴	۱۵۰
		۵۱۲	۰/۸۸۲۲	۰/۸۸۷۶	۰/۸۹۳۵	۲۸۲
۲	۲۵۶	۲۵۶	۰/۸۸۰۴	۰/۸۸۸۹	۰/۸۹۵۳	۳۰۰
		۵۱۲	۰/۸۷۹۴	۰/۸۸۹۶	۰/۸۹۵۲	۵۶۳
		۱۰۲۴	۰/۸۷۸۸	۰/۸۸۷۴	۰/۸۹۴۵	۱۰۸۸

مدل بوتر

در جدول ۹ دقیق دسته‌بندی در پیش‌بینی ریزش مشتری برای بهترین مدل‌های به دست آمده از بخش قبل روی مجموعه اعتبارسنجی ذکر شده است. دقیق‌های گزارش شده به ازای حد آستانه‌ای هستند که به بدست آمدن بیشترین امتیاز F1 برای هر مدل روی مجموعه اعتبارسنجی منجر شده است. باید توجه شود که نتایج به دست آمده برای مدل عمیق به صورت تکی و مدل تقویت گرادیان به صورت جمعی بوده است. عملکرد مدل عمیق پیشنهادشده از تمام روش‌های سنتی به صورت تکی و همچنین دو روش یادگیری جمعی استفاده شده دیگر، یعنی جنگل تصادفی و تقویت تطبیقی، برتر بوده است و با نتایج به دست آمده از الگوریتم تقویت گرادیان رقابت می‌کند. ممکن است با استفاده از مدل عمیق در قالب یادگیری جمعی، بهبود عملکردی بیشتری نیز حاصل شود.

جدول ۹. دقیق دسته‌بندی الگوریتم تقویت گرادیان و مدل عمیق بر روی مجموعه اعتبارسنجی

۲۳۶K	۱۰۶K	مدل
۰/۸۷۲۴	۰/۸۷۴۴	گرادیان تقویت
۰/۸۵۴۱	۰/۸۵۶۲	عمیق

همچنین برای آموزش تقویت گرادیان و سایر مدل‌های جمعی استفاده شده روی یک ماشین، باید کل مجموعه داده در حافظه قرار گیرد؛ اما در خصوص مدل‌های مبتنی بر شبکه عصبی، می‌توان مجموعه داده را در دسته‌های کوچک داخل حافظه بارگذاری کرد و آموزش مدل را با استفاده از دادگان بیشتر انجام داد. در مجموع بهترین عملکرد در پیش‌بینی

ریزش مشتری پس از آموزش روی مجموعه داده‌های $K = 106$ و $K = 236$ از منظر معیارهای مساحت زیر منحنی مشخصه عامل گیرنده و دقت دسته‌بندی نسبت به مجموعه اعتبارسنجی، توسط مدل تقویت گردایان بهدست آمد.

نتیجه‌گیری و پیشنهادها

در این مقاله مسئله پیش‌بینی ریزش مشتریان بانک را مورد بررسی قرار دادیم. هدف از مسئله پیش‌بینی ریزش، شناسایی مشتریانی است که احتمال روی‌گردانی آن‌ها در آینده وجود دارد. این پژوهش روی داده‌های واقعی مشتریان یکی از بزرگ‌ترین بانک‌های ایران انجام شد. با استفاده از حجم زیادی از تراکنش‌های بانکی و تجمعی آن‌ها در سطوح متفاوت، ویژگی‌های مختلفی برای مشتریان استخراج کردیم. به‌منظور پیش‌بینی ریزش، روش‌های یادگیری ماشین را روی ویژگی‌های مشتریان اعمال کردیم.

برای پیش‌بینی رفتار مشتریان، دو بازه زمانی متوالی را در نظر گرفته و با استفاده از رفتار مشتری در بازه اول، تلاش کردیم تا متغیر هدف در بازه دوم را پیش‌بینی کنیم. بر مبنای میزان تغییر میانگین مانده مؤثر مشتری در بازه زمانی دوم نسبت به بازه اول، متغیر هدف را مشخص کردیم. با تجمعی تراکنش‌های مشتریان در سطوح مختلف، آماره‌هایی از رفتار مالی مشتریان را محاسبه کردیم و به عنوان ویژگی‌های رفتاری در نظر گرفتیم. یک مجموعه اعتبارسنجی ثابت برای مقایسه مدل‌ها و تنظیم ابرپارامترها و سه مجموعه آموزش با اندازه‌های مختلف برای یادگیری مدل‌ها جداسازی کردیم. برای مشکل عدم تعادل مجموعه داده، عکس نسبت اندازه دسته‌ها را به عنوان وزن دسته‌ها در فرایند آموزش تأثیر دادیم تا تعادل برقرار شود.

به عنوان مبنای برای مقایسه، تعدادی از الگوریتم‌های پُراستفاده یادگیری ماشین سنتی را روی دادگان اعمال کردیم. علاوه‌بر آن‌ها، تعدادی از روش‌های یادگیری جمعی مرسوم را نیز استفاده کردیم. در ادامه، روش‌های یادگیری عمیق را استفاده کردیم و تعدادی از واحدهای نوین آن را به کار گرفتیم. آزمایش‌های گسترده‌ای را روی روش‌های ذکر شده به طور جداگانه انجام دادیم و نتایج را گزارش کردیم.

در نتایج بهدست آمده مشاهده کردیم که روش‌های یادگیری جمعی و مدل‌های یادگیری عمیق، در پیش‌بینی ریزش مشتری، نسبت به مدل‌های مبنای عملکرد بهتری را ارائه می‌دهند. افزایش اندازه مجموعه آموزش، در بهبود عملکرد مدل‌ها تأثیر دارد. در میان الگوریتم‌های دسته‌بندی یادگیری ماشین سنتی آموزش دیده روی مجموعه دادگان $K = 106$ درخت تصمیم بیشترین مساحت زیرمنحنی مشخصه عامل گیرنده به مقدار $8531/0$ را کسب کرد. در مجموع روی دادگان $K = 106$ مدل تقویت گردایان با $8984/0$ و پس از آن، مدل مبتنی بر یادگیری عمیق با $8825/0$ ، بیشترین مساحت زیرمنحنی مشخصه عامل گیرنده نسبت به مجموعه اعتبارسنجی را بهدست آوردند.

پس از تحلیل رفتار و شناسایی مشتریان در شرف روی‌گردانی، می‌توان از پیشنهادهای کاربردی ذیل برای جلوگیری از ریزش و حفظ مشتری استفاده کرد. تفکیک مشتریان بر اساس سن، شغل، تحصیلات و غیره، به‌منظور ارائه خدمات و تولید محصولات بانکی بر این مبنای ایجاد تنوع در خدمات موجود، ارائه خدمات مورد نیاز مشتریان از طریق

بسترها مجازی و در صورت نیاز در محل فعالیت و زندگی مشتریان، تسهیل در ارائه خدمات مطلوب به مشتریان، افزایش اعتماد مشتریان از طریق ارائه کاربردی، امن و همچنین حفاظت از اطلاعات مشتریان و غیره می‌تواند به حفظ مشتری کمک کند.

به علت حجم زیاد دادگان، محدودیت منابع ساخت افزاری و حافظه ذخیره‌سازی در اختیار، محدودیت زمانی و زمان بر بودن فرایند استخراج، پالایش و بارگذاری دادگان از پایگاه داده بانک، این تحقیق برای مطالعه رفتار مشتریان روی بازه زمانی یک ماهه انجام شده است. در ادامه راه، می‌توان با تخصیص منابع ساخت افزاری بیشتر، پیش‌بینی ریزش مشتری را روی بازه‌های زمانی بلندتری انجام داد. می‌توان علاوه بر تراکنش‌های مالی، از دادگان سایر زمینه‌های بانکی مانند تسهیلات و غیره نیز برای استخراج ویژگی‌های مشتریان استفاده کرد. همچنین، انتظار داریم که چارچوب ایجاد شده و راه کار ارائه شده در این طرح در تحلیل و پیش‌بینی رفتار مشتریان برای سایر مسائل موجود در زمینه مدیریت ارتباط با مشتری نیز کاربرد داشته باشد.

منابع

- احمدی کوشا، آزاده؛ احمدی، فائق؛ رنجبر، محمدحسین و کردنلوئی، حمیدرضا (۱۴۰۳). شناسایی شاخص‌های اعتبارسنجی و رتبه‌بندی مشتریان در تسهیلات خرد در بانک خاورمیانه. *تحقیقات مالی*، ۲(۲۶)، ۴۱۵-۴۳۸.
- احمدی سرتختی، فرشید؛ هژبر کیانی، کامبیز؛ حسینی، سید شمس الدین و معمارنژاد، عباس (۱۴۰۲). طراحی مدلی برای ارزیابی ریسک اعتباری مشتریان ضمانت‌نامه‌های صادر شده توسط صندوق ضمانت صادرات ایران با کمک مدل شبکه عصبی مصنوعی. *تحقیقات مالی*، ۴(۲۵)، ۶۹۵-۷۱۶.
- باجلان، سعید؛ فلاج‌پور، سعید و رئیسی، سارا (۱۴۰۳). بهینه‌سازی پرتفوی اعتباری بانک‌ها با استفاده از رویکرد اکچوئری و شبکه عصبی مصنوعی. *تحقیقات مالی*، ۲۶(۳)، ۷۱۰-۷۳۳.
- رحمی، سعیده؛ رosta، علی و آسایش، فرزاد (۱۴۰۳). ارزیابی ارتباط میان عوامل ارتقای توان رقابت‌پذیری خدمات ارزی مشتریان در صنعت بانکداری. *تحقیقات مالی*، ۲(۲۶)، ۴۳۹-۴۶۲.

References

- Ahmad, A. K., Jafar, A., & Aljoumaa, K. (2019). Customer churn prediction in telecom using machine learning in big data platform. *Journal of Big Data*, 6(1), 1-24.
- Ahmadi Kousha, A., Ahmadi, F., Ranjbar, M. & Kordlouie, M. (2024). Validation Indicator Identification and Customer Ranking in Microloans: A Study at Middle East Bank in Iran. *Financial Research Journal*, 26(2), 399-423. (in Persian)
- Ahmadi Sartakhti, F., Hojabr Kiani, K., Hoseini, S. & Memarnejad, A. (2023). Designing a Model for Credit Risk Assessment of Customers for Guarantees Issued by the Export Guarantee Fund of Iran via Artificial Neural Network Model. *Financial Research Journal*, 25(4), 641- 660. (in Persian)

- Amin, A., Anwar, S., Adnan, A., Nawaz, M., Alawfi, K., Hussain, A. & Huang, K. (2017). Customer churn prediction in the telecommunication sector using a rough set approach. *Neurocomputing*, 237, 242-254.
- Ba, L. J., Kiros, J. R., & Hinton, G. E. (2016). *Layer Normalization*. arXiv preprint:1607.06450. <https://doi.org/10.48550/arXiv.1607.06450>
- Çelik, O. & Osmanoglu, U. O. (2019). Comparing to techniques used in customer churn analysis. *Journal of Multidisciplinary Developments*, 4(1), 30-38.
- De Caigny, A., Coussement, K. & De Bock, K. W. (2018). A new hybrid classification algorithm for customer churn prediction based on logistic regression and decision trees. *European Journal of Operational Research*, 269(2), 760-772.
- Domingos, E., Ojeme, B. & Daramola, O. (2021). Experimental analysis of hyperparameters for deep learning-based churn prediction in the banking sector. *Computation*, 9(3), 34.
- Halibas, A. S., Matthew, A. C., Pillai, I. G., Reazol, J. H., Delvo, E. G. & Reazol, L. B. (2019). Determining the Intervening Effects of Exploratory Data Analysis and Feature Engineering in Telecoms Customer Churn Modelling. *2019 4th MEC International Conference on Big Data and Smart City (ICBDSC)*, 1-7.
- He, K., Zhang, X., Ren, S. & Sun, J. (2016). Deep Residual Learning for Image Recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- Jain, H., Yadav, G. & Manoov, R. (2020). Churn prediction and retention in banking, telecom and IT sectors using machine learning techniques. In *Advances in Machine Learning and Computational Intelligence: Proceedings of ICMLCI 2019* (pp. 137-156). Singapore: Springer Singapore.
- Karvana, K. G. M., Yazid, S., Syalim, A. & Mursanto, P. (2019, October). Customer churn analysis and prediction using data mining models in banking industry. In *2019 international workshop on big data and information security (IWBIS)* (pp. 33-38). IEEE.
- Khan, Y., Shafiq, S., Naeem, A., Ahmed, S., Safwan, N. & Hussain, S. (2019). Customers churn prediction using artificial neural networks (ANN) in telecom industry. *International journal of advanced computer science and applications*, 10(9).
- Kingma, D. P. & Ba, J. (2015). Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015*, Conference Track Proceedings.
- Lalwani, P., Mishra, M. K., Chadha, J. S. & Sethi, P. (2022). Customer churn prediction system: a machine learning approach. *Computing*, 104, 271-294.
- Li, Y., Hou, B., Wu, Y., Zhao, D., Xie, A. & Zou, P. (2021). Giant fight: Customer churn prediction in traditional broadcast industry. *Journal of Business Research*, 131, 630-639.
- Pustokhina, I. V., Pustokhin, D. A., Nguyen, P. T., Elhoseny, M. & Shankar, K. (2023). Multi-objective rain optimization algorithm with WELM model for customer churn prediction in telecommunication sector. *Complex & Intelligent Systems*, 9, 3473–3485.

- Raeesi, S., Bajalan, S. & Fallahpour, S. (2024). Bank's Credit Portfolio Optimization Using Actuarial Approach and Artificial Neural Networks. *Financial Research Journal*, 26(3), 710-733. (in Persian)
- Rahimi, S., Rousta, A. & Asayesh, F. (2024). Evaluating the Relationship between Factors Enhancing the Competitiveness of Customer Foreign Currency Services in the Banking Industry. *Financial Research Journal*, 26(2), 424-446. (in Persian)
- Rahman, M. & Kumar, V. (2020). Machine Learning Based Customer Churn Prediction in Banking. *2020 4th International Conference on Electronics, Communication and Aerospace Technology (ICECA)*, 1196–1201.
- Sayed, H., Abdel-Fattah, M. A. & Kholief, S. (2018). Predicting potential banking customer churn using apache spark ML and MLlib packages: a comparative study. *International Journal of Advanced Computer Science and Applications*, 9(11).
- Spanoudes, P. & Nguyen, T. (2017). *Deep Learning in Customer Churn Prediction: Unsupervised Feature Learning on Abstract Company Independent Feature Vectors*. <https://arxiv.org/abs/1703.03869>
- Umayaparvathi, V., & Iyakutti, K. (2017). Automated feature selection and churn prediction using deep learning models. *International Research Journal of Engineering and Technology (IRJET)*, 4(3), 1846-1854.
- Vafeiadis, T., Diamantaras, K. I., Sarigiannidis, G., & Chatzisavvas, K. C. (2015). A comparison of machine learning techniques for customer churn prediction. *Simulation Modelling Practice and Theory*, 55, 1–9.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*, pp. 5998–6008.
- Vijaya, J., & Sivasankar, E. (2019). An efficient system for customer churn prediction through particle swarm optimization based feature selection model with simulated annealing. *Cluster Computing*, 22, 10757-10768.