

## پیاپی سازی و مقایسه روش های یادگیری ماشین در ارزیابی ریسک اعتباری مشتریان موسسات اعتباری و مالی

زهرا فتحی اقدم<sup>۱</sup>، محمدرضا رسولی<sup>۲</sup>

<sup>۱</sup> کارشناسی، دانشکده مهندسی صنایع، دانشگاه علم و صنعت ایران، تهران، ایران؛ zahraf.7811@gmail.com

<sup>۲</sup> استادیار، گروه مهندسی سیستم های هوشمند، دانشکده مهندسی صنایع، دانشگاه علم و صنعت ایران، تهران، ایران؛ rasouli@iust.ac.ir

### چکیده

عمده ترین ریسک در موسسات مالی، ریسک اعتباری<sup>۱</sup> است که عدم توانایی قرض گیرنده در بازپرداخت تسهیلات و وام در موعد مقرر معنا می شود. امروزه کارشناسان موسسات می توانند با تحلیل داده های مشتریان و استفاده از یادگیری ماشین، برای تخصیص اعتبارات به آنها تصمیم گیری نمایند. بنابراین، اعطا و عدم اعطای وام، بدون قضاوت شخصی و بر مبنای ریاضیات انجام می گیرد. در تحقیق پیشرو، داده های مشتریان موسسه لندینگ کلاب<sup>۲</sup> از وبسایت کگل<sup>۳</sup> جمع آوری شده است و الگوریتم های مختلف یادگیری ماشین برای پیش بینی ریسک اعتباری بررسی شده است. چالش های مهم این پژوهش، انتخاب ویژگی ها و مدلسازی بر کلاس های نامتوازن بود. فلذا، نوآوری کلیدی این پژوهش استفاده از روش های انتخاب ویژگی شامل آزمون تحلیل واریانس<sup>۴</sup>، آزمون اطلاعات متقابل<sup>۵</sup> و روش حذف بازگشتی ویژگی<sup>۶</sup> و ترکیب آنها با الگوریتم های یادگیری ماشین برای پیش بینی ریسک اعتباری می باشد. همچنین برای مدلسازی بر کلاس های نامتوازن<sup>۷</sup>، الگوریتم های حساس به هزینه<sup>۸</sup> مانند درخت تصمیم<sup>۹</sup>، ماشین بردار پشتیبان<sup>۱۰</sup> و الگوریتم های جمعی<sup>۱۱</sup> مثل جنگل تصادفی<sup>۱۲</sup> و تقویت گرادیان سبک<sup>۱۳</sup> استفاده شد. دیگر نوآوری پژوهش استفاده از الگوریتم های جدیدتر حوزه یادگیری ماشین مانند تقویت گرادیان سبک می باشد و نتایج نشان می دهد بهترین طبقه بندی توسط همین الگوریتم و بر ویژگی های منتخب از روش آزمون تحلیل واریانس انجام گرفته است. مساحت زیر منحنی مشخصه عملکرد<sup>۱۴</sup> این مدل برابر با ۰.۷۲۴ و معیار اف ۰.۵ برابر با ۰.۷۸۸ است.

کلمات کلیدی: ریسک اعتباری، طبقه بندی، درخت تصمیم، ماشین بردار پشتیبان، جنگل تصادفی، تقویت گرادیان سبک

## Implementation And Comparison Between Machine Learning Techniques For Assessing Credit Risk Of Customers Of Financial Institutions

Zahra Fathi Aghdam, Mohammadreza Rasouli

Bachelor, Faculty Of Industrial Engineering, Iran University Of Science And Technology, Tehran, Iran

Assistant Professor, Faculty Of Industrial Engineering, Iran University Of Science And Technology, Tehran, Iran

### ABSTRACT

Profitability in financial institutions and banks always depends on granting loans and facilities. Therefore, one of the major risks in these organizations is credit risk, which means the borrower's inability to repay the loan on time. Data mining is one of the most efficient methods in credit risk management as helps relevant experts in institutions to analyze customer data and make appropriate decisions. As a result, granting or not granting loans is done without personal judgment and based on mathematics and statistics and using intelligent systems. In this research, various machine learning algorithms have been developed on Lending Club dataset, collected from Kaggle website, to assess credit risk. Feature selection and modeling on imbalanced classes were considerable challenges in this research. Therefore, innovation of the research is utilizing three different methods of feature selection named analysis of variance (ANOVA), mutual information and Recursive Feature Elimination (RFE) and combining them with machine learning algorithm to predict credit risk. On the other hand, cost-sensitive algorithms such as decision tree (DT) and support vector machine (SVM) and ensemble methods as random forest (RF) and light gradient boosting (LGBM) were implemented to face imbalanced classes. In spite of other research, newer algorithms like LGBM has been implemented to be consider as an innovation. Finally, the best classifier was LGBM built on features selected by ANOVA which resulted ROC AUC and F0.5 equal to 0.724 and 0.788, respectively.

**Keywords:** Credit Risk, Classification, Decision Tree, Support Vector Machine, Random Forest, Light Gradient Boosting

## ۱- مقدمه

سودآوری در موسسات مالی و بانک‌ها همواره به اعطای تسهیلات وابسته است. بنابراین، وام‌ها مهم‌ترین منبع درآمد بانک‌ها و موسسات اعتباری و همچنین بزرگترین ریسک برای این سازمان‌ها هستند [۵]. یکی از مهم‌ترین ریسک‌ها در این موسسات ریسک اعتباری نام دارد که عدم توانایی قرض‌گیرنده در بازپرداخت وام در موعد مقرر معنا می‌شود [۶]. در گذشته، موسسات اعتباری و بانک‌ها افرادی را برای ارزیابی پیشینه متقاضیان و تصمیم‌گیری در راستای اعطا و یا عدم اعطای تسهیلات استخدام می‌کردند [۷] ولی همچنان حجم تسهیلات اعطایی سوخت‌شده و معوقه بانک‌ها زیاد بود. این موضوع عدم وجود روش‌های مناسب اعتبارسنجی و سیستم‌های مدیریت ریسک در شبکه مالی را نشان می‌دهد [۱] که ضرورت بهره‌گیری از روش‌های نوین را ملموس می‌سازد. امروزه با توسعه بانکداری الکترونیکی روزانه داده‌های زیادی تولید می‌شود فلذا این موسسات می‌توانند بدون قضاوت شخصی، با تحلیل داده‌های مشتریان و بر مبنای علم ریاضیات و آمار و با بهره‌گیری از سیستم‌های هوشمند در خصوص اعطای تسهیلات به آن‌ها تصمیم‌گیری نمایند تا سرمایه سازمان از دست نرود و منابع به طور بهینه به وام‌های سودآور اختصاص یابد. در واقع، از دیدگاه یادگیری ماشین می‌توان این مسئله را با بهره‌گیری از الگوریتم‌های طبقه‌بندی دو کلاسه مدلسازی نمود [۸] که از اهداف این پژوهش است. علاوه بر این، شناسایی شاخص‌های مؤثر بر این ریسک نیز به موضوع قابل توجهی تبدیل شده‌است و موسسات اعتباری و بانک‌ها باید با توجه به پیچیدگی فعالیت‌ها و محیط اقتصادی پیرامونشان، معیارهای مناسبی برای ارزیابی ریسک اعتباری مشتریان انتخاب کنند [۲]. از این رو در این پژوهش به دنبال شناسایی تاثیرگذارترین ویژگی‌ها در قدرت بازپرداخت مشتریان خواهیم بود.

مروزی بر تاریخچه ریسک اعتباری نشان می‌دهد که در دهه ۱۶۱۰ چندین بانک اروپایی و آمریکایی ورشکست شدند که علت اصلی آن، اعطای وام‌های با ریسک بالا تحت بحران‌های اقتصادی آن دوران معرفی شد. از این رو برای نخستین بار بیور در سال ۱۶۹۱ مدل "رگرسیون لجستیک" چند متغیره را برای تعیین ورشکستگی شرکت‌ها ارائه کرد [۳]. در پی آن مفاهیم و ایده‌هایی جدید برای تحلیل ریسک اعتباری در سال ۱۹۴۱ ظهور کرد [۹]. با گذر زمان ریسک اعتباری به یک مسأله‌ی مهم داده‌کاوی در حوزه‌ی مالی تبدیل شده و این مهم در سال ۱۹۵۰ معرفی شد [۱۰]. متدولوژی‌های متنوعی برای حل مسأله‌ی طبقه‌بندی ریسک اعتباری وجود دارد. این روش‌ها رگرسیون لجستیک، تحلیل نزدیک‌ترین همسایگی، شبکه بیزین، شبکه عصبی مصنوعی، درخت‌های تصمیم، جنگل تصادفی، الگوریتم ژنتیک، روش‌های تصمیم‌گیری چند معیاره، ماشین بردار پشتیبان و انواع دیگری را شامل می‌شود [۱۱]. جدول زیر الگوریتم‌های به کار گرفته شده در مقالات سال‌های اخیر را نشان می‌دهد.

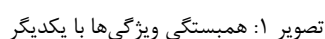
جدول ۱: الگوریتم‌های یادگیری ماشین مورد استفاده در منابع اخیر

شماره مرجع	سال انتشار	رگرسیون لجستیک	ماشین بردار پشتیبان	درخت تصمیم	جنگل تصادفی	تقویت گرادیان	تقویت گرادیان شدید <sup>۱۵</sup>	تقویت انطباقی <sup>۱۶</sup>	شبکه عصبی
۱۳	۲۰۱۷	*	*	*					*
۱۴	۲۰۱۸	*			*	*			
۱۵	۲۰۱۹	*	*	*	*				
۱۶	۲۰۱۹						*		
۱۷	۲۰۲۰	*	*	*	*			*	*
۱۸	۲۰۲۰	*		*	*	*		*	*
۱۹	۲۰۲۱	*			*				*
۲۰	۲۰۲۲		*	*	*	*	*	*	
۲۱	۲۰۲۲	*	*	*	*			*	

بر اساس آنچه در در جدول ۱ مشاهده می‌شود، در این پژوهش نیز از سه الگوریتم درخت تصمیم، جنگل تصادفی و ماشین بردار پشتیبان به عنوان پرتکرارترین الگوریتم‌ها استفاده شده‌است. همچنین مدلسازی با الگوریتم تقویت گرادیان سبک نیز انجام گرفت زیرا این الگوریتم مناسب برای مدلسازی بر روی داده‌های حجم بالا است و در مقالات اخیر نیز استفاده نشده‌است.

## ۲- بیان مسأله

در حال حاضر به دلیل حجم بالای تسهیلات، ریسک باز پرداخت وام‌ها یک چالش بزرگ برای موسسات اعتباری و بانک‌ها می‌باشد فلذا استفاده از روش‌های اعتبارسنجی و در واقع ارزیابی ریسک اعطای تسهیلات به متقاضیان یکی از راه‌های کارآمد جهت مدیریت ریسک می‌باشد [۱۲]. علاوه بر این، وام‌های سوخت شده (وصول نشدن) یکی از مسائل قابل توجه است که از طریق استفاده از مدل‌های اعتباری دقیق‌تر می‌توان تا حدودی بر این مسأله فائق آمد [۴]. هدف از این پژوهش پیاده‌سازی و مقایسه روش‌های مختلف یادگیری ماشین جهت تحلیل و پیش‌بینی ریسک اعتباری مشتریان و طبقه‌بندی آن‌ها است.



### ۳-۳- پیش‌پردازش داده

پیش‌پردازش داده‌ها به جهت بهبود کیفیت داده‌های واقعی برای داده‌کاوی لازم است. بنابراین در ابتدا به جهت نمونه‌گیری از داده‌ها، بر اطلاعات سال ۲۰۱۸ تمرکز شده و همچنین ویژگی هدف به دو حالت "صور در بازپرداخت" و "عدم صور در بازپرداخت" و یا به طور خلاصه وام‌های "خوب" و "بد" خلاصه شده‌است. وام‌های بدون صور ۷۱ درصد حجم موجودیت‌ها را تشکیل می‌دهند که نشان‌دهنده نامتوازن بودن کلاس‌ها می‌باشد. جهت مدیریت داده‌های مفقود، در ویژگی میزان سابقه کار متقاضی، داده‌های مفقود با مقدار صفر جایگزین شده و در سایر ویژگی‌ها با حداکثر مقدار در آن ویژگی جایگزین شده‌اند. انتخاب ویژگی‌ها از اصلی‌ترین چالش‌ها در این پژوهش است فلذا در ابتدا ویژگی‌هایی که بیش از ۲۰ درصد موجودیت‌ها در آن‌ها مفقود هستند حذف شده و همچنین اطلاعات مربوط به تعدادی از ویژگی‌ها در زمان تصمیم‌گیری اولیه موجود نیستند و پس از اعطای وام ایجاد می‌شوند، بنابراین این دسته ویژگی‌ها نیز حذف شده‌اند تا از نشتی اطلاعات<sup>۲۰</sup> و بیش‌برازش جلوگیری شود. در ادامه ویژگی‌های اسمی که تنوع بسیار بالایی دارند و ویژگی‌های عددی که تنها یک مقدار ثابت دارند یا واریانس آن‌ها بسیار کم است نادیده گرفته شده‌اند. علاوه بر آن، گام شناخت داده نشان می‌دهد که نرخ وام‌های بد در مقادیر مختلفی از بعضی ویژگی‌ها مقدار نسبتاً ثابتی است و می‌توان نتیجه گرفت آن ویژگی‌ها تاثیر ویژه‌ای بر ویژگی هدف ندارند پس می‌توان آن‌ها را حذف نمود. نهایتاً ۵۸۰۳۸ موجودیت و ۵۸ ویژگی باقی ماند که با استفاده از سه روش آزمون تحلیل واریانس، آزمون اطلاعات متقابل و حذف بازگشتی ویژگی، ویژگی‌های موثرتر از میان آن ۵۸ ویژگی انتخاب شده‌است. در گام بعدی مدل‌سازی بر روی ویژگی‌های منتخب هر روش به تفکیک شرح داده خواهد شد.

در انتها و پیش از مدل‌سازی و ارائه طبقه‌بندها، ۷۰ درصد داده‌ها برای آموزش و ارزیابی از طریق روش کراس‌ولیدشن<sup>۲۱</sup> جداسازی شده و ۳۰ درصد باقی مانده به عنوان داده‌ی آزمون و نماینده‌ی دنیای واقعی در نظر گرفته شده‌است. همچنین داده‌های عددی نرمال‌سازی شده و داده‌های اسمی به داده‌های موهومی<sup>۲۲</sup> تبدیل شده‌اند.

### ۳-۴- ایجاد و پیاده‌سازی طبقه‌بندها

در این مرحله به انتخاب الگوریتم مدل‌سازی و اجرای آن پرداخته می‌شود. در ادامه‌ی مراحل قبلی، طبقه‌بندی کلاس‌های نامتوازن دومین چالش در این پژوهش است. دو روش موثر بر این چالش، طبقه‌بندی با الگوریتم‌های حساس به هزینه و همچنین الگوریتم‌های جمعی است. بنابراین دو الگوریتم درخت تصمیم و ماشین بردار پشتیبان به عنوان طبقه‌بندهای حساس به هزینه و همچنین جنگل تصادفی و تقویت گرادیان سبک به عنوان نماینده‌ی الگوریتم‌های جمعی مورد استفاده قرار گرفته‌اند. همانطور که پیش‌تر اشاره شد از سه روش برای انتخاب ویژگی استفاده شده‌است و حال بر ویژگی‌های منتخب از هر روش، هر چهار الگوریتم مورد بحث پیاده‌سازی شده‌اند و در مجموع ۱۲ طبقه‌بند حاصل از حالات مختلف انتخاب ویژگی و الگوریتم پیش‌بینی موجود است.

مدل‌سازی با مقادیر پیش فرض در هر الگوریتم موجب بیش‌برازش شد فلذا تنظیم پارامترهای هر الگوریتم انجام گرفت تا از پیچیدگی بیش از حد مدل‌ها جلوگیری شود. از روش گریدسرچ‌سی‌وی<sup>۲۳</sup> تحت حالت ۵ فولد برای انتخاب بهترین مقادیر پارامترها استفاده شده‌است. همچنین پارامترها به گونه‌ای انتخاب شده‌اند تا مدلی پایدار حاصل شود و نه مدلی با مقدار بسیار بالا برای معیار صحت<sup>۲۴</sup> و بسیار پایین برای معیار حساسیت<sup>۲۵</sup> و بالعکس. زیرا بالا بودن هر دوی آن‌ها مطلوب مسأله است. جدول ۲ پارامترهای بهینه معرفی شده توسط گریدسرچ‌سی‌وی در هر الگوریتم را بیان می‌کند. نهایتاً نتایج مدل‌سازی با پارامترهای منتخب در بخش بعدی مطرح خواهد شد.

جدول ۲: پارامترهای منتخب برای هر الگوریتم حاصل از روش گریدسرچ‌سی‌وی

الگوریتم	درخت تصمیم	ماشین بردار پشتیبان	جنگل تصادفی	تقویت گرادیان سبک
آزمون تحلیل واریانس	عمق درخت = ۲ حداقل تعداد نمونه برای اشعاب = ۲ معیار: آنتروپی	سی = ۱ گاما = ۰,۱	عمق درخت‌ها = ۸ حداقل تعداد نمونه برای اشعاب = ۸ تعداد درخت‌ها = ۲۰	عمق درخت‌ها = ۶ نرخ یادگیری = ۰,۱
آزمون اطلاعات متقابل	عمق درخت = ۲ حداقل تعداد نمونه برای اشعاب = ۲ معیار: آنتروپی	سی = ۱ گاما = ۰,۱	عمق درخت‌ها = ۸ حداقل تعداد نمونه برای اشعاب = ۸ تعداد درخت‌ها = ۲۵	عمق درخت‌ها = ۶ نرخ یادگیری = ۰,۱
حذف بازگشتی ویژگی	عمق درخت = ۲ حداقل تعداد نمونه برای اشعاب = ۲ معیار: آنتروپی	سی = ۱ گاما = ۰,۱	عمق درخت‌ها = ۸ حداقل تعداد نمونه برای اشعاب = ۲ تعداد درخت‌ها = ۲۰	عمق درخت‌ها = ۶ نرخ یادگیری = ۰,۰۷

### ۳-۵- ارزیابی

معیارهای متنوعی برای انتخاب بهترین طبقه‌بند وجود دارد. از آنجایی که در این پژوهش کلاس‌ها نامتوازن هستند، معیار دقت مناسب نیست حال آن که دو معیار امتیاز اف و مساحت زیر نمودار منحنی مشخصه عملکرد توانایی و عملکرد مدل را به خوبی بازتاب می‌دهند.

همانطور که پیش تر اشاره شد بالا بودن و نطدیک بودن مقدرا هر دو معیار صحت و حساسیت به طور همزمان مطلوب مسئله است ولی در نگاهی دقیق تر مثبت کاذب از اهمیت نسبی بیشتری نسبت به منفی کاذب برخوردار است، زیرا از اعطای وام های مشکوک به قصور در بازپرداخت جلوگیری می کند و مانع از دست رفتن سرمایه سازمان می شود حال آن که زیاد بودن منفی کاذب سبب فرصت از دست رفته برای سازمان است. بنابراین امتیاز اف ۰.۵، مناسب تر از اف ۱ است زیرا بزرگ بودن صحت را بر حساسیت اولویت می دهد. عبارت ۱ معیار اف ۰.۵ را معرفی می کند.

$$(1) \frac{\text{حساسیت} * \text{صحت} * (1+0.5^2)}{\text{حساسیت} + (\text{صحت} * 0.5^2)}$$

جدول ۳: نتایج حاصل از طبقه بندی بر ویژگی های استخراجی از سه روش مورد بحث، (بهترین طبقه بند، پررنگ نگارش شده است)

روش انتخاب ویژگی	الگوریتم یادگیری ماشین	معیار اف ۰.۵	AUC-ROC
آزمون تحلیل واریانس	درخت تصمیم	۰.۷۵۷۲	۰.۶۱۵۹
	ماشین بردار پشتیبان	۰.۷۷۸۷	۰.۶۹۸۲
	جنگل تصادفی	۰.۷۸۰۴	۰.۷۰۶۷
	تقویت گرادیان سبک	۰.۷۸۸۷	۰.۷۲۴۶
آزمون اطلاعات متقابل	درخت تصمیم	۰.۷۵۷۲	۰.۶۱۵۹
	ماشین بردار پشتیبان	۰.۷۷۶۱	۰.۶۹۱۴
	جنگل تصادفی	۰.۷۷۸۲	۰.۷۰۲۳
	تقویت گرادیان سبک	۰.۷۸۴۴	۰.۷۱۷۷
حذف بازگشتی ویژگی	درخت تصمیم	۰.۷۵۷۲	۰.۶۱۵۹
	ماشین بردار پشتیبان	۰.۷۷۵۹	۰.۶۸۹۹
	جنگل تصادفی	۰.۷۷۸۹	۰.۷۰۴۷
	تقویت گرادیان سبک	۰.۷۸۸۵	۰.۷۲۳۹

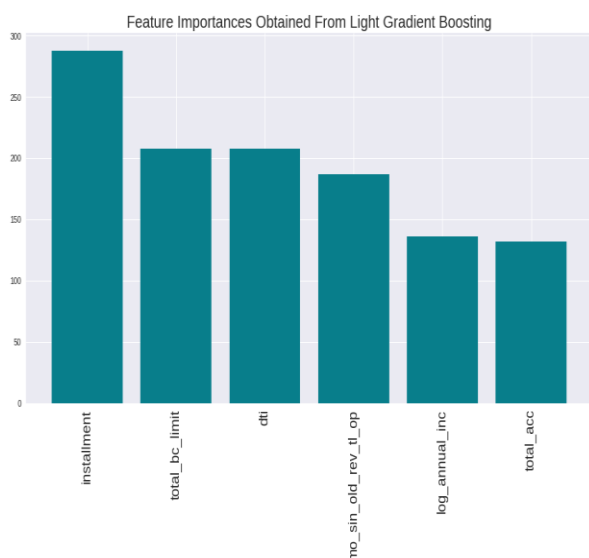
بر اساس نتایج، به طور کلی هر چهار الگوریتم زمانی عملکرد بهتری دارند که ویژگی ها با روش آزمون تحلیل واریانس انتخاب شده اند و نتایج مدلسازی تحت این حالت، نشان می دهد که الگوریتم تقویت گرادیان سبک عملکرد بهتری نسبت به سایرین دارد.

تصویر ۲ مقادیر امتیاز اف ۰.۵ را در ترکیب های مختلف الگوریتم ها و روش های انتخاب ویژگی با یکدیگر مقایسه می کند. تصویر ۳ منحنی های مشخصه عملکرد حاصل از مدلسازی با چهار الگوریتم بر روی ویژگی های منتخب از روش آزمون تحلیل واریانس را نشان می دهد. منحنی تقویت گرادیان سبک بهترین عملکرد و بیشترین مساحت زیر نمودار را نسبت به سایرین دارد و نسبت به خط مبنا بهبود خوبی یافته است. همچنین تصویر ۴ تاثیرگذارترین ویژگی ها را در عملکرد متقاضیان برای بازپرداخت وام به نمایش می گذارد. این نمودار نشان می دهد هر سه دسته ویژگی اشاره شده در قبل، تاثیر به سزایی بر این موضوع دارند. میزان اقساط وام از جمله ویژگی های مربوط به وام، درآمد شخص از جمله ویژگی های مربوط به شخص و نهایتا چهار ویژگی دیگر از سوابق اعتباری فرد استخراج می شوند.

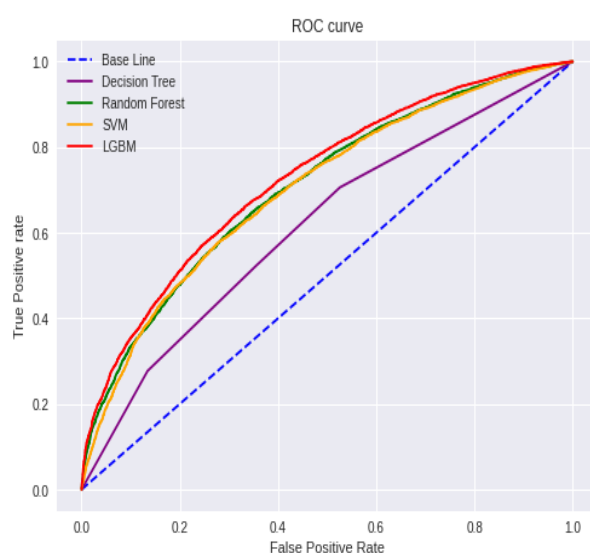


تصویر ۲: مقادیر اف ۰.۵ حاصل از چهار روش طبقه بندی بر روی ویژگی های منتخب استخراجی از روش های انتخاب ویژگی مذکور





تصویر ۴: مهم‌ترین ویژگی‌ها در مدل حاصل از ترکیب آزمون تحلیل واریانس و الگوریتم تقویت گرادیان سبک



تصویر ۳: منحنی مشخصه عملکرد چهار طبقه‌بند بر ویژگی‌های منتخب از روش آزمون تحلیل واریانس

### ۳-۶- توسعه

لیندینگ کلاب می‌تواند این برترین طبقه‌بند را به عنوان ابزاری جهت تصمیم‌گیری برای اعطای تسهیلات به متقاضیان مورد استفاده قرار دهد تا عملکرد کارشناسان مربوطه بهبود و خطای آن‌ها کاهش دهد. در نتیجه با مدیریت درست ریسک اعتباری ضررهای سازمان کاهش میابد.

### ۴- جمع‌بندی و نتیجه‌گیری

بر اساس مروری بر منابع، داده‌کاوی و پیاده‌سازی الگوریتم‌های یادگیری ماشین سبب مدیریت ریسک اعتباری در موسسات اعتباری و بانک‌ها می‌شود. این پژوهش نیز با مطالعه موردی بر داده‌های لیندینگ کلاب نشان داد با تکیه بر متدولوژی کریسپ و توسعه مدل‌ها می‌توان در خصوص متقاضیان هوشمندانه تصمیم‌گیری نمود و ریسک اعتباری را کاهش داد. عملکرد خوب این مدل می‌تواند ضررهای ناشی از اعطای وام‌های مشکوک به قصور در بازپرداخت را در این موسسه اعتباری کاهش دهد. همچنین فرصت‌های از دست رفته به شکل عدم اعطای وام‌های بدون قصور در بازپرداخت را تعدیل نماید.

نوآوری کلیدی در این پژوهش، بر استفاده از روش‌های انتخاب ویژگی و ترکیب آن‌ها با الگوریتم‌های مختلف یادگیری ماشین به جهت مدیریت ریسک اعتباری متمرکز می‌باشد. سه روش آزمون تحلیل واریانس، آزمون اطلاعات متقابل و روش حذف بازگشتی ویژگی برای انتخاب ویژگی‌ها مورد استفاده قرار گرفتند و پیش‌بینی ریسک اعتباری بر روی خروجی‌های روش‌های مذکور به صورت تفکیکی انجام گرفت. همچنین سعی شده از الگوریتم‌های جدیدتر حوزه یادگیری ماشین مانند الگوریتم تقویت گرادیان سبک برای پیش‌بینی ریسک اعتباری استفاده شود.

آنچنان که پیش‌تر اشاره شد دو چالش مهم در این پژوهش انتخاب ویژگی‌ها و مدل‌سازی بر کلاس‌های نامتوازن بود. نتایج نشان می‌دهد بهترین روش برای انتخاب ویژگی در این مجموعه داده، آزمون تحلیل واریانس و بهترین روش برای مقابله با کلاس‌های نامتوازن استفاده از یک الگوریتم جمعی مثل تقویت گرادیان سبک است زیرا الگوریتم‌های پرتکرار در این حوزه در این مطالعه موردی آنچنان موفق نبودند. در نهایت پیشنهاد می‌شود در پژوهش‌های آتی از دیگر الگوریتم‌های یادگیری ماشین و یا دیگر روش‌های انتخاب ویژگی استفاده شود. همچنین استفاده از روش‌های تولید موجودیت<sup>۲۶</sup> و یا کاهش موجودیت<sup>۲۷</sup> در راستای متوازن نمودن کلاس‌ها نیز می‌تواند مفید باشد.

### مراجع

- [۱] کیقبادی، ا. و خدای، ۱۳۹۲، داده‌کاوی صورت‌های مالی جهت اعطای تسهیلات مالی. پژوهش‌های حسابداری مالی و حسابرسی (پژوهشنامه حسابداری مالی و حسابرسی)، ۵(۱۷)، ۲۱۱-۱۷۹.
- [۲] عبدالمی، علی و فرزی زاده، محمد، ۱۳۹۶، ارزیابی مدل دسته‌بندی به کمک داده‌کاوی جهت امتیازدهی مشتریان بانک، چهارمین کنفرانس ملی فناوری اطلاعات، کامپیوتر و مخابرات، مشهد.
- [۳] صفری، زهره و خسروی، حمید، ۱۳۹۵، ارائه مدلی برای اعتبارسنجی مشتریان بانک‌ها با استفاده از الگوریتم‌های داده‌کاوی، دومین کنفرانس ملی رویکردهای نوین در مهندسی کامپیوتر و برق.
- [۴] طلوعی اشقی، عباس و مقدوری شریانی، فرناز و دانشگر، فرید، ۱۳۸۸، امتیاز دهی اعتباری متقاضیان کارتهای اعتباری بانکها با استفاده از تکنیک ماشین بردار پشتیبان، دومین کنفرانس بین المللی شهر الکترونیک، تهران.

- [Δ] Madaan, M., Kumar, A., Keshri, C., Jain, R., & Nagrath, P. (2021). Loan default prediction using decision trees and random forest: A comparative study. In IOP Conference Series: Materials Science and Engineering (Vol. 1022, No. 1, p. 012042).
- [϶] Aslam, U., Tariq Aziz, H. I., Sohail, A., & Batcha, N. K. (2019). An empirical study on loan default prediction models. *Journal of Computational and Theoretical Nanoscience*, 16(8), 3483-3488
- [Υ] Ahmed, M. I., & Rajaleximi, P. R. (2019). An empirical study on credit scoring and credit scorecard for financial institutions. *Int. Journal of Advanced Research in Computer Engineering & Technol. (IJARCET)*, 8, 275-9
- [Α] Zhou, L., & Wang, H. (2012). Loan default prediction on large imbalanced data using random forests. *TELKOMNIKA Indonesian Journal of Electrical Engineering*, 10(6), 1519-1525.
- [Θ] Durand, D. (1941). Risk elements in consumer installment financing. National Bureau of Economic Research, New York.
- [Ϸ] Luo, S., Kong, X., & Nie, T. (2016). Spline based survival model for credit risk modeling. *European Journal of Operational Research*, 253(3), 869-879.
- [ϸ] Zhang, Z., Gao, G., & Shi, Y. (2014). Credit risk evaluation using multi-criteria optimization classifier with kernel, fuzzification and penalty factors. *European Journal of Operational Research*, 237(1), 335-348.
- [Ϲ] Bellotti, T., & Crook, J. (2009). Support vector machines for credit scoring and discovery of significant features. *Expert systems with applications*, 36(2), 3302-3308.
- [Ϻ] Pandey, T. N., Jagadev, A. K., Mohapatra, S. K., & Dehuri, S. (2017, August). Credit risk analysis using machine learning classifiers. In 2017 International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS) (pp. 1850-1854). IEEE.
- [ϻ] Addo, P. M., Guegan, D., & Hassani, B. (2018). Credit risk analysis using machine and deep learning models. *Risks*, 6(2).
- [ϼ] Zhu, L., Qiu, D., Ergu, D., Ying, C., & Liu, K. (2019). A study on predicting loan default based on the random forest algorithm. *Procedia Computer Science*, 162, 503-513.
- [Ͻ] Qiu, W. (2019, July). Credit risk prediction in an imbalanced social lending environment based on XGBoost. In 2019 5th International Conference on Big Data and Information Analytics (BigDIA) (pp. 150-156). IEEE.
- [Ͽ] Tian, Z., Xiao, J., Feng, H., & Wei, Y. (2020). Credit risk assessment based on gradient boosting decision tree. *Procedia Computer Science*, 174, 150-160.
- [Ⓚ] Song, Y., Wang, Y., Ye, X., Wang, D., Yin, Y., & Wang, Y. (2020). Multi-view ensemble learning based on distance-to-model and adaptive clustering for imbalanced credit risk assessment in P2P lending. *Information Sciences*, 525, 182-204.
- [Ⓛ] Moscato, V., Picariello, A., & Sperl , G. (2021). A benchmark of machine learning approaches for credit score prediction. *Expert Systems with Applications*, 165, 113986.
- [Ⓜ] Dornigg, T. (2022). Credit risk modeling-predicting customer loan defaults with machine learning models (Doctoral dissertation).
- [Ⓝ] Laulajainen, M. (2022). Case Study on Utilizing Machine Learning in Corporate Default Risk Prediction: A practical Implementation to Credit Risk Management Process.

<sup>Ϸ</sup> Credit Risk

<sup>ϸ</sup> Lending Club

<sup>Ϲ</sup> [www.kaggle.com](http://www.kaggle.com)

<sup>Ϻ</sup> Analysis of Variance (ANOVA)

<sup>ϻ</sup> Mutual Information

<sup>ϼ</sup> Recursive Feature Elimination (RFE)

<sup>Ͻ</sup> Imbalanced Classes

<sup>Ͽ</sup> Cost-Sensitive

<sup>Ⓚ</sup> Decision Tree (DT)

<sup>Ⓛ</sup> Support Vector Machine (SVM)

<sup>Ⓛ</sup> Ensemble Models

<sup>Ⓛ</sup> Random Forest (RF)

<sup>Ⓛ</sup> Light Gradient Boosting (LGBM)

<sup>Ⓛ</sup> Area Under The Receiver Operating Characteristic (AUC-ROC)

<sup>Ⓛ</sup> Extreme Gradient Boosting (XGBoost)

<sup>Ⓛ</sup> Adaptive Boosting (Adaboost)

<sup>Ⓛ</sup> Crisp-DM (Cross-industry standard process for data)

<sup>Ⓛ</sup> SEMMA (Sample, Explore, Modify, Model, Assess)

<sup>Ⓛ</sup> KDD (Knowledge Discovery in Databases)

<sup>Ⓛ</sup> Data Leakage

<sup>Ⓛ</sup> Cross Validation

---

<sup>۲۲</sup> Dummy Variables

<sup>۲۳</sup> Grid Search CV

<sup>۲۴</sup> Precision

<sup>۲۵</sup> Recall

<sup>۲۶</sup> Oversampling

<sup>۲۷</sup> Undersampling