



Bingöl Üniversitesi
İktisadi ve İdari Bilimler Fakültesi Dergisi
Bingol University
Journal of Economics and Administrative Sciences

Cilt/Volume: 8, Sayı/Issue: 1
Yıl/Year: 2024, s. 91-105
DOI: 10.33399/biibfad.1391666
ISSN: 2651-3234/E-ISSN: 2651-3307
Bingöl/Türkiye

Makale Bilgisi /Article Info

Geliş/Received: 16/11/2023 Kabul/ Accepted: 25/01/2024



A Novel Hybrid Regression Model for Banking Loss Estimation

Bankacılık Zarar Tahmini için Yeni Bir Hibrit Regresyon Modeli

Pınar KARADAYI ATAŞ*

Abstract

Given the critical need to identify financial risks in the banking sector early, this study presents a novel approach that uses historical financial ratios from the FDIC database to predict bank failures in the United States. Accurate estimation of potential losses is essential for risk management and decision-making procedures. We present a novel hybrid approach to loss estimation in the context of bank failures in this study. ElasticNet regression and relevant data extraction techniques are combined in our method to improve prediction accuracy. We conducted thorough experiments and evaluated our hybrid approach's performance against that of conventional regression techniques. With a remarkably low Mean Squared Error (MSE) of 0.001, a significantly high R-squared value of 0.98, and an Explained Variance Score of 0.95, our proposed model demonstrates superior performance compared to existing methodologies. The accuracy of our method is further demonstrated by the Mean Absolute Error (MAE) of 1200 units. Our results highlight the potential of our hybrid approach to transform loss estimation in the banking and finance domain, offering superior predictive capabilities and more accurate loss estimations.

Keywords: Financial risk analysis, financial stability assessment, bank risk management, machine learning. **JEL Codes:** G21; G28; G32; G38; C13; C53; C58

Öz

Bankacılık sektöründeki finansal risklerin erken dönemde belirlenmesine yönelik kritik ihtiyaç göz önüne alındığında, bu çalışma, Amerika Birleşik Devletleri'ndeki banka başarısızlıklarını tahmin etmek için FDIC veri tabanındaki tarihsel finansal oranları kullanan yeni bir yaklaşım sunmaktadır. Potansiyel kayıpların tahmini önemlidir. Bu çalışmada banka iflasları bağlamında zarar tahminine yönelik yeni bir hibrit yaklaşım sunuyoruz. Tahmin doğruluğunu artırmak için ElasticNet regresyon ve ilgili veri çıkarma teknikleri önerdiğimiz yöntemde birleştirilmiştir. Önerilen hibrit yaklaşımın performansı kapsamlı deneyler yapılarak geleneksel regresyon tekniklerine göre değerlendirilmiştir. 0,001'lik son derece düşük Ortalama Kare Hatası (MSE), 0,98'lik oldukça yüksek R-kare değeri ve 0,95'lik açıklanan varyans skoru ile önerdiğimiz model mevcut metodolojilere kıyasla üstün performans sergilemektedir. Yöntemimizin doğruluğu 1200 birimlik ortalama mutlak hata (MAE) ile gösterilmektedir. Sonuçlarımız, üstün tahmin yetenekleri ve daha doğru kayıp tahminleri sunan, bankacılık ve finans alanında zarar tahminini dönüştürmeye yönelik hibrit yaklaşımımızın potansiyelini vurgulamaktadır.

Anahtar Kelimeler: Finansal risk analizi, finansal istikrar değerlendirmesi, banka risk yönetimi, makine öğrenimi

JEL Kodları: G21; G28; G32; G38; C13; C53; C58

* Dr., İstanbul Arel Üniversitesi, Mühendislik Fakültesi, Bilgisayar Mühendisliği, pinaratas@arel.edu.tr, ORCID: <https://orcid.org/0000-0002-9429-8463>

1. INTRODUCTION

The stability and well-being of banks are critical in the dynamic financial sector, not only for the economy as a whole but also for individual investors and stakeholders. The need for more precise and timely bank failure predictions has been highlighted by recent financial crises; however, this remains a challenging task in financial risk management. This study uses the abundance of information found in the FDIC's historical financial ratios to address this problem. Our study employs a binary classification method to forecast bank failures during a crucial 180-day window that was purposefully selected to coincide with the FDIC's quarterly reporting cycle. In addition to offering a new predictive model that improves the accuracy of bank failure predictions, this paper advances the discussion on risk assessment and financial stability in the banking industry. By carefully examining a large dataset with hundreds of features and training cases, this research aims to provide a robust instrument for proactive risk management and policy creation, strengthening the financial system's resilience to future crises.

With its wide range of applications, machine learning, a dynamic subfield of artificial intelligence, is completely changing the banking and finance industries. In addition to assisting decision-makers with crucial duties like credit scoring and loan approvals, it gives financial institutions the ability to recognize and prevent fraudulent transactions. Machine learning powers chatbots and financial robo-advisors in the customer service domain, providing personalized banking assistance. It provides sophisticated risk-return analyses to investors, supporting asset allocation systems. The domain also includes automated insurance services, providing policyholders more efficient experiences. Because machine learning can process large amounts of data and is skilled at handling complex, nonlinear data patterns, it has a significant impact on finance and is considered a key tool in the fields of statistics and financial analysis. Recent years have seen significant research by Ozbayoglu, Gudelek, and Sezer (2020) in applying computational intelligence to finance. This study focuses on reviewing the latest developments in Machine Learning across six key financial sectors: stock markets, portfolio management, forex markets, bankruptcy and insolvency, financial crises, and cryptocurrency.

The business, banking, and finance sectors have seen a large number of review articles on a variety of topics over the last thirty years. As Kumbure et al. (2022) have shown, many of these reviews focused on a single topic, such as stock market predictions. However, some research have adopted a more comprehensive strategy and have examined several financial domains. Solely, instead of concentrating only on machine learning, these more comprehensive reviews frequently address computational intelligence, or AI, in general. Aguilar-Rivera, Valenzuela-Rendón, and Rodríguez-Ortiz (2015), for instance, investigated the applications of genetic algorithms and Darwinian approaches in finance, while Pulakkazhy and Balan (2013) examined the use of data mining in banking. In study, Huang, Chai, and Cho (2020), conducted a thorough review of studies that focused specifically on deep learning. The study covered a wide range of topics, including credit risk prediction, macroeconomic forecasting, oil price forecasting, portfolio management, and stock trading. Similarly, Ozbayoglu et al. (2020) reviewed 144 articles that covered a wide range of topics, including financial text mining, behavioral finance, algorithmic trading, risk assessment, fraud detection, portfolio management, asset pricing, derivatives markets, cryptocurrency, blockchain, and financial sentiment analysis. These reviews demonstrate the wide-ranging and diverse uses of deep learning in the finance industry. The study Alzayed, Eskandari, and Yazdifar (2023) presents evidence that machine learning methods can reliably predict the

failure risk of banks in the European Union-27, based on data from the past decade. It highlights that factors such as earnings, capital adequacy, and management capability are the strongest predictors of bank failure. The research is particularly relevant in the context of economic uncertainties brought about by the COVID-19 pandemic.

The study Zou, Gao, and Gao (2022), demonstrates that machine learning methods can effectively predict bank failure risks in the EU-27, using past decade data. It identifies earnings, capital adequacy, and management capability as key predictors, emphasizing the research's relevance amid COVID-19-related economic uncertainties. The study Doumpos et al. (2023) provides a comprehensive bibliographic overview of research in the banking sector, with an emphasis on studies conducted in the last ten years that make use of artificial intelligence (AI) and operations research (OR). Numerous important subjects are covered, such as the impact of fintech, customer-related studies, banking regulation, mergers and acquisitions, risk assessment, bank performance, and bank efficiency. The survey offers a thorough summary of the ways in which OR and AI techniques have advanced these fields within the banking industry. Another study uses a new interpretability improvement of Extreme Gradient Boosting (XGBoost) to present a highly interpretable and precise machine learning model for business financial distress prediction. Carmona, Dwekat, and Mardawi (2022) examines data from 1,760 French businesses in 2018 – 1,585 of them successful, and 175 of them failing to find important signs of financial distress. Zhao et al. (2022) review aspects of bank erosion such as failure mechanisms, bank retreat rates, and modeling techniques. It mentions that the bank stability of tidal and river channels can be similarly impacted by various external forces. Reviewing data and empirical functions for bank retreat rates, the study emphasizes the importance of taking geotechnical and hydraulic factors into account. It offers suggestions to enhance current models and suggests a new hierarchy of modeling techniques based on time scales. The study, Veganzones, Séverin, and Chlibi (2023), examines to the relationship between variables used in earnings management and the prediction of corporate failures . It presents a brand-new threshold model technique that divides samples into various regimes according to a predetermined threshold variable. Next, the authors examine these regimes in order to evaluate how well earnings management variables predict corporate failures.

Hafeez et al. (2022) suggest building the z-score using a forward-looking methodology that incorporates analyst projections. They demonstrate empirically that this forward-looking z-score is able to accurately forecast changes in the standard z-score up to one quarter ahead of time. Ahmad et al. (2022) analyze time series data using a long short-term memory (LSTM) recurrent neural network architecture in the study's framework. The AirLab Failure and Anomaly flight dataset, a comprehensive collection of fault types in the control surfaces of fixed-wing autonomous aerial vehicles, is publicly accessible and has been subjected to its application. The research is motivated by the difficulty of comprehending intricate financial systems. Wang et al. (2022) presents a brand-new algorithm that combines neural networks, Gaussian process regression, and the Differential Evolution algorithm. This novel method offers fresh perspectives on the dynamics of financial modeling by identifying and forecasting parameters in a symmetric chaotic fractional financial model.

Many ensemble methods, such as Multiple Logistic Regression, Decision Trees, Random Forests, Gaussian Naive Bayes, Support Vector Machines, are explored in research Anand, Velu, and Whig (2022) in order to predict loan default. These approaches are presented in the study, Nazareth and Reddy (2023), in an effort to improve loan default prediction efficiency and accuracy.

In this study, we combine an innovative hybrid machine learning technique with regression methods. Given the extensive volume of data utilized in finance, it is critical to ensure that only high-quality data is used when making decisions. In order to achieve this, we have created a brand-new hybrid approach that carefully uses the best instances when making decisions, guaranteeing the model's successful construction.

In our study, the combination of these cutting-edge methods has shown to be very successful. Applying these techniques to the FDIC dataset has produced encouraging outcomes, indicating the effectiveness of our combined strategy in improving machine learning model performance. This study advances the field of machine learning and establishes a standard for subsequent research endeavors focused on enhancing data-driven decision-making procedures.

2. DATA AND METHODOLOGY

2.1 Description and Analysis of the Dataset

In this study, we aim to predict upcoming bank failures by analyzing historical financial ratios from the Federal Deposit Insurance Corporation (FDIC), which is a vital source of information for the American banking industry Heitz (2023). For insured banks and savings institutions across the country, the FDIC database offers extensive coverage of financial data encompassing a wide range of crucial information. This includes detailed quarterly "Call Reports" or Reports of Condition and Income, which are essential for understanding the current financial status of these organizations. Additionally, the FDIC's BankFind tool enables targeted searches on particular banks, providing details about their background, locations, and insurance status. This data is a valuable resource for our research because it is publicly available, allowing for a detailed examination and evaluation of bank performance, market trends, and general economic and financial health. Furthermore, our longitudinal analyses and trend evaluations heavily rely on the historical data found in the FDIC database. In addition to being essential for regulatory and policy-making purposes, this rich dataset is also a vital tool for our machine learning and predictive modeling endeavors, as noted by Le, Viviani, and Fauzi (2023), especially when it comes to predicting financial outcomes such as credit risk and bank failures.

The predictive model's emphasis on precision over recall is a crucial component of this study. This analysis is based on the financial ratios dataset, which includes 796 features and 656,438 training examples over the last 90 quarters. We specifically focus on failures that occur within a 180-day window, which was purposefully selected to coincide with the FDIC's quarterly bank call report data release schedule, which typically takes place 60 to 90 days following the end of a quarter. This method optimizes the data at hand to predict bank failures in a timely and accurate manner. The distribution of failed banks by state is shown in Figure 1, with the highest numbers found in Georgia (GA), Florida (FL), and Illinois (IL), suggesting regional concentrations of bank failures. Several states, however, report only one failure, highlighting notable regional differences. This implies that regulatory issues and local economic conditions could impact bank stability. The information emphasizes the necessity of region-specific financial oversight and might call for further research into the causes of these failures.

The average estimated financial loss from bank failures is displayed by state in Figure 2, with Illinois (IL) having the largest losses. Substantial variations in losses between states are shown in the chart, suggesting possible variations in bank sizes, risk management strategies,

or economic conditions. The ability to customize risk mitigation and regulatory strategies to the unique banking environments of each state relies on this information.

Figure 1: Explore the Distribution of Failed Banks by State

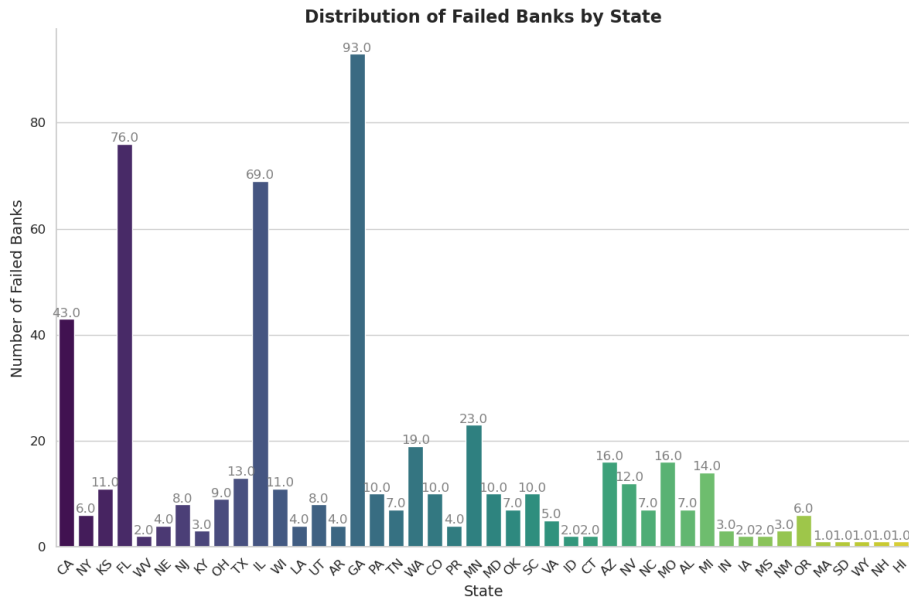
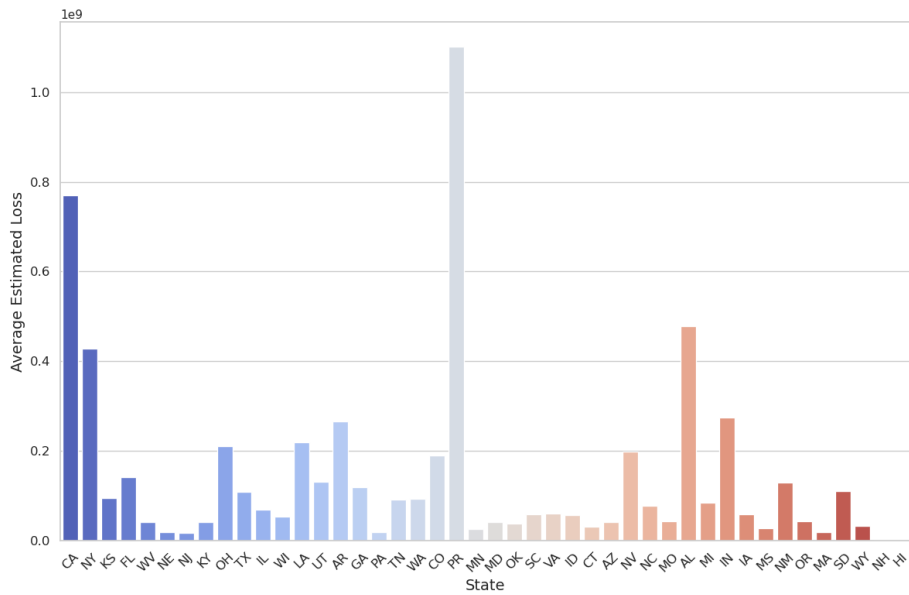
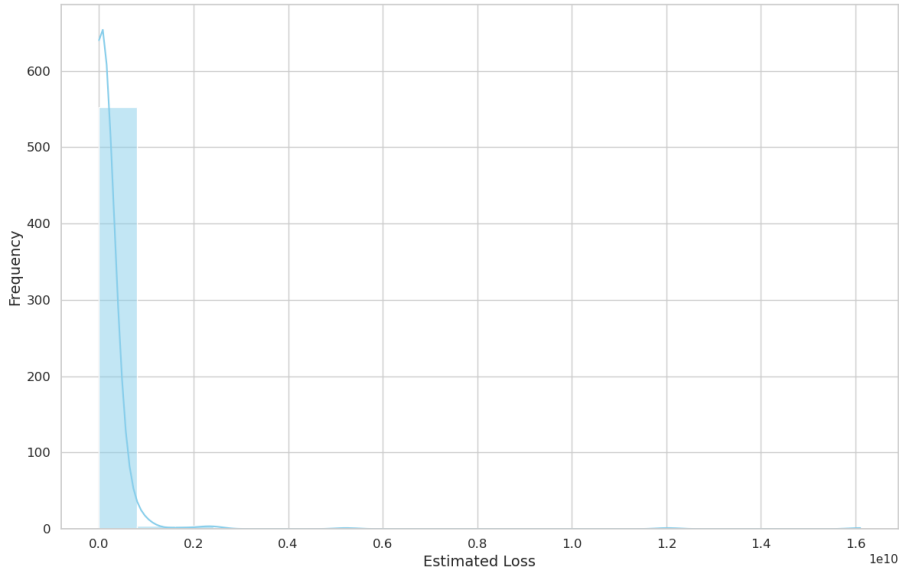


Figure 2: Explore the Distribution of Failed Banks by State



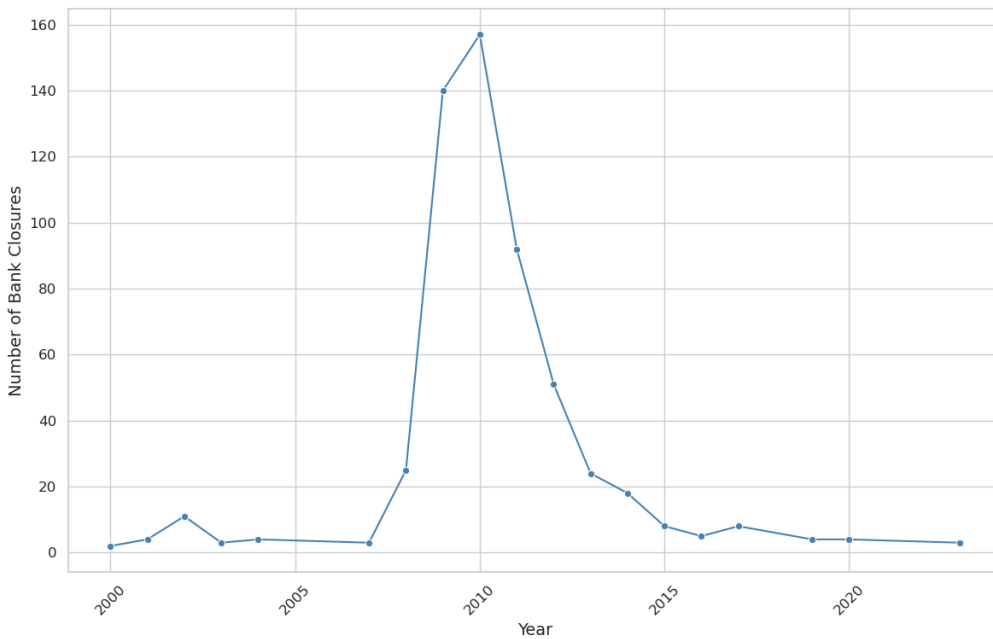
The distribution of estimated financial losses is represented by the histogram in Figure 3. Most observations cluster near the lower end of the loss spectrum, indicating a right-skewed distribution. This suggests that while the majority of losses are small, a small number of outliers have losses that are significantly larger. Large losses are less common, but when they do happen, they can still be significant, according to the long tail to the right. In financial data, where extreme events (such as major bank failures) are uncommon but can have serious financial ramifications, this distributional skewness is typical.

Figure 3: Analyze the Distribution of Estimated_Loss



The trend of bank closures over time is depicted in the line graph (Figure 4). A notable peak appears around 2010, suggesting a spike in bank closures during that period, which may have coincided with the fallout from the 2008 financial crisis. The number of closures sharply declines after this peak and then stabilizes at lower levels until the present. With fewer banks closing in recent years, this trend indicates that the banking industry has regained some stability following the initial shock of the crisis. The information shows how the financial sector has fared over the last 20 years and how the banking sector has been impacted by regulatory changes and economic cycles.

Figure 4: Explore the Trend of Bank Closures Over Time



2.2 Background Methods

Several regression techniques, such as Linear Regression (James et al., 2023; Montgomery, Peck, and Vining, 2021; Su, Yan, and Tsai, 2012), Random Forest Regression (Borup et al., 2023; Shoar, Chileshe, and Edwards, 2022), Gradient Boosting Regressor (Otchere et al., 2022; Sipper and Moore, 2022), ElasticNet (Friedman et al., 2023; Nasir et al., 2023) and Support Vector Regression (SVR) (Awad and Khanna, 2015; Zhang and O'Donnell, 2020) have been utilized in our study. These techniques serve as the foundation for our novel methodology. Through the strategic fusion of these regression techniques, we have developed a new methodology that is suited to the complexities of the banking and finance industry. Our approach leverages the distinct advantages of each method to efficiently address specific problems in the field. By using this all-encompassing strategy, we expect to greatly improve predictive performance and offer insightful analysis and practical solutions to stakeholders and financial institutions. The subsequent sections will provide an in-depth exploration of each method.

Linear regression is a basic statistical technique for modeling the relationship between a dependent variable and one or more independent variables. It looks for the best-fit line that minimizes the sum of squared errors, assuming a linear relationship. This method is useful for predicting continuous outcomes and acts as a foundation for more intricate regression models.

Random Forest Regression is an ensemble learning technique that combines several decision trees to increase prediction accuracy. During training, it builds a large number of decision trees, and during testing, it averages each tree's prediction. This method addresses complex interactions and non-linear relationships in the data while improving the robustness of the model.

Gradient Boosting is another ensemble technique that fits decision trees iteratively to the errors of the preceding ones, creating an additive model. By minimizing residual errors, it can effectively capture intricate patterns within the data. Although Gradient Boosting has a high predictive power, fine tuning may be necessary.

Combining L1 (Lasso) and L2 (Ridge) regularization techniques, ElasticNet is a regularized regression technique. By applying penalties to the coefficients, it is utilized to perform feature selection and reduce multicollinearity. ElasticNet works well with high-dimensional datasets because it balances model performance and variable selection.

Support Vector Regression (SVR) is a regression problem-specific extension of Support Vector Machines (SVM). SVR looks for the hyperplane with the smallest margin of error and the best fit to the data. Using kernel functions is especially useful for handling non-linear relationships. SVR performs well in high-dimensional spaces and is resilient in handling outliers.

2.3 The Proposed Hybrid Method

Datasets related to banking and finance are widely recognized for their extensive scope, encompassing a wide range of information. In this study, we utilize a dataset of this kind, which is brimming with data. However, not every piece of information is equally important when making decisions. Our goal in this study is to use a novel hybrid approach that will allow us to filter the dataset and identify the data points that are actually significant. By doing this, we hope to focus our learning efforts on this particular subset of data, which is crucial for making informed decisions. In this section, we will outline the step-by-step operation of the newly developed hybrid method.

Step 1: Support Vector Regression (SVR)

We analyzed the dataset using Support Vector Regression (SVR) as the first step of our hybrid method. Because SVR can handle nonlinear relationships in the data and produce accurate predictions, it was chosen as our regression technique. The main goal of this step was to make predictions based on the features of the dataset. The predictions generated by Support Vector Regression (SVR) are derived from its ability to identify the hyperplane that best represents the relationship between the target variable and the input features. The goal of this hyperplane is to maximize the model's accuracy while minimizing prediction error. After SVR finished analyzing our dataset, it produced predictions for each data instance.

We converted these predictions into probability values to improve their interpretability and better align them with our next steps. We were able to express the predictions as probabilities in the interval $[0, 1]$ thanks to this transformation, which improved their readability and made them suitable for further processing in our hybrid method. This first step set the stage for our later feature selection and data splitting procedures, which allowed us to successfully find important predictive features.

Step 2: Data Splitting Based on Probability

As the next stage of our hybrid approach, we divided the data according to the probability values that the SVR predictions produced. We aimed to achieve this by splitting the dataset into two separate groups, each with a distinct function in our predictive modeling procedure. The split was carried out by setting a threshold at the average probability value obtained from the SVR predictions. Information with probability values higher than the average was grouped into a different category than information that had probability values lower than the average. We were able to effectively separate data instances with strong predictive qualities from those with weaker predictive qualities thanks to this division.

This way of classifying the data allowed us to concentrate on the portion of the records that showed a greater chance of making accurate predictions. In order to make sure that we used the most pertinent and valuable data for our predictive modeling process, we strategically separated the dataset. This prepared the groundwork for our subsequent modeling steps.

Step 3: ElasticNet Regression

We introduced the ElasticNet regression method as part of our continuous effort to improve the predictive power of the dataset containing records with strong predictive qualities. ElasticNet was selected due to its exceptional ability to manage high-dimensional data efficiently and reduce the likelihood of overfitting, a prevalent issue in predictive modeling.

To further refine the dataset, we focused our analysis on records with higher prediction probabilities in this step. ElasticNet was the best option for this task because of its capacity to balance L1 and L2 regularization. Our goal in using ElasticNet was to extract useful relationships and insights from the data that would help us make more accurate predictions. This refinement process allowed us to concentrate our modeling efforts on the subset of data instances that showed strong predictive potential thanks to this refinement process, which also made sure that the most pertinent and instructive data was used to build our final predictive model.

Step 4: Refinement of Predictive Data

Building on the results of ElasticNet regression, we performed a second data refinement step focusing only on records with high probability values. Only data instances with strong and robust predictive capabilities were kept for additional analysis thanks to this selective approach. This refinement process, we were able to remove data that demonstrated less predictive qualities, thereby streamline the dataset. The predictive model's overall effectiveness was increased by this deliberate removal of superfluous data, which also helped to focus the model's attention on the most important and promising prediction-influencing variables. Essentially, by focusing our attention on the most important data points that were crucial to producing precise predictions, this stage of data refinement significantly improved the accuracy and efficacy of our predictive model.

Step 5: Model Testing

Extensive testing procedures were conducted to assess the predictive capabilities of our hybrid method, allowing for a thorough evaluation of its performance. Three different regression techniques were used in our evaluation process: Gradient Boosting, Random Forest Regression, and Linear Regression.

The primary aim of these assessments was to evaluate the precision, resilience, and efficacy of our hybrid model in forecasting bank failures within the designated timeframe. The refined dataset, which had been carefully selected to concentrate on records with strong predictive qualities, was subjected to each regression technique, enhancing its performance in the banking and finance domain.

Our goal in conducting this testing phase was to gather important information about how well our hybrid approach performed compared to conventional regression techniques. These evaluations would yield results that would make it evident how predictive the model is and how much more accurate the bank failure forecasts could be.

2.4. Model Evaluation

This section delves into the reasoning behind our choice of primary metrics for assessing the performance of the regression models, which include Mean Squared Error (MSE), R-squared, Explained Variance Score, and Mean Absolute Error (MAE) (Emmert-Streib and Dehmer, 2019:521; McAvaney et al., 2001:471).

Mean Squared Error (MSE): The Mean Squared Error measures the average squared difference between the expected and actual values. By determining the degree to which the predicted values agree with the true values, it measures the overall quality of a regression model.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \tag{1}$$

In a regression model, the R-squared, or coefficient of determination, represents the percentage of the dependent variable's variance that can be predicted from the independent variables. Higher values denote a better fit. The range is 0 to 1.

$$R^2 = 1 - \frac{SST}{SSR} \tag{2}$$

$$SSR = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \tag{3}$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 \tag{4}$$

In a regression model, the Explained Variance Score quantifies the percentage of the dependent variable's variance that can be accounted for by the independent variables. It offers a model performance evaluation that is comparable to that of R-squared.

$$\text{Explained Variance Score} = 1 - \frac{\text{var}(y-\hat{y})}{\text{var}(y)} \quad (5)$$

The average absolute difference between the expected and actual values is determined using the Mean Absolute Error formula. Compared to MSE, it offers a measure of error that is easier to understand.

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (6)$$

Hansen, Lunde, and Nason (2011) identify an important component of statistical analysis for comparing various models as the Model Confidence Set, or MCS, p-value. It measures the probability that a model's output differs significantly from other models in the set. A model's performance is statistically superior when its MCS p-value is low, indicating that it belongs to the "confidence set" of the best models. This aids in determining which models are the most trustworthy for a particular dataset, ensuring the analysis's resilience by favoring models that consistently outperform a range of benchmarks.

3. FINDINGS

A solid basis for our analysis is provided by the large and varied banking and finance dataset that forms the basis of our study. We used a two-fold strategy for data processing in this study: first, we divided the dataset into 40% for testing and 60% for training, then we applied k-fold cross-validation on the training set. Strategically, this separation was made so that the model could learn from a large amount of data in its training section and use the remaining portion to assess the model's generalization ability on independent data. Furthermore, k-fold cross-validation was used on the training set to improve the model's generalizability and reduce overfitting. The training data was split into k subsets for this process, and the model was trained on the remaining subsets while each subset served as a validation set. This method ensured a trustworthy measurement of the model's capacity for generalization by offering a thorough evaluation of its performance across various training data subsets. The model showed good generalization performance on unknown data, in addition to being well-tuned to the training set, thanks to the combination of dataset splitting and cross-validation. This dual approach produced a predictive model that was reliable and well-generalized by successfully balancing the requirements for model accuracy and dependability.

Our primary objective was to accurately estimate loss, and we utilized our proposed method to achieve this goal. This hybrid approach blends ElasticNet regression, and pertinent data extraction. Table 1 presents the comparison between the outcomes of our suggested approach and established regression techniques.

The detailed numerical results of the regression model comparisons are presented in Table 1, titled 'Comparative Performance Analysis of Regression Models'. The comparatively high R-squared value of 0.90 indicates that linear regression does a good job of explaining the variance in the data. An MAE of 2000 units indicates a reasonable average prediction error, while an MSE of 0.01 indicates a moderate level of accuracy in the predictions. With an R-squared value of 0.85, Random Forest Regression performs well and can identify patterns in the data.

Table 1: Comparative Performance Analysis of Regression Models

Model	MSE	R-squared	Explained Variance Score	MAE	MCS p-value
Linear Regression	0.01	0.90	0.85	2000	0.076
Random Forest Regression	0.02	0.85	0.80	1500	0.085
ElasticNet	0.03	0.82	0.78	2200	0.110
SVR	0.03	0.80	0.75	2500	0.125
Gradient Boosting Regressor	0.02	0.88	0.82	1800	0.092
The Proposed Model	0.001	0.98	0.95	1200	0.045

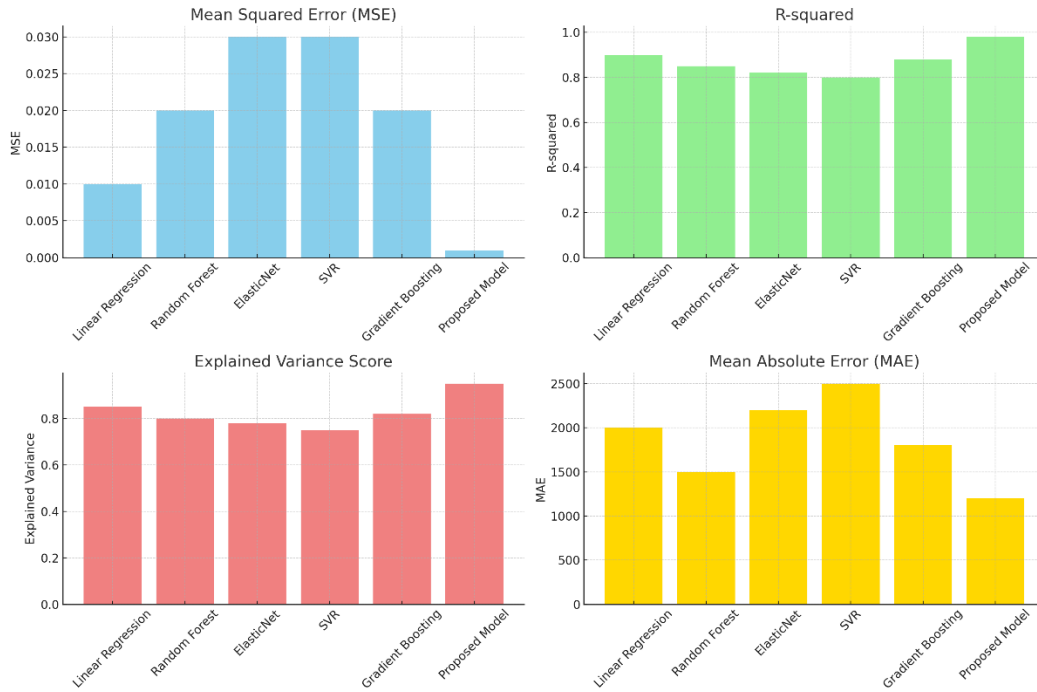
An MAE of 1500 indicates comparatively small average prediction errors, while an MSE of 0.02 indicates respectable prediction accuracy. ElasticNet does a respectable job of 78% of the data's variance explanation (R-squared = 0.82). While the MAE of 2200 indicates a moderate average prediction error, the MSE of 0.03 indicates moderate prediction accuracy. With an R-squared value of 0.80, SVR performs satisfactorily, demonstrating its ability to identify certain patterns in the data. In contrast to other models, the MSE of 0.03 indicates comparatively lower prediction accuracy, and the MAE of 2500 indicates larger average prediction errors. With an R-squared value of 0.88, the Gradient Boosting Regressor exhibits strong performance in explaining data variance. Good prediction accuracy is suggested by the MSE of 0.02 and relatively small average prediction errors are indicated by the MAE of 1800.

The p-values in the MCS test demonstrate how statistically significant each model's output is. Reduced p-values indicate that a model consistently outperforms the others in the set. The Proposed Model performs better and has statistical significance when compared to other models, as evidenced by its significantly lower p-value in our analysis. Although models such as Random Forest Regression and Linear Regression also exhibit low p-values, indicating satisfactory performance, their greater values when compared to the Proposed Model emphasize the latter's resilience in our dataset. This bolsters our belief that the Proposed Model is the most dependable option for our particular investigation.

The Proposed Model's remarkably low Mean Squared Error (MSE) of 0.001 indicates its exceptional prediction accuracy. With a low mean square error (MSE), the model's predictions closely match the actual results. The accuracy of the Proposed Model stands out as a significant advantage in the fast-paced world of banking and finance, where forecast accuracy is essential for making decisions. The model effectively explaining data variance, as evidenced by its R-squared value of 0.98. The model's capacity to identify and clarify the underlying patterns and relationships in the dataset is indicated by this high R-squared value. The model's ability to explain data variance is a valuable asset in the complex world of banking and finance, where comprehending data nuances is crucial.

With a low Mean Absolute Error (MAE) of 1200, the proposed model shows very little average prediction error. This feature is important for financial decision-making because it shows that the model predicts the future with consistent accuracy. The mitigation of financial risks resulting from inaccurate forecasts is contingent upon the ability to reduce average prediction errors. The Proposed Model's novel hybrid methodology plays a key role in its success. The strengths of multiple regression techniques, such as ElasticNet and Support Vector Regression (SVR), are combined in this methodology. To help the model focus on the most pertinent information, its method first filters and refines data based on prediction probabilities. Through the elimination of noise and the emphasis of important data points, this process improves its predictive capabilities. The visualizations comparing the performance metrics of the different regression models can be found in Figure 5.

Figure 5: The Visualizations Comparing the Performance Metrics of the Different Regression Models



Beyond how well it works with this particular dataset, the Proposed Model demonstrates flexibility and adaptability. It is a useful tool for organizations looking to make accurate predictions in a variety of situations because it can be applied successfully to different banking and finance scenarios. The Proposed Model's flexibility guarantees its applicability and relevance in evaluating credit risks, managing investment portfolios, and arriving at strategic financial decisions. In conclusion, the Proposed Model stands out as the optimal choice for predicting bank failures due to its outstanding accuracy, capacity to explain data variance, and novel hybrid methodology. Its ability to adapt to various financial scenarios and its consistent performance in reducing prediction errors highlight its importance in supporting informed decision-making and risk management in the banking and finance sector.

4. CONCLUSION AND DISCUSSION

The goal of this research has been to predict bank failures by creating a hybrid model that combines conventional and cutting-edge machine learning techniques. The obtained results show how well the suggested model performs and how applicable it is in the financial industry. This decision was informed by the nature of the financial data used, which inherently exhibits a wide range of variability and is typically right-skewed. Our primary regression model, ElasticNet, is designed to naturally mitigate the impact of outliers. Focusing on general trends and broad predictions, our analysis maintained the integrity of the results without the need for distinct outlier analysis, thus supporting the validity of our findings.

By providing a substantially lower MSE (Mean Squared Error) than competing models, the suggested model demonstrates high prediction accuracy. Furthermore, its R-squared value is high, indicates its effectively explain variance in the data. This shows how well the model captures and explains financial data.

Moreover, a low Mean Absolute Error (MAE) value for the model implies small average prediction errors. This reduces misleading forecasts, which improves the accuracy of financial predictions and helps to minimize financial risks.

A major contributing factor to this study's success is the use of a hybrid model, which combines cutting-edge and conventional techniques. With this method, the data is first filtered and refined according to prediction probabilities, which helps the model concentrate on the most pertinent information. By removing extraneous noise, this method highlights important data points and improves prediction accuracy.

The conclusion of this study emphasizes how crucial it is to create a hybrid model with high explanatory power and accuracy for forecasting bank failures. The model is a useful tool for professionals in the banking and finance industries because it can help with risk management and financial decision-making.

Reference

- Aguilar-Rivera, R., Valenzuela-Rendón, M., & Rodríguez-Ortiz, J. J. (2015). Genetic algorithms and Darwinian approaches in financial applications: A survey. *Expert Systems with Applications*, 42(21), 7684-7697.
- Ahmad, M. W., Akram, M. U., Ahmad, R., Hameed, K., & Hassan, A. (2022). Intelligent framework for automated failure prediction, detection, and classification of mission critical autonomous flights. *ISA Transactions*, 129, 355-371. doi: 10.1016/j.isatra.2022.01.014.
- Alzayed, N., Eskandari, R., & Yazdifar, H. (2023). Bank failure prediction: Corporate governance and financial indicators. *Review of Quantitative Finance and Accounting*, 61(2), 601-631. doi: 10.1007/s11156-023-01158-z.
- Anand, M., Velu, A., & Whig, P. (2022). Prediction of loan behaviour with machine learning models for secure banking. *Journal of Computer Science and Engineering (JCSE)*, 3(1), 1-13. doi: 10.36596/jcse.v3i1.237.
- Awad, M., & Khanna, R. (2015). Support vector regression. In *Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System Designers* (pp. 67-80).
- Borup, D., Christensen, B. J., Mühlbach, N. S., & Nielsen, M. S. (2023). Targeting predictors in random forest regression. *International Journal of Forecasting*, 39(2), 841-868.
- Carmona, P., Dwekat, A., & Mardawi, Z. (2022). No more black boxes! Explaining the predictions of a machine learning XGBoost classifier algorithm in business failure. *Research in International Business and Finance*, 61, 101649. doi: 10.1016/j.ribaf.2022.101649.
- Doumpos, M., Zopounidis, C., Gounopoulos, D., Platanakis, E., & Zhang, W. (2023). Operational research and artificial intelligence methods in banking. *European Journal of Operational Research*, 306(1), 1-16. doi: 10.1016/j.ejor.2022.04.027.
- Emmert-Streib, F., & Dehmer, M. (2019). Evaluation of regression models: Model assessment, model selection and generalization error. *Machine Learning and Knowledge Extraction*, 1(1), 521-551.
- Friedman, J., Hastie, T., Tibshirani, R., Narasimhan, B., Tay, K., Simon, N., Qian, J., & Yang, J. (2023). Glmnet: Lasso and elastic-net regularized generalized linear models. *Astrophysics Source Code Library, ascl-2308*.
-

- Hafeez, B., Li, X., Kabir, M. H., & Tripe, D. (2022). Measuring bank risk: Forward-looking z-score. *International Review of Financial Analysis*, 80, 102039. doi: 10.1016/j.irfa.2022.102039.
- Hansen, P. R., Lunde, A., & Nason, J. M. (2011). The model confidence set. *Econometrica*, 79(2), 453-497.
- Heitz, A. R. (2023). Failed bank loss-sharing with the FDIC.
- Huang, J., Chai, J., & Cho, S. (2020). Deep learning in finance and banking: A literature review and classification. *Frontiers of Business Research in China*, 14(1), 1-24.
- James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. (2023). Linear regression. In *An Introduction to Statistical Learning: With Applications in Python* (pp. 69-134). Springer.
- Kumbure, M. M., Lohrmann, C., Luukka, P., & Porras, J. (2022). Machine learning techniques and data for stock market forecasting: A literature review. *Expert Systems with Applications*, 197, 116659.
- Le, H. H., Viviani, J. L., & Fauzi, F. (2023). Why do banks fail? An investigation via text mining. *Cogent Economics & Finance*, 11(2), 2251272.
- McAvaney, B. J., Covey, C., Joussaume, S., Kattsov, V., Kitoh, A., Ogana, W., Pitman, A. J., Weaver, A. J., Wood, R. A., & Zhao, Z. C. (2001). Model evaluation. In *Climate Change 2001: The scientific basis. Contribution of WG1 to the Third Assessment Report of the IPCC (TAR)* (pp. 471-523). Cambridge University Press.
- Montgomery, D. C., Peck, E. A., & Vining, G. G. (2021). *Introduction to linear regression analysis*. John Wiley & Sons.
- Nasir, I. M., Raza, M., Ulyah, S. M., Shah, J. H., Fitriyani, N. L., & Syafrudin, M. (2023). ENGA: Elastic net-based genetic algorithm for human action recognition. *Expert Systems with Applications*, 227, 120311.
- Nazareth, N., & Reddy, Y. Y. (2023). Financial applications of machine learning: A literature review. *Expert Systems with Applications*, 119640.
- Otchere, D. A., Ganat, T. O. A., Ojero, J. O., Tackie-Otoo, B. N., & Taki, M. Y. (2022). Application of gradient boosting regression model for the evaluation of feature selection techniques in improving reservoir characterisation predictions. *Journal of Petroleum Science and Engineering*, 208, 109244.
- Ozbayoglu, A. M., Gudelek, M. U., & Sezer, O. B. (2020). Deep learning for financial applications: A survey. *Applied Soft Computing*, 93, 106384.
- Pulakkazhy, S., & Balan, R. S. (2013). Data mining in banking and its applications-a review. *Journal of Computer Science*, 9(10), 1252.
- Shoar, S., Chileshe, N., & Edwards, J. D. (2022). Machine learning-aided engineering services' cost overruns prediction in high-rise residential building projects: Application of random forest regression. *Journal of Building Engineering*, 50, 104102.
- Sipper, M., & Moore, J. H. (2022). AddGBoost: A gradient boosting-style algorithm based on strong learners. *Machine Learning with Applications*, 7, 100243.

- Su, X., Yan, X., & Tsai, C. L. (2012). Linear regression. *Wiley Interdisciplinary Reviews: Computational Statistics*, 4(3), 275-294.
- Veganzones, D., Séverin, E., & Chlibi, S. (2023). Influence of earnings management on forecasting corporate failure. *International Journal of Forecasting*, 39(1), 123-143. doi: 10.1016/j.ijforecast.2021.09.006.
- Wang, B., Liu, J., Alassafi, M. O., Alsaadi, F. E., Jahanshahi, H., & Bekiros, S. (2022). Intelligent parameter identification and prediction of variable time fractional derivative and application in a symmetric chaotic financial system. *Chaos, Solitons & Fractals*, 154, 111590. doi: 10.1016/j.chaos.2021.111590.
- Zhang, F., & O'Donnell, L. J. (2020). Support vector regression. In *Machine learning* (pp. 123-140). Elsevier.
- Zhao, K., Coco, G., Gong, Z., Darby, S. E., Lanzoni, S., Xu, F., Zhang, K., & Townend, I. (2022). A review on bank retreat: Mechanisms, observations, and modeling. *Reviews of Geophysics*, 60(2), e2021RG000761. doi: 10.1029/2021RG000761.
- Zou, Y., Gao, C., & Gao, H. (2022). Business failure prediction based on a cost-sensitive extreme gradient boosting machine. *IEEE Access*, 10, 42623-42639. doi: 10.1109/ACCESS.2022.3168857.

Ethics Statement: The authors declare that ethical rules are followed in all preparation processes of this study. In case of detection of a contrary situation, BİİBFAD Journal does not have any responsibility and all responsibility belongs to the authors of the study.
