



Analysis of Credit Risk Using Machine learning

P. R ASHALATHA

Senior Grade Lecturer in CSE
Government Polytechnic K.R.Pete
ashl12.pr@gmail.com

CHANNAPPA A.

Senior Grade Lecturer in CSE
Government Polytechnic Kudligi
channappajeevitha@gmail.com

SURESH B

Senior Grade Lecturer in CSE
Government Polytechnic Koppal
suresh.arb@gmail.com

Abstract

One of the main study areas in the banking industry is the analysis of credit scoring, which is an efficient method for assessing credit risk. In the banking industry, machine learning is utilised for a wide range of tasks, including data analysis. Using categorization approaches, modern techniques such as machine learning have created a self-regulating mechanism for data analysis. The classification method uses supervised learning, in which the computer classifies the new dataset using the knowledge it has gained from the input data. This research article compares different machine learning approaches used to assess credit risk. Various machine learning algorithms are learned and applied on the dataset to credit transactions that need to be approved or denied. The German Credit dataset, which comprises 1000 instances and 21 attributes and is taken from the UCI repository, is used to put the ideas into practice. Depending on these attributes, the transactions are either approved or rejected. The results of this study's comparison of various algorithms, including Neural Network, Logistic Regression, Naive Bayes, Random Forest, and Stacking ensemble model, reveal that stacking method was more accurate at predicting credit risk.



Keywords: Credit Risk evaluation, Neural Network, Logistic Regression, Naive Bayes, Random Forest, Stacking.

I. INTRODUCTION

Since the previous ten years, credit risk appraisal has grown in significance across various industries, and the field of banking risk management has flourished. For the banking industry, the introduction of credit scoring and credit risk assessment was quite advantageous. Credit risk evaluation can be defined as the identification of the risk levels associated with the credit transaction as to whether the party will adhere to the agreed terms. Credit scoring is a statistical analysis done by financial institutions and lenders to predict the potential risk, corresponding to a transaction. There are two types of credit risk evaluations. Candidates are divided into "good" and "poor" credit risks in the first group. This process of grouping the data based on financial data is known as application scoring. The applicant's payment history, payment patterns, and other factors are taken into account in the second category. The term for this is behavioural scoring[1].

The banking industry has created a number of sophisticated tools in recent years to evaluate the credit risk associated with a few aspects of their operations. Better risk estimation and management as a result has helped to facilitate successful corporate transactions. When the risk is correctly evaluated, a sizable number of business partners may employ it.

The goal of this study is to compare and contrast the performance of several credit risk assessment methods, including Neural Network, Logistic Regression, Naive Bayes, Random Forest, and stacking algorithm. The strategies are evaluated using a genuine dataset of German credit data with the aim of evaluating them, and the findings are used to examine the superior technique that could be applied to precisely assess the credit risk of the banking industry.

II. LITERATURE REVIEW

Various machine learning classification techniques have been presented in recent years to assess credit risk. This section offers an overview of the many methods used to assess credit risk that are described in the papers that follow.

The appraisal of credit risk gained importance as Basel II got more and more widespread. Basel II's requirement has led to an increase in loan applications from clients and competition among



financial institutions over who will approve loans. Support Vector Machine (SVM) was used by Li, Shiue, and Huang to evaluate the customer loans, although only on a small sample of data. For a tiny portion of the sample data, a non-linear nonparametric model was built using machine learning methods to anticipate the consumer credit risk. [2] discussed that the activity related status data will be communicated consistently and shared among drivers through VANETs keeping in mind the end goal to enhance driving security and solace. Along these lines, Vehicular specially appointed systems (VANETs) require safeguarding and secure information correspondences. Without the security and protection ensures, the aggressors could track their intrigued vehicles by gathering and breaking down their movement messages. A mysterious message confirmation is a basic prerequisite of VANETs. To conquer this issue, a protection safeguarding confirmation convention with expert traceability utilizing elliptic bend based chameleon hashing is proposed. Contrasted and existing plans Privacy saving confirmation utilizing Hash Message verification code, this approach has the accompanying better elements: common and unknown validation for vehicle-to-vehicle and vehicle-to-roadside interchanges, vehicle unlinkability, specialist following capacity and high computational effectiveness

The model was built using the Classification and Regression Trees (CART) algorithm, which was discussed by Khandani, Kim, and Lo[3]. CART is a commonly used analytical technique in which an output variable is related to a set of input variables through a series of binary relations that are repeated.

The requirement for an efficient technique to assess client credit risk was covered by Devasena[4], and Huang, Liu, and Ren. The author of these publications discussed several supervised learning classification algorithms that were applied to various data sets. To compare various methods, a variety of metrics were employed. Memory-based classifiers including the IBk classifier, Kstar classifier, and LWL classifier that were used in the German credit data set were discussed by Devasena [4]. Different approaches, including Random Trees, basic CART, PART, J48, Fuzzy, and NBTrees, were analyzed by Gulsoy and Kulluk[6]. The criteria of precision, number of rules, Kappa statistics, recall, and accuracy were utilised to assess credit risk. According to Huang, Liu, and Ren, the examination of the Chinese enterprise dataset revealed that the Probabilistic Neural Network (PNN) has the lowest error rate and highest AUC value.

In order to assess credit risk, Khashman[7] looked at the Emotional Neural Network (EmNN) model with 12 neurons applied to the Australian credit dataset. EmNN is based on a learning algorithm for emotional back propagation. Using the ensemble model applied to the German credit dataset, Wang, Yu, and Ji[8] compared the various classification strategies (Random Forest, Naive Bayes, XGBoost, and RF-Bagging). The management of credit risk is one of the



major issues in a corporate credit rating, according to Zhong, Miao, Shen, and Feng [9], and scorecards are frequently used to handle this problem. Artificial intelligence techniques like ANN and SVM do amazingly well in autonomous credit scoring despite the substantial reliance on user involvement. They also talked about how different methods, including BP, ELM, SVM, and SLNF, were used on preprocessed and normalized real-world financial data. Different machine learning techniques, including Logistic Regression, Neural Networks, Bayesian Networks, and Random Forest, were presented by authors [10].

III. PROPOSED METHODOLOGY

Dataset

German credit data [17], which classifies consumers as "good" or "poor" credit risks based on a series of criteria, served as the dataset for this study. This dataset is freely accessible and may be found in the UCI Machine Learning Repository [18]. 1000 occurrences of each of 21 attributes, including a categorization attribute for each instance, are included in the dataset. This dataset has already been employed in credit scoring and credit risk assessment, which is why it was selected for this investigation. The purpose of this study is to determine whether each instance's outcome will be good or bad using a data set that has been classified based on features or attributes. Various machine learning classification algorithms have been used on the same data set to compare how well they perform in this task. Each instance is described by the first 20 attributes utilised in this analysis, and the final attribute is used to determine if a transaction is good or poor. The details of the attributes are tabulated in table I.



Table I. List attributes in the data set

Attributes	Type
1. Creditability 2. Account Balance 3. Status of payment 4. Purpose 5. Savings in cost 6. Current employment period 7. Sex and Marital Status 8. Guarantors 9. Most precious resources 10. Simultaneous loans 11. Type of house 12. Employment 13. Number of dependents 14. Telephone 15. Foreign Workers	Categorical
16. Credit length (in months) 17. Credit Amount 18. Current address duration 19. Lifespan 20. Amount of loans from this bank	Numeric

Methodology

Machine learning classification algorithms namely Neural Network, Logistic Regression, Naive Bayes, and Random Forest were applied on the dataset. Later the models are aggregated and stacking is done.



The dataset is divided into sections for testing and training. Due of its improved accuracy, the K-Fold-Cross-Validation technique is frequently employed. The data is divided into M folds using the K-Fold-Cross-Validation technique, where M folds are used for learning and the Mth fold is used for testing the data. The data operates on various percentage divisions. By setting the sub-samples number to 10, this technique was utilised to confirm the analysis's standard. By contrasting the outcomes of the procedures using different parameters, the system is assessed. The random forest is implemented by taking about 10 trees and neural network model is built with 100 neuron hidden layers, ReLu activation function. [5] discussed about a method, Sensor network consists of low cost battery powered nodes which is limited in power. Hence power efficient methods are needed for data gathering and aggregation in order to achieve prolonged network life. However, there are several energy efficient routing protocols in the literature; quiet of them are centralized approaches, that is low energy conservation. This paper presents a new energy efficient routing scheme for data gathering that combine the property of minimum spanning tree and shortest path tree-based on routing schemes. The efficient routing approach used here is Localized Power-Efficient Data Aggregation Protocols (L-PEDAPs) which is robust and localized. This is based on powerful localized structure, local minimum spanning tree (LMST). The actual routing tree is constructed over this topology. There is also a solution involved for route maintenance procedures that will be executed when a sensor node fails or a new node is added to the network.

Stacking: The aim of stacking is to explore a space of different models for the same problem. The concept is that you can approach a learning problem with several sorts of models that can learn a portion of the problem but not the entire problem space. As a result, you may create a variety of learners and utilise them to create intermediate predictions, one for each learnt model. Then you include a new model that picks up on the same target from the intermediate forecasts. The name of the last model comes from the claim that it is layered on top of the others. As a result, you may enhance your general performance, and frequently you produce a model that is superior to any individual intermediate.



Different settings and parameters were used to evaluate the credit risk assessment approaches. Common metrics including F1 score, specificity, accuracy, sensitivity, error rate, and precision were used to compare these approaches. The results of the various methods used to assess the credit risk are listed in Table II. The ROC analysis is also demonstrated in figure 1.

Table II: Performances of the prediction model for all the features

Model	AUC	Accuracy	F1	Precision	Recall
Stack	0.80	0.76	0.84	0.79	0.90
Random Forest	0.76	0.75	0.84	0.78	0.91
Neural Network	0.75	0.74	0.82	0.79	0.85
Naive Bayes	0.79	0.75	0.83	0.82	0.84
Logistic Regression	0.79	0.75	0.83	0.79	0.87

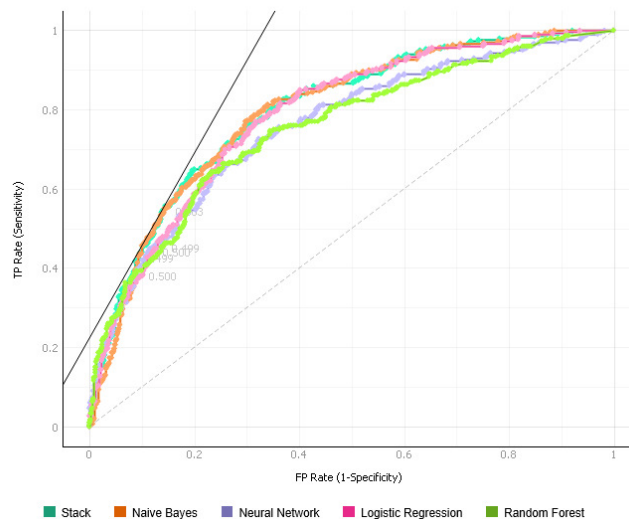


Fig1: ROC analysis of prediction models

IV. CONCLUSION AND FUTURE WORK

There have been several attempts to create rapid and effective ways to foresee the future because the global markets are rife with danger. The development of credit ratings and the assessment of



credit risk were very advantageous for the banking industry. In this study, credit risk in the German credit dataset was assessed using a variety of machine learning algorithms. These have been put into practise and evaluated on different classification techniques, including LR, BN, NN, RF and stacking. By using them on a pre-existing dataset called German credit, which has a thousand transactions each day, the strategies are evaluated. According to the analysis above, the Random Forest methodology increases the accuracy of credit risk assessment. Future research can evaluate other deep learning algorithms to determine if accuracy improves.

Reference

- [1] A. Khashman, "Neural networks for credit risk evaluation: Investigation of different neural models and learning schemes," *Expert Systems with Applications*, vol. 37, pp. 6233-6239, 2010.
- [2] Christo Ananth, Dr.S. Selvakani, K. Vasumathi, "An Efficient Privacy Preservation in Vehicular Communications Using EC-Based Chameleon Hashing", *Journal of Advanced Research in Dynamical and Control Systems*, 15-Special Issue, December 2017, pp: 787-792.
- [3] A. E. Khandani, A. J. Kim, and A. W. Lo, "Consumer credit-risk models via machine-learning algorithms," *Journal of Banking & Finance*, vol. 34, pp. 2767-2787, 2010.
- [4] C. L. Devasena, "Adeptness Evaluation of Memory Based Classifiers for Credit Risk Analysis," in *2014 International Conference on Intelligent Computing Applications*, 2014, pp. 143-147.
- [5] Christo Ananth, S.Mathu Muhila, N.Priyadharshini, G.Sudha, P.Venkateswari, H.Vishali, "A New Energy Efficient Routing Scheme for Data Gathering ", *International Journal Of Advanced Research Trends In Engineering And Technology (IJARTET)*, Vol. 2, Issue 10, October 2015), pp: 1-4.
- [6] N. Gulsoy and S. Kulluk, "A data mining application in credit scoring processes of small and medium enterprises commercial corporate customers," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 9, p. e1299, 2019.
- [7] A. Khashman, "Credit risk evaluation using neural networks: Emotional versus conventional models," *Applied Soft Computing*, vol. 11, pp. 5477-5484, 2011.
- [8] M. Wang, J. Yu, and Z. Ji, "Personal credit risk assessment based on stacking ensemble model," in *International Conference on Intelligent Information Processing*, 2018, pp. 328-333.
- [9] H. Zhong, C. Miao, Z. Shen, and Y. Feng, "Comparing the learning effectiveness of BP, ELM, I-ELM, and SVM for corporate credit ratings," *Neurocomputing*, vol. 128, pp. 285-295, 2014.
- [10] E. Caldeira, G. Brandao, H. Campos, and A. Pereira, "Characterizing and evaluating fraud in electronic transactions," in *2012 Eighth Latin American Web Congress*, 2012, pp. 115-122.