

Predictive Lending AI

Kommaraju Rupa Kusuma¹, Dr. Manne Suneetha².

¹M.Tech Student, Data Science, Dept. of Information Technology, V R Siddhartha Engineering College, Kanuru, India (Email: rupakommaraju23@gmail.com)

²Head of the Dept, Professor, Dept. of Information Technology, V R Siddhartha Engineering College, Kanuru, India (Email: hodit@vrsiddhartha.ac.in)

Article Info

Page Number: 335-344

Publication Issue:

Vol. 72 No. 1 (2023)

Article History

Article Received: 12 October 2022

Revised: 24 November 2022

Accepted: 18 December 2022

Abstract: The product Predictive Lending AI is a AI based Automatic Lending Solution for Non-Banking Financial Company (NBFC). It helps a global lending firm to reduce the loan processing time by 40%, minimize business risk and increase revenues with effective up-selling of lending products. There are three use cases (Application scoring, Defaulter Prediction and Churn Prediction). The first use case is Application Scoring in which we calculate a score based on the input and approve a loan based on that score. The second use case is about Loan Defaulter Prediction. The person who takes loan from an organization and doesn't repay the loan amount is called Loan defaulter. The third use case is about Churn Prediction. Churn is defined as the movement of customers from one provider to another.

Keywords: Prediction, Non-Banking Financial Company (NBFC), Defaulters, Application Scoring, Churn Prediction.

Introduction

As, the count of the people applying for the loans has been increasing for various reasons in the recent years. The bank employees are unable to analyze or predict whether the customer can payback the amount or not (good customer or bad customer) for the given interest rate. In order to succeed in the stream of banking, one has to have an idea about the behavioral patterns of various customers based on their transaction history. This is what, our model Lending AI is doing, by predicting the cases of customers who may finally end up as a loan defaulter or not. To predict the credit default, several methods have been proposed. The use of method depends on the complexity of banks and financial institutions, size, and type of the loan. The commonly used method has been discrimination analysis. By using scoring models that are AI-based and use deep learning, banks and financial institutions can access more realistic predictions on credit risk, using customers' credit history and the power of big data. This way credit can be approved to the right people and better pricing options offered to people who deserve it. The output of the model will generate a binary value that can be used as a classifier that will help banks to identify whether the borrower will default or not default. As the final step in the direction, linear regression method is also going to be performed on the dataset.

Related Works

In the previous version the framework is developed to effectively identify the Probability of Default of a Bank Loan applicant. The metrics derived from the predictions reveal the high

accuracy and precision of the built model. The model proposed in an effective prediction model for predicting the credible customers who have applied for bank loan. Decision Tree is applied to predict the attributes relevant for credibility.

This prototype model can be used to sanction the loan request of the customers or not. The model proposed in has been built using data from banking sector to predict the status of loans. This model uses three classification algorithms namely j48, Bayes Net and naive Bayes. The model is implemented and verified using Weka. The best algorithm j48 was selected based on accuracy. An improved Risk prediction clustering Algorithm that is multi-dimensional is implemented to determine bad loan applicants.

In this work, the Primary and Secondary Levels of Risk assessments are used and to avoid redundancy, Association Rule is integrated. In a decision tree model was used as a classifier and for feature selection genetic algorithm is used. The model was tested using Weka.

Support Vector Machine, Decision Tree, Logistic Regression, Neural Network, Perceptron model, all these techniques are combined in this model. The effectiveness of applying the above techniques on credit scoring is studied. The analysis results show the performance is outstanding based on accuracy. The aim of the study is to introduce a discrete survival model to study the risk of default and to provide the experimental evidence using the Italian banking system.

Data mining in banking Due to tremendous growth in data the banking industry deals with, analysis and transformation of the data into useful knowledge has become a task beyond outstanding based on accuracy. The aim of the study is to introduce a discrete survival model to study the risk of default and to provide the experimental evidence using the Italian banking system.

Data mining in banking Due to tremendous growth in data the banking industry deals with, analysis and transformation of the data into useful knowledge has become a task beyond human ability. Data mining techniques can be adopted in solving business problems by finding patterns, associations and correlations which are hidden in the business information stored in the databases.

Proposed Methodology

In this paper different models like Random Forest, Logistic Regression and Decision Tree are built on the dataset. The input data is loaded from the database that is stored and the model metrics should be identified such as whether the data is structured or unstructured and understanding of features. The data analysis includes what are the important features that are required for prediction and any features that should be imputed or removed, etc. And the target variable should be identified. Next, the model is built by choosing the right Algorithm for prediction. The model should be able to predict for the incoming new data. After the model is built the leads that are qualified. To increase the qualified leads, we have again feed the model with data and can increase the performance.

The Historical data from NBFC's is used to build a model in the Rapid Miner tool. After completing EDA, a model is built using ML algorithms and predicts the results for

unseen data using that model.

RapidMiner tool is used for providing a solution to this use case. RapidMiner is a data science software platform developed by the company of the same name that provides an integrated environment for data preparation, machine learning, deep learning, text mining, and predictive analytics. RapidMiner's data science platform delivers lightning-fast business impact for over 40,000+ organizations in every industry to drive revenue, reduce costs, and avoid risks.



Fig.1. Architecture Diagram

The product Lending AI is an AI based Automatic Lending Solution for Non-Banking Financial Company (NBFC). It helps a global lending firm to reduce the loan processing time by 40%, minimize business risk and increase revenues with effective up-selling of lending products.

A. NBFC

A Non-banking financial company or non-banking financial institution is a financial institution that does not have a full banking license or is not supervised by a national or international banking regulatory agency. It is an institution, which is a company and has principal business of receiving deposits under any scheme or arrangement in one lump sum or in instalments by way of contributions or in any other manner, is also a non-banking financial company (Residual non-banking company).

B. EDA

EDA is a phenomenon under data analysis used for gaining a better understanding of data aspects like, it identifies the main features of the data and also it identifies which variables are important for our problem.

C. Data Pre-processing

Machine Learning algorithms don't work so well with processing raw data. Before we feed such data to an ML algorithm, we must pre-process it. In other words, we must apply some transformations on it. With data pre-processing, we converted raw data into a clean data set. Some ML models need information to be in a specified format. To pre-process data, we will use the library `scikit-learn`.

D. Data Visualization

Data visualization is the discipline of trying to understand data by placing it in a visual context so that patterns, trends and correlations that might not otherwise be detected can be exposed. Python offers multiple great graphing libraries that come packed with lots of different features are Histograms, Bar Charts, Scatter Plots, Using Seaborn, Horizontal Bar charts, Staticmaps, Network diagrams.

Data Cleaning And Preparation

A. Data Description

The dataset from `kaggle.com` is used for statistical modelling. The dataset is divided into training dataset and testing dataset. The Training dataset contains of 4,96,307 examples with 1 special attribute and 146 regular attributes. The Testing dataset contains of 55,143 examples with 146 regular attributes.

B. Data Preprocessing

In this stage, identified the attributes with more than 50% missing values and removed those attributes. There are many observations they are there are 42 attributes with more than 50% missing values, there are 6 attributes which have a greater number of values, the target label loan status has a high-class imbalance problem, there are no duplicate data points in the given dataset.

C. Data Cleaning

I) Imputing Missing values

In this process the missing values in Numerical attributes are replaced by mean and the missing values in Nominal attributes are replaced by mode.

II) Converting Nominal to Numerical

In this process the Nominal attributes are converted to Numerical attributes by dummy encoding.

D. Data Sampling

In case of splitting the data, Stratified sampling technique is used. Stratified sampling is a type of sampling method in which the total population is divided into smaller groups or strata to complete the sampling process. The strata is formed based on some common characteristics in the population data. After dividing the population into strata, the researcher randomly selects the sample proportionally.

E. Model Building

After splitting the data in the 70:30 ratio, the training data is sent to the classification model. The model which is built is applied on the testing data. Then check the actual results with predicted values. Keep on changing the attributes till best performance is achieved.

Once when the best performance is achieved, then we can predict the results on the unseen data.

F. Experimental Setup

The attribute selection is done, and the attribute selection reduces memory requirements and increases the accuracy of the model. The random Forest, Decision Tree, Logistic regression, Support Vector Machine algorithms are applied on the same dataset for building the model, and then the comparison is done between the four algorithms for identifying the best approach.

I) Random Forest

The most important aspect of the random forest algorithm is the variable importance ranking. It creates recursive partitioning trees using a majority vote. A number 'm' is specified which is much smaller than the total number of attributes. At each node, 'm' variables are selected at random out of the total number of attributes, and then split is performed.

II) Decision Tree

The Recursive binary splitting technique can be used to perform split at a node. In this method, all the attributes are taken into consideration and various split points are tried and tested. They are tested using a cost function and the split with the best cost is selected.

III) LogisticRegression

A logistic regression is run using each variable against the binary target variable for the result. ROC curve for each variable is plotted. The variable containing the largest area under the curve has the largest relevancy and contributes the most for the result. The feature containing the largest Information gain ratio has the lowest importance. The subset of optimal features is arranged in descending order to obtain the highest relevancy features of the dataset

IV) SupportVectorMachine

SVM is supervised machine learning model with learning algorithms which examine the data and uses that data for regression and classification. This model uses a technique namely a kernel trick to transform the data and based on these transforms of data, it finds the best optimum results. It is not considered as better as than the other machine learning models because it works on less data set.

Result and Analysis

In this we calculate the results of predictions on all the four algorithms that I have used on the dataset in model building. The dataset is split into the training dataset and testing dataset in the 70:30 ratio. The training dataset consists of the 70 % of the data from the dataset and the testing dataset consist of the 30% of the data from the dataset. The performance is calculated based on the F1- Score, Gini Index, Accuracy, ROC. The f1- Score of the Logistic Regression, random forest, Decision Tree, Support Vector Machine are 96.1%, 92.5%, 86.5%, 89.1% respectively.

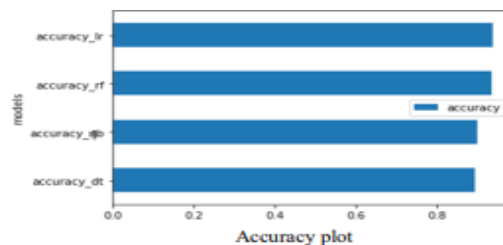


Fig.2. Accuracy Plot

Conclusion

The Logistic Regression is the best approach with 96% F1- Score and great accuracy. The defaulter prediction is used finally on applicants' data and the loan defaulters are identified and the loan amount is safeguarded. The defaulter prediction use-case prevents the banking sector from huge loss and downfall by being alert before sanction the loan amount to the loan defaulters. This use-case is very useful.

accuracy: 99.00%

	true Approved	true Rejected	class precision
pred. Approved	30491	541	98.26%
pred. Rejected	25	25559	99.90%
class recall	99.92%	97.93%	

Fig.3. Confusion Matrix

References

1. Altman, E.I., P. Narayanan, An International Survey of Business Failure Classification Models, Financial Markets, Instruments and Institutions, 1997
2. Altman, E.I., Saunders A, Credit risk measurement: Development over the last 20 years, Journal of Banking and Finance, 1998, 20: 1721~1742
3. Shi Xiquan, Zou Xinyue, The Application of canonical discriminant analysis in credit risk evaluation of enterprise, The Study of Finance and Economics, 2001, 27 (10) : 53~57
4. Zhang Aimin, Zhu Chunshan, Xu Danjian, Principle components prediction model and the empirical study of financial failure of public company, Journal of Finance, 2001, 3 : 10~25
5. A.B. Hussain, and F.K.E. Shorouq, "Credit risk assessment model for Jordanian commercial banks: Neural scoring approach", Review of Development Finance, Elsevier, vol.4, pp.20-28, 2014.
6. A. Blanco, R. Mejias, J. Lara, and S. Rayo, "Credit scoring models for the microfinance industry using neural networks: evidence from Peru", Expert Systems with Applications, vol.40, pp.356-364, 2013.
7. T. Harris, "Quantitative credit risk assessment using support vector machines: Broad versus Narrow default definitions", Expert Systems with Applications, vol. 40, pp. 4404-4413, 2013.
8. A. Abhijit, and P.M. Chawan, "Study of Data Mining Techniques used for Financial Data Analysis", International Journal of Engineering Science and Innovative Technology, vol. 2(3), pp.503-509, 2013.
9. D. Adnan, and D. Dzenana, "Data Mining Techniques for Credit Risk Assessment Task", in Proceedings of the 4th International Conference on Applied Informatics and Computing Theory (AICT'13), 2013, p. 105-110.
10. G. Francesca, "A Discrete-Time Hazard Model for Loans: Some Evidence from Italian
11. J H Cheon, D Kim, Y Kim and Y Song, "Ensemble Method for Privacy-Preserving Logistic Regression Based on Homomorphic Encryption", *IEEE Access*, vol.6, pp.46938-46948, 2018.
12. S Guo, H He and X Huang, "A Multi-Stage Self-Adaptive Classifier Ensemble Model With Application in Credit Scoring", *IEEE Access*, vol.7, pp.78549-78559, 2019.
13. Dayu Xu, Xuyao Zhang, Junguo Hu and Jiahao Chen, "A Novel Ensemble Credit Scoring Model Based on Extreme Learning Machine and Generalized Fuzzy Soft Sets", *Hindawi Mathematical Problems in Engineering*, 2020.

14. S.M.SandR.S.T, "Loan Credibility Prediction System Based on Decision Tree Algorithm", *International Journal of Engineering Research and*, vol.V4, no.09, 2018.
15. P. M. Addo, D. Guegan and B. Hassani, "Credit Risk Analysis Using Machine and Deep Learning Models", *SSRNElectronic Journal*, 2018.
16. S Tabik, R.F. Alvear-Sandoval, M.M. Ruiz, J.L. Sancho-Gómez, A.R. Figueiras-Vidal and F.Herrera, "A Tutorial On Ensembles And Deep Learning Fusion With Mnist As Guiding Thread: A Complex Heterogeneous Fusion Scheme Reaching 10 Digit error", *arxiv.org*, January 2020.
17. <https://docs.rapidminer.com/latest/studio/getting-started/ui-overview.html>
18. A.Saxena and M.Prasad, A review of clustering techniques and developments publication *Neurocomputing*, Elsevier, 2017
19. M.Ala'raj and M.Fabbod, "Classifiers consensus system approach for credit scoring", *Knowledge-Based Syst*, vol.104, pp. 89-105, Jul. 2018.
20. Adnan Dželihodžić and Dženana Đonko, "Comparison of Ensemble Classification Techniques and Single Classifiers Performance for Customer Credit Assessment", *Modeling of Artificial Intelligence*, vol.11, no.3, pp.140-150, 2018.
21. Okokpuije Okesola, O.Kennedy, Adewale and Samuel N. John, "An improved Bank Credit Scoring Model A Naïve Bayesian Approach", *International Conference on Computational Science and Computational Intelligence*, pp.228-233, 2017.
22. H He, W. Zhang and S. Zhang, "A novel ensemble method for credit scoring: Adaption of different imbalance ratios", *Expert Syst. Appl*, vol.98, pp.105-117, May 2018.
23. MI Wang, H Wang, J Wang, H Liu, R Lu, T Duan, et al., "A novel model for malaria prediction based on ensemble algorithms", *PLOS One*, 2019.
24. N. Mohammad and A. Onni, "Credit Risk Grading Model and Loan Performance Of Commercial Banks In Bangladesh", *European Journal of Business Management*, vol.7, no.13, Yuan Sun, Weifeng Jian, Yufeng Fu, "A new perspective of credit scoring for small and medium-sized enterprises based on invoiced data", *IEEE Access*, 2021
25. Wenyu, Zhang, "A Novel Multi-Stage Ensemble Model With a Hybrid Genetic Algorithm for Credit Scoring on Imbalanced Data", *IEEE Access*, 2021
26. Adel Hassan, Rashid Jayousi, "Financial Services Credit Scoring System Using Data Mining", *IEEE Access*, 2020
27. A Safiya Parvin, B Saleena, "An Ensemble Classifier Model to Predict Credit Scoring",
28. *IEEE Access*, 2020
29. Dawn Iris Calibo, Melvin A. Ballera, "Variable Selection for Credit Risk Scoring on Loan Performance Using Regression Analysis", *IEEE Access*, 2019
30. Olatunji J. Okesola, Kennady O. Okokpuije, deyinka A. Adewale, "An Improved Bank Credit Scoring Model", *IEEE Access*, 2018
31. Weisong Chen, Liang Shi, "Credit scoring with F-score based on support vector machine", *IEEE Access*, 2018
32. Radha Vadala, Bandaru Rakesh kumar, "An application of Credit Scoring in E-

- Lending Platform”, IEEE Access, 2019
33. Syed Zamil Hasan Shoumo, Mir Ishrak Maheer Dhruva, Sazzad Hossain, “Application of Machine Learning in Credit Risk Assessment”, IEEE Access, 2019
 34. Yashaswi Raj Suresh Kumar, Brijesh Kumar Singh “Loan Risk Prediction Using Transaction Information” Journal of Analysis and Computation (JAC), 2021
 35. B Spoorthi, Swetha S. Kumar, “Comparative Analysis of Bank Loan Defaulter Prediction Using Machine Learning Techniques”, IEEE, 2021
 36. Jaya Sinha, Rani Asthya, “Machine Learning based Loan Allocation Prediction System for Banking Sector”, IEEE, 2021
 37. P. Maheswari, CH. V. Narayana, “Prediction of Loan Defaulter - A Data Science Perspective”, IEEE, 2020
 38. Apurva Datkhile, Komal Chandak, “Statistical Modelling on Loan Default Prediction Using Different Models”, International Journal of Research in Engineering, Science and Management, 2020
 39. Abhishek Shivanna, Dharma P Agrawal, “Prediction of Defaulters using Machine Learning on Azure ML”, IEEE, 2020
 40. Hafiz Ilyas Tariq, Asim Sohail, “Loan Default Prediction Model Using Sample, Explore, Modify, Model, and Assess (SEMMA)”, Journal of Computational and Theoretical Nanoscience, 2019
 41. T. Harris, “Quantitative credit risk assessment using support vector machines: Broad versus Narrow default definitions”, Expert Systems with Applications, vol. 40, pp. 4404–4413, 2013.
 42. A. Abhijit, and P.M. Chawan, “Study of Data Mining Techniques used for Financial Data Analysis”, International Journal of Engineering Science and Innovative Technology, vol. 2(3), pp. 503-509, 2013.
 43. D. Adnan, and D. Dzenana, “Data Mining Techniques for Credit Risk Assessment Task”, in Proceedings of the 4th International Conference on Applied Informatics and Computing Theory (AICT '13), 2013, p. 105-110. [1] D. J. Hand, G. Blunt, M. G. Kelly and N. M. Adams, “Data Mining For Fun and Profit”, *Statistical Science*, vol. 15, pp. 111-131, 2020.
 44. Dayu Xu, Xuyao Zhang, Junguo Hu and Jiahao Chen, “A Novel Ensemble Credit Scoring Model Based on Extreme Learning Machine and Generalized Fuzzy Soft Sets”, *Hindawi Mathematical Problems in Engineering*, 2020.
 45. S Tabik, R.F. Alvear-Sandoval, M.M. Ruiz, J.L. Sancho-Gómez, A.R. Figueiras-Vidal and F. Herrera, “A Tutorial On Ensembles And Deep Learning Fusion With Mnist As Guiding Thread: A Complex Heterogeneous Fusion Scheme Reaching 10 Digit error”, *arxiv.org*, January 2020.
 46. Yuan Sun, Weifeng Jian, Yufeng Fu, “A new perspective of credit scoring for small and medium-sized enterprises based on invoice data”, IEEE Access, 2021
 47. Wenyu, Zhang, “A Novel Multi-Stage Ensemble Model With a Hybrid Genetic Algorithm for Credit Scoring on Imbalanced Data”, IEEE Access, 2021
 48. Adel Hassan, Rashid Jayousi, “Financial Services Credit Scoring System Using

Data Mining”, IEEE Access, 2020

49. A Safiya Parvin, B Saleena, “An Ensemble Classifier Model to Predict Credit Scoring”, IEEE Access, 2020