# ISSN: 2321-2152 **IJMECCE** International Journal of modern electronics and communication engineering

# E-Mail editor.ijmece@gmail.com editor@ijmece.com

www.ijmece.com



ISSN 2321-2152 <u>www.ijmece.com</u> Vol 13, Issue 1, 2025

# LOAN PREDICTION DATASET USING MACHINE LEARNING WITH DATA NALYSIS

<sup>1</sup>B.Sravanthi, <sup>2</sup>T.Shalini, <sup>3</sup>P.Naveen, <sup>4</sup>B.Anusha, <sup>5</sup>J.S.Radhika,

<sup>1,2,3,4,</sup> U.G.Scholor, Department of IT, Sri Indu College Of Engineering & Technology, Ibrahimpatnam, Hyderabad. <sup>5</sup>Assistant Professor, Department of IT, Sri Indu College Of Engineering & Technology, Ibrahimpatnam, Hyderabad.

**Abstract.** Every day, the demand for loans in banks becomes greater and higher as people's demands continue to rise. The tedious and time-consuming process of screening and validating an applicant's eligibility is usually the last step before a bank processes a loan application. Banks could incur losses due to applicants who fail to repay their loans. Implementing machine learning technologies that employ classification algorithms to forecast qualified loan applicants is a great way to reduce human effort and make successful decisions in the loan approval process.

# **1** Introduction

Machine Literacy is a subset of artificial intelligence that allows computer programs to automatically learn from former tasks. It works by analysing the data, relating patterns, and incorporating minimum mortal intervention. Nearly any workthat can bedone using a datadescriptionpatternorsetofrulescanbedoneusing a machine learning machine. This allows companiesto modifyprocesses that preliminarily only humans could make hypotheticals for client service calls, accounts, and reviews.

Loan distribution is the middle enterprise of virtually everybank.Loandistributionisthemiddleenterpriseof nearly every bank. Utmost of a bank's means come directlyfromthegainsitmakesfromtheloansitmakes.

Themainthingofthebankingterrainistoputwealthin a safe place. Currently, many banks / financial companies approve loans after a verification and validation redemption process, but it is not yet clear if the selected applicant is correct among all applicants. Through this system, it is possible to predictwhether a particular applicant is safe and the entire process of verifying the characteristics will be automated by machinelearningtechnology.Creditforecastingisvery useful for both bank employees and applicants.

Thegoalofthissystemistoprovideaquick,immediate and easy way to select good applicants. It can offer banks special benefits. The credit forecasting system can automaticallycalculatetheweightsforeachfeaturethat participates in credit processing, and the new test data

willprocessthesamefeaturesfortheassignedweights. Themodelcansetadeadlineto seeifthe applicantcan approve the loan. Credit analysis allows to jump to specific applications and check according to priority. Thissystemisexclusivelyforbank/financialcompany management authorities, the entire forecasting process is carriedout privately and nostakeholders can change the process. The results of a particular creditID can be senttovariousdepartmentsofthebanksothattheycan take appropriateaction on demand. Thishelpsallother departments handle other paperwork.

# 2 LiteratureSurvey

According to the authors, the forecasting process begins with data clean-up and processing, missing value substitution, data set experimental analysis, and modelling, and continues to model evaluation and test data testing. A logistic regression model has been executed. The highest accuracy obtained with the originaldatasetis0.811.Modelsarecomparedbasedon performance measurements such as sensitivity and specificity. As a result of analysing, the following conclusionsweredrawn.However,othercharacteristics of customers that play a very important role in lending decisions and forecasting defaulters should also be evaluated. Some other traits, such as gender and marriage history, do not seem to be considered by the company [1]. A credit credibility soothsaying system that helps companies make the right opinions to authorizeorrejectthecreditclaimsofguests.Thishelps the banking assiduity to open effective distribution channels. This means that if the customer has a minimum repayment capacity, their system can avoid futurerisks.Includingothertechniques(usingtheWeka tool) thatarebetterthanthe generaldata mining model has been implemented and tested for domains [2]. The authorsuggeststhat, acreditstatusmodel for predicting loan applicants as valid or standard customers. The proposed model shows a score of 75.08 when classifyingloanaspirantsusingR-Package.Lenderscan use this interpretation to make mortgage choices for mortgage operations. In addition, comparative studies were conducted at different iterative levels. The replication position is a 30- grounded ANN model that offers a more advanced delicacy than other situations. This model can be used to avoid large losses in marketable banks [3]. Six machine learning classification models we reused to predict Android



applications. The model is available in open-source software R. This application works well and meets the requirements of all banks. The downside of this model is that it gives each element a different weight, but in reality, it may be possible to approve aloanonly based onasinglepowerfulelement, which is not possible with thissystem. This component can be easily connected to many other systems. There are cases of computer failure, and the most important weights of contenterrors and features are fixed by the automatic prediction system, and soon, so-called software may be safer, more reliable and more [4].Risk assessment and forecasting is an important task in the banking industry in determining whether a good and lazy loan applicant is applicable. To improve the accuracy of risk, risk assessments are conducted in primary and secondary education. Customer data is extracted and related attributes are selected using information gain theory. Ruleforecastingisperformedforeachcredittypebased onpredefinedcriteria. Approved and rejected applicants are considered"Applicable" and evaluated as"Not Applicable". Corresponding experimental results have shown that the method proposed predicts better accuracy and takeslesstime thanexistingmethods[5]. The main purpose of this design is to prognosticate whichcustomerswillberepaid withaloanbecausethe lenderneedstoanticipatetheproblemthattheborrower won't be suitable to repay the threat. Studies of three models show that logistic regression with a rating is superior to other models, random forests, and decision trees. Poor credit seekers aren't accepted, presumably because they have the option of not paying. In utmost cases, high-value appliers may be eligible for a reduction that may repay the loan. Certain sexual orientations and marriage statusappearto be outof the reach of the company [6].

# **3** FeatureEngineering

Predicated on the field knowledge, this system can developnewfeaturesthatcanaffectthetargetvariables. Created three new functions:

**Total Income**–In Fig 1, as explained in the bivariateanalysis, combined the income of the applicant with the income of the co-applicant. The higher the total income, the more likely there are to get loan approval.



Fig.1.DensityvsTotalIncomeLog

**EMI**– In Fig 2, EMI is theyearly volumethat the seeker must pay to reimburse the loan. The model behind this variableisthat people with lofty. EMI may possess challenges with the prepayment of their loans. EMI can exist figured by taking the rate of the loan volumetothemajority theloanvolumerateoftheloan volume to the majority of the loan volume.





**Balance Income**–In Fig 3,this is the return deserted over after compensating the EMI. The model behind creating this variable is that the advanced the valuation,additionallyprobableapersonistoreimburse the loan and thusadditionally probableit'sto authorize the loan.



Fig.3.DensityvsBalanceIncome

# 4 ProposedFramework

#### **Businessunderstanding:**

In the early stages, the base is on deriving the design fromacustomoutlookandrephrasingthatloreintodata mining challenge delineations and primary designs.

#### Dataconvention:

Thedataconventionaspectfocusesontheoriginaldata library,datacommand,relatingdatarateoutcomes,and a subset of stake for undertaking retired data.

#### Dataprocessing:



Thedataprocessingaspectincludesalltheconditioning to produce the concluding dataset.

#### Modelling:

Thealgorithmwhichwillbeusedfordatamodellingis Logistic Regression using stratified k-folds crossvalidation and Random Forest.

Logistic Regression using stratified k-folds crossvalidation: This system uses validation to see how robust the model is against hidden data. This is an approach for booking distinct exemplifications of reports that don't train the model. Latterly, the system tests the model in this illustration and also finalize it. Some of the generally applied confirmation styles are theconfirmationincubatepath,k-foldcross-validation,

Leaveoneoutcross-validation(LOOCV), and stratified kfold cross-validation. In Table 1, analyzed the mean validation and f1-score of Logistic Regression with the kfolds model.

Table1.AccuracyforLogisticRegressionmodel

The mean validation accuracy for this model turns out to be	0.7214
Themeanvalidationfl scoreforthis model turns out to be	0.8279

InFig4, Visualized theroccurve:



Fig.4.AUCvalueof0.5626

Random Forest: This system was tested to reduce the exactnessbyconformingtothehyperparametersofthis model. Themodelusedgridsearchtomasteroptimized valuations for hyperparameters. Grid - search is a way to elect the stylish one from the family of hyperparameters parameterized by the parameter grid. adapted the max\_depth and n\_estimators' parameters. max\_depth determines the maximum depth of the tree andn\_estimatorsdeterminethenumberoftreesusedin therandomforestmodel.InTabel2,generatedthemean validation accuracy for the hyperparameters.

Table2.MeanValidationAccuracy	ofHyperparameters
MeanValidationAccuracy	0.7947



Fig5.Featureimportancepredictingthetargetvariable

In Fig 5, Credit\_History is the most major point succeeded by Balance Income, Total Income, EMI. Accordingly, feature engineering assisted the model in forecasting the target variable.

### **5** Conclusion

Borrowers use a loan application to qualify for a mortgage. The above research employs a logistic regressionalgorithm-basedpredictionmodel.Tocreate a logistic classification model that predicts loan status, over600sampledatawerecollectedandevaluated.The algorithm can obtain a maximum accuracy of about 82 percent and regression models are used to obtain such precision. The model can anticipate outcomes and is quickly adaptable to a wide range of inputs. Also, this strategy saves the banking industry and its staff a significant amount of time.

# References

- M.Sheikh,A.Goel,T.Kumar,"AnApproach for Prediction of Loan Approval using Machine Learning Algorithm," International Conference on Electronics and Sustainable Communication Systems (ICESC), (2020).
- [2] S. M S, R. Sunny T, "Loan Credibility Prediction System Based on Decision Tree Algorithm," International Journal of Engineering Research & Technology (IJERT) Vol. 4 Issue 09, (2015).
- [3] A. Kumar, I. Garg and S. Kaur, "Loan Approval Prediction based on Machine Learning Approach," IOSR Journal of Computer Engineering, (2016).
- [4] Dr K. Kavitha, "Clustering Loan Applicants based on Risk Percentage using K-Means ClusteringTechniques,"IJARCSSE-Volume 6, Issue 2, (2016).
- [5] P. Dutta, "A STUDY ON MACHINE LEARNING ALGORITHM FOR ENHANCEMENT OF LOAN PREDICTION",InternationalResearch

ISSN 2321-2152 www.ijmece.com

Vol 13, Issue 1, 2025



Journal of Modernization in Engineering Technology and Science, (2021).

- [6] G.Arutjothi,DrC.Senthamarai,"Predictionof Loan Status in Commercial Bank using Machine Learning Classifier," Proceedings of the International Conference on Intelligent Sustainable Systems, (2017).
- [7] P. Supriya, M. Pavani, N. Saisushma, N. Kumari and K. Vikas, "Loan Prediction by using Machine Learning Models," International Journal of Engineering and Techniques, (2019).
- [8] R. Salvi, R. Ghule, T. Sanadi, M. Bhajibhakare, "HOME LOAN DATA ANALYSISANDVISUALIZATION," International Journal of Creative Research Thoughts (IJCRT), (2021).
- [9] B.Srinivasan,N.Gnanasambandam,S.Zhao, R. Minhas, "Domain-specific adaptation of a partial least squaresregression model for loan defaults prediction," 11th IEEE International Conference on Data Mining Workshops, (2011).
- [10] M. V. Reddy, Dr B. Kavitha, "Neural NetworksforPredictionofLoanDefaultUsing Attribute Relevance Analysis," International Conference on Signal Acquisition and Processing, (2010).
- [11] G.Chornous, I.Nikolskyi, "Business-Oriented Feature Selection for Hybrid Classification Model of Credit Scoring," IEEE Second International Conference on Data Stream Mining & Processing August (2018).