# Predicting Credit Card Approval of Customers Through Customer Profiling using Machine Learning

**Arokiaraj Christian St Hubert, R. Vimalesh, M. Ranjith, S. Aravind Raj**

***Abstract**: In the banking sector, every banking infrastructure contains an enormous dataset for customers' credit card approval which requires customer profiling. The customer profiling means collection of data related to what customers need. It depends on customers' basic information like field of work, address proof, credit score, salary details, etc. This process mainly concentrates on predicting approval of credit cards to customers using machine learning. Machine Learning is the scientific study of algorithms and statistical models that computers use to perform specific tasks without any external instructions or interference. In the current trend this process is possible using many algorithms like "K-Mean, Improved K-Mean and Fuzzy C-Means". This helps banks to have an high profitability to satisfy their customers. However, the currently prevailing system shows an accuracy percentage of about 98.08%. The proposed system aims at improvising the accuracy ratio while using only few algorithms.*

*Keywords: Credit Score, Machine Learning, Supervised Learning, Unsupervised Learning.*

## I. INTRODUCTION

**M**achine Learning (ML) means a detailed study of algorithms and scientific models in a scientific manner. This ML is used by computer systems in order to perform a specific task without using outside instructions or explicit instructions. It relies on patterns and inference. Machine Learning is a subset of Artificial Intelligence (AI) which provides systems the capability to learn automatically and improve from experience without any explicit instruction. It gives importance for developing computer programs which can access data and use those data to learn for themselves. Real Life Examples include image recognition, speech recognition, medical diagnosis, prediction, financial services, etc. There are different types of learning in machine learning. They are:

- o Supervised Learning
- o Unsupervised Learning
- o Reinforcement Learning

**Arokiaraj Christian St Hubert**, Computer Science and Engineering, Sri Manakula Vinayagar Engineering College, Puducherry, India.

**R. Vimalesh**, Computer Science and Engineering, Sri Manakula Vinayagar Engineering College, Puducherry, India.

**M. Ranjith**, Computer Science and Engineering, Sri Manakula Vinayagar Engineering College, Puducherry, India.

**S. Aravind Raj**, Computer Science and Engineering, Sri Manakula Vinayagar Engineering College, Puducherry, India.

Supervised Learning means providing a full set of data collected from various sources while training an algorithm. Unsupervised Learning means providing the machine with a trained dataset without a specific needed outcome or any explicit instructions. Reinforcement Learning means providing the machine with a game like situation. The system puts forward trial and error to come up with a solution to the problem.

### A. Supervised Learning

Supervised learning is a task of machine learning where a function is being learnt by the machine which then maps an input to an output based on input-output pairs. This function is inferred from a labelled and trained data consisting of a set of training examples. This example examines the training data and produces an inferred function. This function can be used for mapping new examples. An optimal scenario is required for the algorithm to accurately determine the class labels in the case of unseen instances which in turn requires the learning algorithm to generalize from the training data in a reasonable manner. Supervised learning constructs a model which can predict based on evidence and proof even in the presence of uncertainty. Supervised Learning is divided into two types:

- ➢ **Classification**- This method separates the data and it provides a fixed output which may be a 'yes' or a 'no' or binary values such as '0' or '1'. Example: The working status of a person which can be either employed or not employed.
- ➢ **Regression**- This method fits the data and it gives continuous random values. Example: Prediction of a weather change.

### B. Classification

Classification is a technique where the response value can be predicted by separating the data into classes. This technique aims to reproduce class assignments. In order to extract models which describe important data classes or to predict future data trends, data analysis is used as two forms namely classification and prediction. Classification is a data mining technique in machine learning where group membership for data instances can be predicted.

**Examples are:**

- •Recognition of a type of car in a photo
- •Finding whether the mail is a mail spam or a message from a friend.

• Predicting the weather condition.

There are various classification algorithms. The classification algorithms are:

• Linear Classifiers
• Nearest Neighbor

• Support Vector Machines
• Decision Trees
• Boosted Trees
• Random Forest
• Neural Networks

## II. EXISTING SYSTEM

Credit card has evolved to a great level in banking industry. Each banking system consists of an enormous number of datasets to carry customer's transactions of their credit cards. So, banks would be in need of customer profiling. Customer Profiling [1], [2], [3], [11] in banks cognizes the issuer's decisions about whom to give banking facilities and what credit limit to be provided. It also helps the issuers to have a better understanding over their potential and current customers. In previous researches, profiling mainly depended on transaction data or demographic data, but in this research, both transaction and demographic data are merged in order to get more accurate results and minimize the possibility of risk occurrence.

By using the best techniques, it leads to improvement in accuracy and helps banks to have high profitability through customer satisfaction by focusing on the valuable customer (companies) which are considered as the main engine in the bank's profitability. This study used k-mean, improved k-mean, fuzzy c-means and neural networks. The used dataset is labeled and for neural network classification creating a new label as a target becomes the main aspect of this study, which helps to reduce the execution time of clustering process and provide the best results with accuracy. Finally, by comparing the accuracy ratio the results show that the neural network is the best clustering technique which could give an accuracy percentage of about 98.08%.
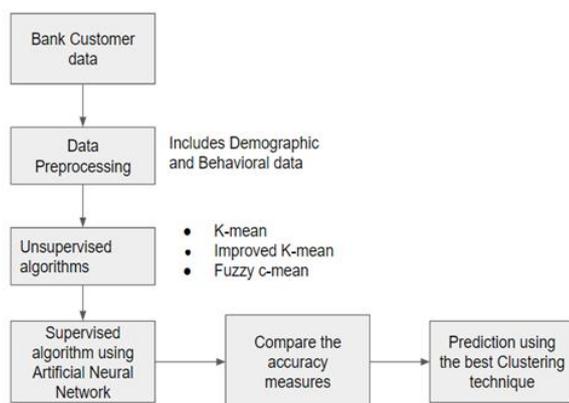


### Fig. 1. Architecture of the existing system

"Fig. 1" displays the architecture of the existing system. Initially, the customer data is collected from a reliable source. In the existing system, the dataset was collected from the archive of the University of California, Irvine. The dataset contains about 23 variables and there are no missing data in it. Then these data are preprocessed where both demographic and behavioral or the data used for transactions are considered. Then these data are trained using unsupervised algorithms namely k-mean, improved k-mean and fuzzy c-mean where each algorithms trains data in the form of customer segments. Then, again these data are trained using supervised technique namely artificial neural network. Finally, the results from all the algorithms are considered and compared and the prediction occurs using the best clustering technique. The results showed that the artificial neural network showed the highest accuracy than the other unsupervised algorithms by creating a new label target for the dataset.

The drawback is that the effectiveness and performance of this approach are yet to be improved by incorporating some deep learning algorithms especially in the field of medical informatics.

## III. PROPOSED SYSTEM

The proposed work aims at improving the current technology using algorithms like decision tree algorithm and k nearest neighbor algorithm. The proposed system is about the maintenance and analysis of customer profiles to approve credit cards. This process happens in an automated process. Automated process means it involves the assistance of machine learning. By manually collecting large sets of data from various customers who are working in various fields and already having credit cards, these data are taken into account and trained and tested in the machine. The trained and tested data are initially separated randomly through algorithms. It is then implemented using python programming language with the help of anaconda tool. The text editor used is Jupyter notebook. After all the trainings and testings of data, the machine will predict whatever data it is being presented at that time, based on customers' credit score, it will predict whether credit cards can be approved or not. These predictions happen through two algorithms namely:

(i) Decision Tree algorithm
(ii) K-Nearest Neighbor algorithm (KNN)
(iii) Logistic Regression

### A. Decision Tree Algorithm

Decision Tree algorithm is an algorithm in supervised learning. This algorithm can be used for solving both classification and regression problems unlike other supervised learning algorithms. The main aim of this tree is that it can create training models which will then be used to predict class or value of target variables by learning decision rules inferred from training data.

The understanding level of this algorithm is so easy compared to other classification algorithms. Using tree representation, this algorithm tries to solve the problem. The tree contains internal nodes which correspond to an attribute and leaf nodes which

correspond to a class label. So, for predicting a class label for a record, the process starts from the root of the tree, then the values of the root attribute are compared with the values of the record's attribute and finally on the basis of comparison, the branch corresponding to that value is followed and then jumped to the next node.

### B. K-Nearest Neighbor Algorithm

K – Nearest Neighbor (KNN) algorithm is a supervised machine learning algorithm which can be easily implemented and also can be used to solve both classification and regression problems. This algorithm assumes that similar things exist in close proximity or in other words near to each other. KNN captures this idea with some mathematics learned by everyone in childhood for calculating the distance between points on a graph. This idea is about similarity which can also be called distance, proximity or closeness. The distance can be calculated in other ways also and depending on the problem to be solved, any one of the ways can be adopted. One of the most known and familiar choices is the calculation of the straight-line distance which is called Euclidean distance.

### C. Logistic Regression

Logistic Regression is a technique applied in ML to classify problems based on the probability concepts. The logistic regression hypothesis limits the cost function between 0 and 1. In logistic regression, the graphs are plotted based on the features or variable used in the dataset. The independent variables are declared in the x-axis and the target or deciding variable is declared in the y-axis. This produces a graph which may be sinusoidal in nature or just a normal linear one.

Finally, the working of this process can be shown through by creating a web page like an online banking page where all the customer details are entered and the trained and tested customer data are stored in the form of database in the back end. Also, after all the customer details are entered and uploaded in the web page, the final process is the calculation of the credit score [4] of the customers.

## IV. SYSTEM DESIGN

There are six modules in total in the proposed system where the "Machine learning models" and "Training using various Classifiers" is a single module. The first five modules are implemented using several procedures in machine learning and final module is deployed in a web page using flask.
The data collected is sent and processed through each module like initially filling the missing data through different techniques which are discussed in detail in the modules portion. Then these data are split randomly with the help of a python library for sending them to training and testing purposes. After the data are trained and tested, they are sent to machine learning algorithms using some important variables to find out the accuracy and finally they are compared and the best one is used for future data analysis purpose. Finally, in the prediction module, the user's data is entered and the details are stored in the database. "Fig. 2"

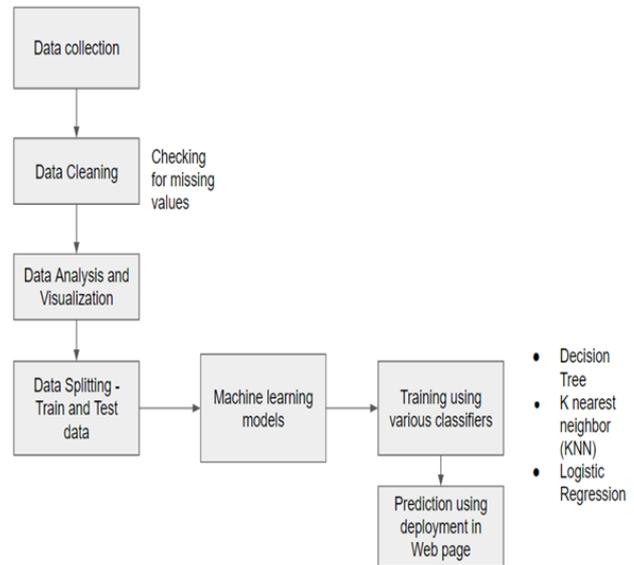portrays the architecture of the proposed system.



**Fig. 2. Architecture of the proposed system**

## V. MODULES DESCRIPTION

### A. Data Collection

Data about the customers should be collected. Collection of data can be done in two ways and it is shown in "Fig. 3".
(i) **Primary Data** - Data which has to be collected by the bank employees manually about the customers.
(ii) **Secondary Data** - Data which can be collected from other banks or other sources who might have already stored about the customers.



**Fig. 3. Data collection**

## B. Data Cleaning

```
gender = {'Male' : 1, 'Female' : 0}
# Making the actual convertion and replacing the values in the table
data['gender'] = data['gender'].map(gender)
data.head()
```

| Id | name | dob | district | state | pincode | phone number | email | gender |
|---|---|---|---|---|---|---|---|---|
| 1 | A Mohammed Irfan | 10/23/1993 | Chennai | Tamil nadu | 600005 | 9884839572 | irfantheoneandonly@gmail.com | 1 |
| 2 | Abirami D | 01/26/1992 | Madurai | Tamil nadu | 625018 | 8940525256 | sabari206@gmail.com | 0 |
| 3 | Aishi Neya S J | 03/23/1999 | Madurai | Tamil nadu | 625702 | 9789643264 | aamirthacable@gmail.com | 0 |
| 4 | Ajay A | 08/10/1998 | Trichy | Tamil nadu | 620001 | 7708955436 | ajaynirmal1998@outlook.com | 1 |
| 5 | Ajay Shankar S | 11/17/1995 | Vellore | Tamil nadu | 635810 | 9159568638 | ajayshankar1711@gmail.com | 1 |

**Fig. 4. Data cleaning**

The above "Fig. 4" represents the process of data cleaning. Data cleaning [5], [7] is done to fill the missing values. The missing values can be done in different ways.

(i)  The missing values can simply be entered as NAN (Not A Number).

(ii)  If the data does not contain important features such as an Aadhar card number or a PAN, that particular row or column can be deleted because without the unique data that decides the final decision, the process cannot be proceeded further.

(iii)  If the data contains missing feature such as an income of a customer, it can be filled by considering the other features of the customer and based on that, an average value will be calculated and filled.

(iv)  One hot encoding is a procedure where the data which can give only single decision (categorical data) are converted into numerical data in the form of 0's and 1's. For example: Credit card approval, gender of a person.

## C. Data Analysis and Visualization

The data to be analysed [10] is presented in the form of graphs.

The graph is plotted with two important deciding features plotted as X and Y axes against each other. For example: The credit score and the final decision are plotted against each other. If the credit score is above 700, credit card can be approved which is represented as '1' and if the credit score is not above 700, credit card cannot be approved which is represented as 0.
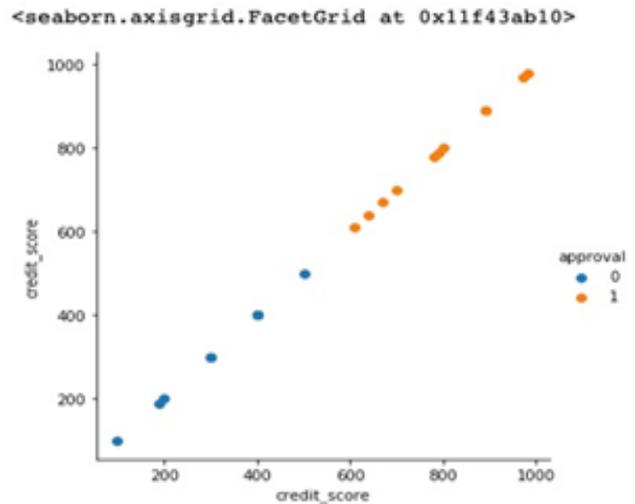


**Fig. 5. Data analysis and visualization**

The above illustration "Fig. 5" is for the credit score analysis where the dots in blue represent the data that are not approved and the orange dots represent the data that are approved.

## D. Data Splitting

Data Splitting means dividing a large set of data and sending some data for training and some data for testing. The training and testing Data can be split randomly based on a "Scikit-learn" library from Python which is called "Train_Test_Split library".

The data is split in the ratio of 8:2. For example: If there are 1000 data, 800 data will be used for training and the remaining data will be used for testing. Sometimes, the split ratio can change like 75% data can be used for training and the remaining 25% for testing, so it cannot be said accurately on how the data is divided respectively. This process is shown in "Fig. 6".

```
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(
        data.loc[:, ['cibil score','employed?']],
        data.loc[:, 'approval'])

print("{0:0.2f}% in training set".format((len(X_train)/len(data.index)) * 100))
print("{0:0.2f}% in test set".format((len(X_test)/len(data.index)) * 100))

74.93% in training set
25.07% in test set

X_train.head

<bound method NDFrame.head of       cibil score  employed?
Id
447           734          1
1030          825          1
844           625          1
923           833          1
617           654          1
...           ...        ...
849             0          0
102             0          0
384           718          1
135           830          1
510           800          1

[801 rows x 2 columns]>
```

**Fig. 6. Data splitting**

### E. Machine Learning Models

Machine learning [6] models include training the data under classification algorithms. Classification [8], [9] is an approach in supervised learning where computer programs learn from the input data provided to them and then uses them for classification of new observations. The classification algorithms used are: Decision Tree, K-Nearest Neighbor -and Logistic Regression.

### F. Prediction

After rigorous training and testing of sample data, it is now ready to predict any data given to the machine. The prediction will be made with the help of a web application using Flask. Flask is a Python based microframework which deploys machine learning with the web applications.

The user data entered in the web page will be stored in the database which is the backend. The backend is designed using XAMPP Server database importing MySQL. The front end is coded using HTML, CSS, Java Script and jQuery for validation purposes. The HTML, CSS and My SQL are interfaced through PHP (Hypertext Preprocessor). The final prediction will be provided in the web page whether to approve credit card or not which is shown in "Fig. 7".
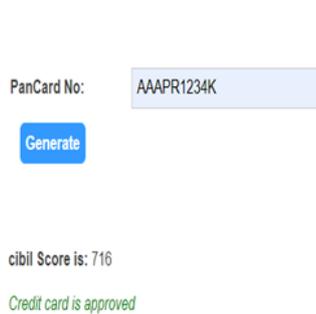


**Fig. 7. Prediction**

### VI. RESULTS

The trained and tested data are processed through three different algorithms for getting the best accuracy result. It shows that both the decision tree and k nearest neighbor algorithms gave the same training and testing accuracies of about 99.7% and 99.6% respectively since they used only limited variables for determining the final decision. However, the logistic regression algorithm showed a training accuracy of about 90.7% and a testing accuracy of about 91.4% even though it also used the same number of variables. Thus, the final accuracy ratio has been improved compared to that of the one in the existing system.
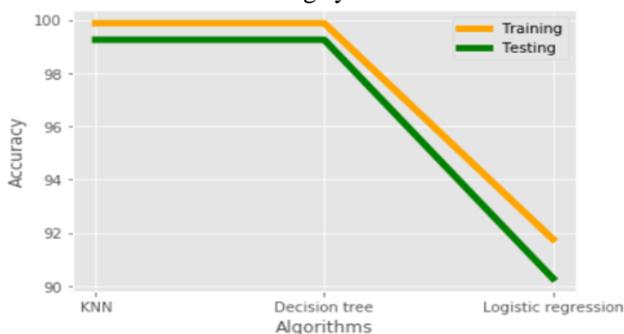


**Fig. 8. Comparison of accuracies of different algorithms**

A graph with all the algorithms in the x-axis and their corresponding accuracies in the y-axis is shown in "Fig. 8":

### VII. CONCLUSION

(i) The datasets of customers are completely collected, analyzed and trained.
(ii) These trained datasets serve as a helping factor in predicting the approval of credit cards for customers.
(iii) Both Decision Tree and KNN algorithms would provide good results after continuous training of different sets of collected data.
(iv) In real time when more dataset is trained and tested and when the variables for final decision are increased, both the decision and knn algorithms will show a change in their training and testing accuracies.
(v) The accuracy percentage improved compared to the existing system.

### REFERENCES

1. Dawood, E. A. E., Elfakhrany, E., & Maghraby, F. A. (2019). Improve profiling bank customer's behavior using machine learning. IEEE Access, 1–1.
2. Mondal, P. C., Deb, R., & Huda, M. N. (2016). Know your customer (KYC) based authentication method for financial services through the internet. 2016 19th International Conference on Computer and Information Technology.
3. Vaidya, A., & Diwakar, H. (2008). Customer profiling for business advantage using an Indian bank data. 2008 International Conference on Computer and Communication Engineering.
4. Huang, J., Wang, H., Wang, W., & Xiong, Z. (2013). A Computational Study for Feature Selection on Customer Credit Evaluation. 2013 IEEE International Conference on Systems, Man, and Cybernetics.
5. Karimov, J., Ozbayoglu, M., Tavli, B., & Dogdu, E. (2015). Generic menu optimization for multi-profile customer systems. 2015 IEEE International Symposium on Systems Engineering (ISSE).
6. S.-Schwartz, S. Ben-David, Understanding Machine Learning: From Theory to Algorithms, Cambridge, U.K.:Cambridge Univ. Press, 2014.
7. M. Ayoubi, "Customer segmentation based on CLV model and neural network", *Int. J. Comput. Sci.,* vol. 13, no. 2, pp. 31, 2016.
8. M. Sharahi, M. Aligholi, "Classify the data of bank customers using data mining and clustering techniques (case study: Sepah bank branches Tehran)", *J. Appl. Environ. Biol. Sci.,* vol. 5, pp. 458-464, 2015.
9. Aryuni, M., Didik Madyatmadja, E., & Miranda, E. (2018). Customer Segmentation in XYZ Bank Using K-Means and K-Medoids Clustering. 2018 International Conference on Information Management and Technology (ICIMTech).
10. Liu, K., & Liu, P. (2009). Visual Analysis of Customer Data in Commercial Banks. 2009 International Conference on Business Intelligence and Financial Engineering.
11. Kasa, N., Dahbura, A., Ravoori, C., & Adams, S. (2019). Improving Credit Card Fraud Detection by Profiling and Clustering Accounts. 2019 Systems and Information Engineering Design Symposium (SIEDS).

### AUTHORS PROFILE

**Mr. Arokiaraj Christian St Hubert** has completed both his Bachelor of Technology and Master of Technology in colleges that are affiliated to Pondicherry University, Puducherry, India. He is working as an Assistant Professor in the Department of Computer Science and Engineering at Sri Manakula Vinayagar Engineering College, Puducherry affiliated to Pondicherry University, Puducherry, India. He also has a membership in ISTE.

**Predicting Credit Card Approval of Customers Through Customer Profiling using Machine Learning**

**R. Vimalesh** is pursuing Bachelor of Technology in the stream of Computer Science and Engineering at Sri Manakula Vinayagar Engineering College, Puducherry affiliated to Pondicherry University, Puducherry, India. His research interests include in the field of Machine Learning, Database management system and Web page development.

**M. Ranjith** is pursuing Bachelor of Technology in the stream of Computer Science and Engineering at Sri Manakula Vinayagar Engineering College, Puducherry affiliated to Pondicherry University, Puducherry, India. His research interests include in the field of Database management system and Web page development.

**S. Aravind Raj** is pursuing Bachelor of Technology in the stream of Computer Science and Engineering at Sri Manakula Vinayagar Engineering College, Puducherry affiliated to Pondicherry University, Puducherry, India. His research interests include in the field of Database management system and Web page development.

557