



Machine learning models for Predicting Bank loan eligibility

By Artificial neural network

G. Siva Kumari¹, I. Susmitha², D. Ashok Kumar³, A. Shanmukh

Sivananda Reddy⁴, Mrs. Venu Kumari⁵

Summary

Artificial Neural Networks play an increasingly important role in financial applications for such tasks as pattern recognition, classification, and time series forecasting. This study develops a proposed model that identifies artificial neural network as an enabling tool for evaluating credit applications to support loan decisions in the Jordanian commercial banks. A multi-layer feed-forward neural network with back propagation learning algorithm was used to build up the proposed model. Different representative cases of loan applications were considered based on the guidelines of different banks in Jordan, to validate the neural network model. The results indicate that artificial neural networks are a successful technology that can be used in loan application evaluation in the Jordanian commercial banks.

Key words: *Personal Loans, Interpreting Model, Classification Model, Artificial Neural Network, Predicting Model, Feature Selection*

1. Introduction

Artificial neural networks have been fruitfully used in a variety of business fields including

marketing, accounting, management information systems, and production management. Most of the studies have used neural networks for predicting future stock behavior, financial crises, bankruptcy, exchange rate, and detecting credit card fraud. The granting of loans by banks is one of the key areas concerning decision problems that need subtle care. Neural Networks have successfully provided effective credit evaluations for supporting granting loans. Researchers are currently focusing on using neural network classification models and particularly backpropagation neural networks in classifying loan applications into good and bad ones. This research attempts to explore whether using neural networks will provide more accurate personal loan decisions in the Jordanian commercial banks. It will also propose a neural computing model as the basis for a decision support tool in granting or rejecting bank loans. The model will apply banking credit standards to determine the customer who will be eligible for credit approval.

Despite the increase in consumer loans defaults and competition in the banking market, most of the Jordanian commercial banks are reluctant to use artificial intelligence software technologies in their decision-making routines.



Generally, bank loan officers rely on traditional methods to guide them in evaluating the worthiness of loan applications. Generally, loan applications evaluations are based on

a loan officers' subjective assessment. Such judgment is inefficient, inconsistent, and non-uniform. Therefore, a knowledge discovery tool is needed to assist in decision making regarding the application. Furthermore, the complexity of loan decision tools and variation between applications is an opportunity for a neural-computing technology to provide learning capability that does not exist in other technologies.

Neural networks with their capability of capturing nonlinear and complex relationships are a powerful alternative to the conventional forecasting and classification methods. Neural networks are consistent paradigms of the nonparametric approach in financial modeling due to their ability to correctly classify and predict consumer loan defaults.

1.1. Literature Review:

We studied various research papers on loan prediction models. [1] A research paper by G. Arutjothi and Dr. C. Senthamaria explained that predicting credit defaulters is a complex task so there is a need for a machine learning model for this to save time and resources. Using the R software they have proposed that the combination of Min-Max normalization and K-Nearest Neighbor (K-NN) classifier will be good to accurately predict loan approvals. [2]

Aboobyda and Tagir from University of Khartoum, Sudan used j48, bayesNet and naiveBayes algorithms for loan prediction and concluded that j48 would be best for the accurate prediction of credit approvals. They have used Weka application for implementation and testing of the model and compared the results of all the three algorithms.[3] A research paper by Kumar Arun, Garg Ishan and Kaur Sanmeet explained the use of various classification models such as Random Forest, SVM, LM, Nnet and ADB for the prediction of the loan approvals.[4] Glorfeld and Hardgrave had projected a high-performance model with optimum design for the use of neural network thus estimating the credit value of the applications of loans. 75% of the loan applicants were correctly predicted by their designed model. [5] Andy Liaw and Matthew Wiener in their research paper told about the classification and regression by Random Forest. [6] Stephan Dreiseitl and Lucila Ohno-Machado told about the artificial neural network classification and logistic regression models, How logical regression is useful in building various systems for predictions. Machine learning predictive models are also a good choice for loan predictions in this market scenario[7][8][9].

Wang et al. (2023) introduced a stacking-based model to approve financial institution risks, selecting the best model by comparing performance. They also built a bank approval model using deep learning on imbalanced data, utilizing CNN for feature extraction and counterfactual augmentation for balanced sampling. Optimizing the auto finance prediction



model based on bank model features led to around a 6% increase in joint loan approval, as demonstrated in experiments on real data.[10]

1.2. Research Problem:

Lending markets have rapidly expanded over the 20th century, and managerial decisions have become increasingly complex, requiring sophisticated decision support systems based on models beyond simple linear relationships to meet the increasing demand from population expansion. AI has made significant progress in the past few years, which has led to the creation of various professional financing applications. AI is widely believed to eventually replace humans in the finance industry. Financial institutions need to balance increasing liquidity and reducing NPL. It is a trade-off, as the growth of either metric is never a good sign. Loans are based on bank guidelines and regulations, standard analyses, scoring systems, and expert opinions. These procedures are tedious and time-consuming, and they can result in biased decisions. As a result, established banks and startups constantly seek ways to innovate, deal with nonlinear data, and discover valuable relationships or patterns in a large set of historical data. Artificial neural networks (ANNs) might be the best solution because of their unique learning capabilities, which are unavailable with other technologies. Although ANNs are powerful prediction tools, they lack the meaningful interpretations that traditional intelligence models can provide. It is

because their black-box nature makes them challenging to interpret. For instance, when running a computation on a network, the program might end up with a different weight matrix. This is because the random values are changed as the program executes. Additionally, it is unclear how ANNs operate.

1.3. Objective:

This study aims to develop a loan decision support system using ANNs to assist users (e.g., loan officers, borrowers, higher management, and financial instructors) in predicting and providing meaningful personal loan decisions.

The general objectives of this study are as follows:

- To predict the loan eligibility of borrowers using an ANN classification model.

2. Dataset:

The dataset used in this study is the historical dataset 'Loan Eligible Dataset,' available on Kaggle [7] and licensed under Database Contents License (DbCL) v1.0. Table 1 below gives a brief description of the dataset attributes.

Table 1: Dataset description

Variable Name	Description	Data Type
Loan_ID	Loan reference number (Unique I.D.)	Numeric
Gender	Applicant gender	Categorical
Married	Applicant marital status	Categorical
Dependents	Number of family members	Numeric
Education	Applicant educational qualification (graduate or not graduate)	Categorical
Self_Employed	Applicant employment status (yes for self-employed, no for employed/others)	Categorical
Applicant_Income	Applicant's monthly salary/income	Numeric
Coapplicant_Income	Additional applicant's monthly salary/income	Numeric
Loan_Amount	Loan amount	Numeric
Loan_Amount_Term	The loan's repayment period (in days)	Numeric
Credit_History	Records of applicant's credit history (0: bad credit history, 1: good credit history)	Numeric
Property_Area	The location of the applicant's home (Rural/Semi-urban/Urban)	Categorical
Loan_Status	Status of loan (Y: accepted, N: not accepted)	Categorical

3.1 Logistic Regression:

Logistic regression is a popular statistical method used for binary classification tasks in machine learning. Despite its name, it is actually a classification algorithm rather than a regression algorithm. In logistic regression, the dependent variable is binary in nature 1 refers to true and 0 refers to false. Goal is to find the best fitting model for independent and dependent variable relationship. Independent variable can be continuous or binary. It is also called Logistic Regression Logistic regression is used in machine learning

3.2 Decision Trees:

One of the most easy and famous classification algorithms is Decision Tree Algorithm. This algorithm helps interpreting and understanding better. Decision tree algorithms are one of the supervised learning algorithms. The decision tree algorithm is capable of solving both classification and regression problems, which distinguishes it from the rest of the supervised learning algorithms.

3.3 Random Forest:

Ensemble Learning is an approach which tells not to be reliant only on one model to make prediction. Rather take into consideration a number of models, and on the basis of outputs of all such models, come to a conclusion. The prediction made using this approach is far more accurate than it would have been considering only one model for predicting. Random forest is kind of an ensemble classifier which is using decision tree algorithm in a randomized fashion. Similarly, generate a new bootstrap data set for each decision tree to be built. Build a number of decision trees in the same way using the subset of total variables present in the bootstrap data set. The prediction will definitely be far more accurate than it would have been using only one decision tree.

3.4 KNN:

KNN is a rule which learns by memorizing. This algorithm prerequisites containing of the training data. The neighbors are found at the time of testing of the already stored training data.



In KNN algorithm, Euclidean distance is measured between the training and test data. Euclidean distance is calculated as square root of summation of difference between the data to be tested and training data. Suppose we are provided with a data set having a number of values. Let us say we have are given with a value whose category is to be tested. We are also provided with a value K, which tells us the number of neighbours¹ which are closest to the test value. In order to do so, Euclidean distance between all the similar values from the table are to be calculated. K values having least distance with the test value are considered and their category is checked. Checking the categories of K nearest values, the value of the test data is predicted.

3.5 Support vector machine:

SVMs are commonly used within classification problems. They distinguish between two classes by finding the optimal hyperplane that maximizes the margin between the closest data points of opposite classes. The number of features in the input data determine if the hyperplane is a line in a 2-D space or a plane in a n-dimensional space. Since multiple hyperplanes can be found to differentiate classes, maximizing the margin between points enables the algorithm to find the best decision boundary between classes. This, in turn, enables it to generalize well to new data and make accurate classification predictions. The lines that are adjacent to the optimal hyperplane

are known as support vectors as these vectors run through the data points that determine the maximal margin.

4. Artificial Neural Network Foundations and Techniques:

An ANN is a computing system inspired by the biological networks that govern animal brains. The first known cases of ANNs can be traced back to the previous work. Their model established the basis for studying biological processes and their link to AI. An ANN is a system that learns and collects knowledge by studying patterns and relationships in data. It is comprised of a large number of single units or processing elements (PEE) that are connected to weights. There are three main types of neural networks: recurrent neural networks (RNNs), ANNs, and convolution neural networks (CNN). Lately, ANN models have become famous for exploring the general relationship patterns between linear and nonlinear input-outputs to evaluate testing sets. Successful prediction is a measurement of ANN performance.

ANNs are typically used to solve predictive tasks, including classification, regression, and time series analysis. They are 6 popularly known as universal function approximates, and many types exist. The most popular methods of ANNs are generalized regression neural networks (GRNNs), multilayer perceptron (MLP), and artificial neural fuzzy inference systems (ANFISs), which are combinations of fuzzy logic decision systems and ANNs. Most studies employing ANN tools have demonstrated their capacity to represent noisy,

incomplete, and inconsistent data. With their ability to learn from previous experiences and provide an improved output when taught using newer examples, ANNs can minimize the probability of biased choices. It is comprised of a large number of discrete PEE that are connected to weights. Recurrent neural networks (RNNs), convolution neural networks (CNNs), and ANNs are the three main varieties of neural networks. This reduces subjectivity in decision-making since they can manage linear and nonlinear relationships between input and output variables. ANNs are self-adaptive, and completing their objective requires a representative training set. However, ANN failure to explain the reason behind the results causes users not to trust their output. Thus, the training data must follow a universal set of parameters, indicating the importance of data preparation before modeling.

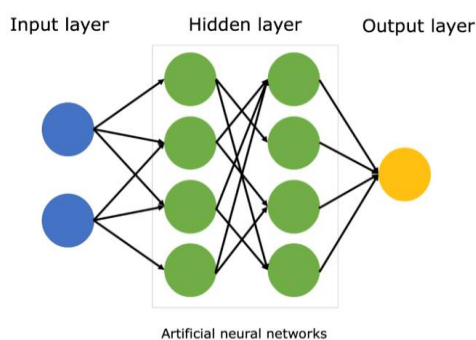


Fig 1 Artificial neural network

Some networks want a neuron to inhibit the other neurons in the same layer. This is called lateral inhibition. The most common use of this is in the output layer. For example, in text recognition if the probability of a character

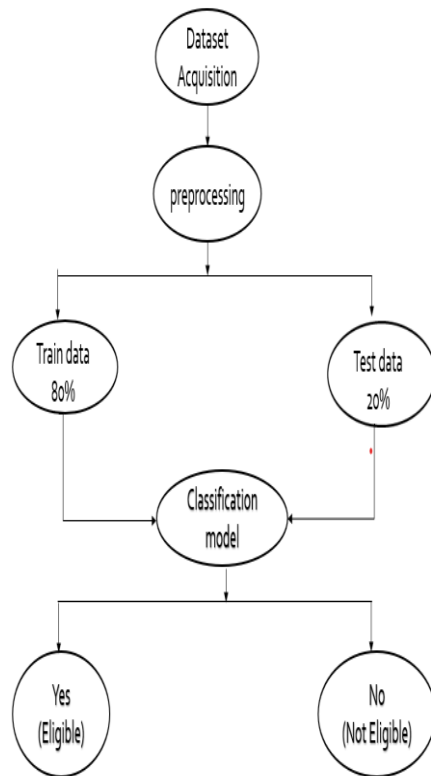
being a "P" is .85 and the probability of the character being an "F" is .65, the network wants to choose the highest probability and inhibit all the others. It can do that with lateral inhibition. This concept is also called competition. Another type of connection is feedback. This is where the output of one-layer routes back to a previous layer. An example of this is shown in Figure 2.

4.1 Training an Artificial Neural Network:

Once a network has been structured for a particular application, that network is ready to be trained. To start this process the initial weights are chosen randomly. Then, the training, or learning, begins. There are two approaches to training - supervised and unsupervised. Supervised training involves a mechanism of providing the network with the desired output either by manually "grading" the network's performance or by providing the desired outputs with the inputs. Unsupervised training is where the network has to make sense of the inputs without outside help. The vast bulk of networks utilize supervised training. Unsupervised training is used to perform some initial characterization inputs. However, in the full-blown sense of being truly self-learning, it is still just a shining promise that is not fully understood, does not completely work, and thus is relegated to the lab.

4.2 Applications of Artificial Neural Networks:

Artificial neural networks (ANNs) find several applications in determining bank loan eligibility due to their ability to learn complex patterns and relationships from data. Here's how ANNs can be applied in this context:



Credit Scoring: ANNs can analyze historical data on loan applicants, including their credit history, income, employment status, and other relevant factors, to predict the likelihood of default or timely repayment. This helps banks assess the creditworthiness of applicants more accurately.

Risk Assessment: By analyzing a wide range of variables, including economic indicators and market trends, ANNs can assist banks in assessing the overall risk associated with lending to certain individuals or businesses. This helps in making more informed decisions about loan approvals and setting interest rates.

Fraud Detection: ANNs can be trained to detect fraudulent loan applications by identifying suspicious patterns or anomalies in

the data. This helps banks minimize losses due to fraudulent activities and maintain the integrity of their lending processes.

Customer Segmentation: ANNs can segment customers based on their financial behavior, risk profile, and other characteristics, allowing banks to tailor their loan products and marketing strategies to different segments more effectively.

Automated Decision Making: ANNs can be integrated into automated loan approval systems, where they analyze loan applications in real-time and make decisions on whether to approve or reject them based on predefined criteria. This streamlines the loan approval process and improves efficiency.

Personalized Recommendations: By analyzing customer data and preferences, ANNs can generate personalized loan recommendations, such as suggesting the most suitable loan products or terms for individual customers based on their financial situation and goals.

Overall, ANNs offer powerful capabilities for banks to enhance their loan approval processes, mitigate risks, improve customer experience, and make more data-driven decisions in lending operations.

5. Proposed Neural Network Model for Loan Decisions:

Designing a neural network successfully relies on a clear understanding of the problem, and on deciding upon most influential input variables. The procedure of designing a neural network model is a logical process. This research followed the step designing methodology presented in figure (1). The process was not a single-pass one, but it

required going back to previous steps several times

5.1 Banking Data Sets:

The data set used in this research was divided into training and testing data sets. There were 94 cases used in the training and 46 in the testing. Both training and testing data sets contained half-good applications and half-bad applications. The training set cases approximately cover the input data space. All training cases were set by default taking into account the banks' guidelines for personal credit approval in the Jordanian banks.

5.2 Specifying Loan Variables:

In relation to loan applications evaluation, a set of decision variables that determine the credit worthiness of an application were used in this current research. These variables were taken from the guidelines that credit officers use in Jordanian commercial banks. The information that is considered as significant includes: the applicant's nationality, residency, companies' type, guarantor, job experience, and the DBR (i.e. debt balance ratio that measures the applicant's repaying ability). High DBR ratio indicates high credit risk. On the other hand, low DBR ratio indicates a good credit application. Therefore, there are 13 input neurons for the network; each represents an independent variable. They are referred to as respectively. On the other hand, the output in this research is 1 for good application or 0 for

bad application; a single output neuron is needed to produce the output.

Defining the Network Parameters:

Deciding on the optimal parameters for the MLFN model is a critical issue. The best model is one that has the combination of parameters that minimizes the mean squared error. At the beginning, the researcher selected an initial configuration (one hidden layer with a number of hidden neurons, and one neuron in the hidden layer). After conducting a number of experiments with each model, the researcher retained the best network. In brief, a series of repeated trials were applied in order to reach the optimal set of parameters. Furthermore, the new set of trials increased the number of hidden layers to two and three hidden layers with a different number of neuron combinations in order to decide on the best set of hidden layers and neurons. A multi-layer feed-forward neural network with two hidden layers with ten neurons in the first hidden and seven neurons in the second hidden had the best speed time at the fixed threshold level of 0.05. Table (1) illustrates the optimal set of the proposed Neural Network parameters.

Deciding on the optimal parameters for the MLFN model is a critical issue. The best model is one that has the combination of parameters that minimizes the mean squared error. At the beginning, the researcher selected an initial configuration (one hidden layer with a number of hidden neurons, and one neuron in the hidden layer).

After conducting a number of experiments



with each model, the researcher retained the best network. In brief, a series of repeated trials were applied in order to reach the optimal set of parameters. Furthermore, the new set of trials increased the number of hidden layers to two and three hidden layers with a different number of neuron combinations in order to decide on the best set of hidden layers and neurons. A multi-layer feed-forward neural network with two hidden layers with ten neurons in the first hidden and seven neurons in the second hidden had the best speed time at the fixed threshold level of 0.05. Table (1) illustrates the optimal set of the proposed Neural Network parameters.

Data acquisition:

The process of collecting and gathering data from various sources for analysis or processing. In the context of bank loan eligibility, data acquisition involves gathering relevant information about loan applicants to assess their creditworthiness and make informed lending decisions.

Processing:

Once data is acquired for bank loan eligibility assessment, it undergoes several processing steps to extract valuable insights and make informed lending decisions.

Splitting Data:

Splitting data into training and testing categories is critical for reducing the risk of

bias in the assessment and validation process. The model is trained or fitted using the training set; however, the testing set is required to evaluate the final model's performance objectively.

Train Data:

Training data refers to the subset of data that is used to train a machine learning model. In the context of bank loan eligibility assessment, training data typically consists of historical data on loan applicants for whom the outcomes (whether the loan was approved or denied) are known. First the networks' weights are initialized, consequently the network is prepared for training.

The training cases are used to adjust the weights through minimizing the prediction made by the network. The back-propagation gradient descent-learning algorithm uses the performance function in order to trace the best set of weights that minimizes the average mean squared error. The algorithm allocates the error backward through the network's layers successively. Furthermore, the algorithm develops through a number of epochs and uses the error to adjust the weights in the direction in which the performance function decreases quickly.

Test Data:

Test data, also known as validation data or evaluation data, is a separate subset of data that is used to assess the performance of a machine learning model after it has been trained on the training data. minimizes the average mean squared error. The algorithm allocates the error backward

through the network's layers successively.

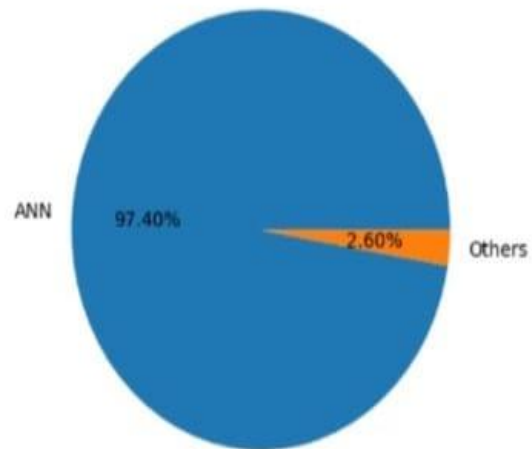
Testing the network is necessary in order to examine its ability to classify the testing data correctly. Testing starts after the training has been completed. The network was simulated on the testing set (i.e. cases the network has not seen before). The results were very good; the network was able to classify 95% of the cases in the testing set. The results are shown in the appendix. A pair wise t-test was used to compare between the actual neural network output and the target output of the testing set. The null hypothesis.

Classification Model:

The context of bank loan eligibility assessment, a classification model is a type of machine learning model that predicts whether a loan applicant should be approved or denied based on various input features. Classification models are trained on historical data where the outcome (loan approval yes or no) is known, and they learn to classify new applicants into one of the predefined classes (e.g., approved or not approved) based on patterns and relationships in the data.

Result:

This section presents the results obtained from data preparation and feature selection. It also explains the algorithms for developing and evaluating the proposed models' performance.



Discussion:

This study used historical data on loan repayment statuses and interest rates of prior borrowers to create AI models, in contrast to most previous studies that relied on guidelines or grading systems to make predictions and choices.

According to quantitative and qualitative analysis, it can be deduced that the loan amount, annual income, evaluation of credit risk, debt-to-income ratio, number of satisfactory accounts, home ownership, loan purpose, repayment term, duration of employment, and the sum of the current balances of all accounts are crucial variables to take into account when designing loan eligibility and interest rate models. Although the accuracy of loan eligibility prediction using a classification neural network model was weak, the confusion matrix indicates that neural networks can help loan officers make better informed application decisions, thus minimizing NPLs and financial losses. Additionally, providing the reasons behind the results of an AI model fosters trust in the model's decision for loan officers and borrowers. Moreover,

determining the appropriate interest rate for the borrower minimizes the risk that the borrower may not repay the loan.

QUALITATIVE ANALYSIS:

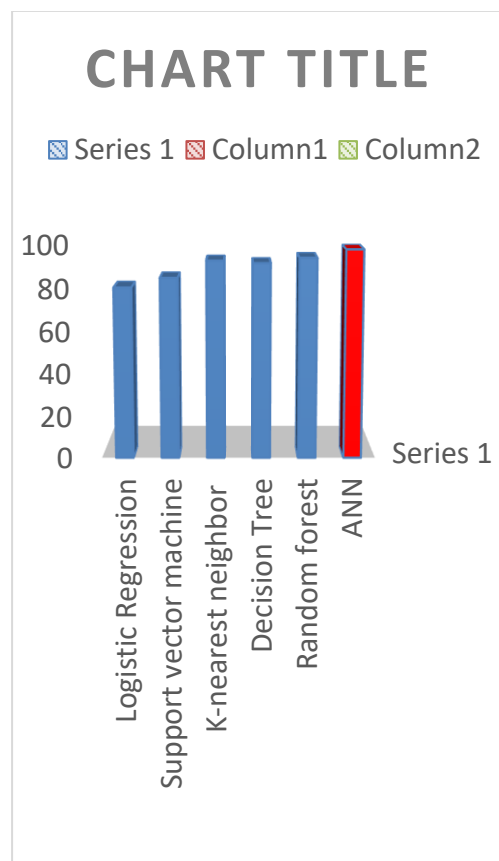
The fig 4 describes all the various features in the dataset and how many people are eligible and how many persons are not eligible can be clearly seen in this graph.



Fig 5 Bar graph of algorithms

QUANTITATIVE ANALYSIS:

Table 2: Tables for accuracy of algorithms



Conclusion:

ALGORITHM	ACCURACY
Logistic Regression	80%
Support vector machine	84.4%
K-nearest neighbor	92.3%
Decision Tree	91.11%
Random Forest	93.56%
Artificial Neural network	97.40%

This research project developed a decision support tool to aid financial institutions in determining both eligibility and interest rates for personal loans. The main objective is to help financial institutions make well-informed decisions and to advise potential borrowers about the dangers that can be involved with the loans they choose. The system divides users into eligible groups and those who are not, then computes appropriate interest rates for those judged suitable. This research explores several frequently ignored aspects of predicting loan

eligibility, such as model interpretability and addressing issues related to class imbalance, in contrast to many prior studies that heavily relied on predetermined guidelines or scoring systems.

REFERENCES:

- [1] G. Arutjothi and Dr. C. Senthamaria. "Prediction of Loan Status in Commercial Bank using Machine Learning Classifier", International Conference on Intelligent Sustainable Systems (ICISS2017). doi:10.1109/iss1.2017.8389442.
- [2] Aboobyda Jafar Hamid and Tarig Mohammed Ahmed. "DEVELOPING PREDICTION MODEL OF LOAN RISK IN BANKS USING DATA MINING", Machine Learning and Applications: An International Journal (MLAIJ) Vol.3, No.1, March 2016.
- [3] Kumar Arun, Garg Ishan, Kaur Sanmeet. "Loan Approval Prediction based on Machine Learning Approach", National Conference on Recent Trends in Computer Science and Information Technology (NCERT CSIT-2016).
- [4] Louis W. Glorfeld and Bill C. Hardgrave. "An improved method for developing neural networks: The case of evaluating commercial loan creditworthiness", Computers & Operations Research, Volume 23, Issue 10, October 1996.
- [5] Andy Liaw and Matthew Wiener. "Classification and Regression by randomForest", ISSN 1609-3631, Vol. 2/3, December 2002.
- [6] Stephan Dreiseitl and Lucila Ohno-Machado. "Logistic regression and artificial neural network classification models: a methodology review", Journal of Biomedical Informatics, Volume 35, Issues 5–6, October 2002.
- [7] T. Choudhury, G. Dangi, T. P. Singh, A. Chauhan and A. Aggarwal, "An Efficient Way to Detect Credit Card Fraud Using Machine Learning Methodologies," 2018 Second International Conference on Green Computing and Internet of Things (ICGCIoT), 2018, pp. 591-597, doi: 10.1109/ICGCIoT.2018.8753077.
- [8] S. Taneja, D. Garg, M. V. Tarun Kumar and T. Choudhury, "The Machine Predicted Market," 2018 International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS), 2018, pp. 256-260, doi: 10.1109/CTEMS.2018.8769306.
- [9] Sharma, H.K., Choudhury, T., Toe, T.T. (2022). Machine Learning Based Predictive Analytics: A Use Case in Insurance Sector. In: Jeyanthi, P.M., Choudhury, T., Hack-Polay, D., Singh, T.P., Abujar, S. (eds) Decision Intelligence Analytics and the Implementation of Strategic Business Management. EAI/Springer Innovations in Communication and Computing. Springer, Cham. https://doi.org/10.1007/978-3-030-82763-2_14.
- [10] Koulouridi, E., Kumar, S., Nario, L., Pepanides, T., & Vettori, M. (2021). Managing and monitoring credit risk after the Covid-19 pandemic. McKinsey & Co., available at <https://www.mckinsey.com/business-functions/risk/our-insights/managing-and-monitoring-credit-risk-after-the-covid-19-pandemic#>. (Accessed 29 December 2021).

About Authors:



Mrs. P. Venu Kumari is Currently working as Assistant Professor in Tirumala Engineering College. She received her M.Tech (VLSI) Degree from Jawaharlal Nehru Technological University Hyderabad. She is a life member of technical association in IETE.



G. siva kumari currently studying B.Tech (Electronics & Communication Engineering) in Tirumala Engineering College, Jawaharla Nehru Technological University Kakinada, Andhra Pradesh in the year 2024. She completed her intermediate in **Morning star College**.



I. Susmitha currently studying B.Tech (Electronics & Communication Engineering) in Tirumala Engineering College, Jawaharlal Nehru Technological University Kakinada, Andhra Pradesh in the year 2024. She completed her intermediate in Bhavana Junior College.



D. Venkata Ashok Kumari currently studying B.Tech (Electronics & Communication Engineering) in Tirumala Engineering College, Jawaharlal Nehru Technological University Kakinada, Andhra Pradesh in the year 2024. He completed his intermediate in Vagdevi Junior College.



A. Shanmukha Sivananda Reddy currently studying B.Tech (Electronics & Communication Engineering) in Tirumala Engineering College, Jawaharlal Nehru Technological University Kakinada, Andhra Pradesh in the year 2024. He completed his intermediate in Jupiter Junior College.