

## BANKRUPTCY PREDICTION OF INDIAN BANKS USING ADVANCED ANALYTICS<sup>3</sup>

*The banking sector in India plays a crucial role in economic growth. A bank provides an opportunity for investments to encourage economic growth and the potential to yield higher returns. In this study, we develop a bankruptcy prediction model by using machine learning (ML) techniques, namely logistic regression, random forest, and AdaBoost, and compare these models with those developed using deep learning (DL) techniques, namely the artificial neural network (ANN). ANN results in the highest accuracy and the most favourable prediction model for bankruptcy. Data used in this study are collected from survived and failed private and public sector banks from India from March 2001 to March 2018. For bankruptcy prediction, we use the bank's macroeconomic and market structure-related features. The feature selection technique 'Relief algorithm' is used to select useful features for the bankruptcy prediction model. Because failed banks in comparison with survived were less in the dataset, the issue of imbalanced cases may have arisen, in which case most ML and DL techniques do not perform well. Thus, we convert the dataset into a balanced form by using the synthetic minority oversampling technique (SMOTE). The results of this study can help in performing financial analyses of banks and thus have significant implications for their stakeholders.*

*Keywords: Bankruptcy; Imbalanced Data; SMOTE; Relief Algorithm; Deep Learning; Artificial Neural Network*

*JEL: G21; G28; G34*

### 1. Introduction

Being a central player within a nation's economy, banks control the supply of money in circulation and stimulus. The banking sector in India is sufficiently capitalized and well-regulated by the Reserve Bank of India (RBI), which is India's central banking agency<sup>4</sup>. In recent times, the banking sector in India has introduced innovative models such as payments banks and small finance banks. Banking institutions worldwide have been undergoing dynamic changes where their survival depends on the quality of services they offer to their customers. These financial organizations are encountering stifling competition because of

<sup>1</sup> Institute of Management Technology Nagpur, sarbjitoberoi@gmail.com.

<sup>2</sup> Institute of Management Technology Nagpur, sayansir@gmail.com.

<sup>3</sup> This paper should be cited as: Oberoi, S. S., Banerjee, S. (2023). Bankruptcy Prediction of Indian Banks Using Advanced Analytics. – *Economic Studies (Ikonicheski Izsledvania)*, 32(4), pp. 22-41.

<sup>4</sup> The RBI is India's central bank that controls the flow and supply of the Indian rupee.

increasing consumer demands, rapid growth in technological infrastructure, and continual changes in banking regulations and policies. The success of these financial organizations largely depends on how they leverage resources such as technological infrastructure, the quality of services they offer to their customers, and governing policies.

Bankruptcy forecasting is a crucial concern in the financial sector and has attracted increasing attention from both academic researchers and industry practitioners. Because of the burgeoning development in the power of computing, researchers have even attempted to use machine learning (ML) and deep learning to forecast the plaguing challenge of ‘bankruptcy’ (Altman, et al., 2020). Most banks and financial organisations still prefer using a traditional technique to evaluate their performance (Qu, et al., 2019; Altman et al., 2020). However, the methodological limitations of these techniques and approaches should be considered (Qu, et al., 2019).

Bankruptcy prediction has been in vogue in the research community for approximately five decades now (Makinen, Solanko, 2018). Bankruptcy prediction remains a vital factor because it helps measure the financial health of a firm before it becomes bankrupt. Studies have widely used statistical (SL) and ML techniques to build bankruptcy prediction models (Appiah, 2015). To develop unbiased and generalised bankruptcy prediction models, specific features that can effectively describe the status of a bank should be selected (Liang et al. 2016). We attempt to build an efficient bankruptcy prediction model that can solve the issue of ‘imbalanced classes’.

We initially assume that one of the following conditions triggers the failure of a bank: dissolution, negative total assets, state intervention, and merger and acquisition (Pappas, et al., 2017). We collect data from 58 Indian public and private sector banks that have been categorized as ‘failed’ or ‘survived’ as per conditions indicated by Pappas et al. (2017). Because of an imbalanced dataset, we use the synthetic minority oversampling technique (SMOTE) method to transform data in the balanced form (Fernández, et al., 2018). Moreover, a relief algorithm is used to select crucial features for bankruptcy prediction. These selected features are then fed into different ML and DL techniques to develop the most efficient and generalised bankruptcy prediction model. We randomly divide the whole dataset into training and test datasets, accounting for 80% and 20% of the data, respectively. A different bankruptcy prediction model is developed using ML techniques, such as logistic regression, random forest, and Adaboost, and DL techniques. Finally, all these results are compared based on the model accuracy to derive the best-generalized bankruptcy prediction models.

The remaining study is organized as follows. Section 2 provides the literature review. Section 3 describes the data, descriptive statistics, and methodology used to build the model. Section 4 explains empirical results. Section 5 provides implications of the research study.

## **2. Literature Review**

This study focuses on bankruptcy forecasting which has been a trending topic in recent times. Statistical techniques have been mainly used for bankruptcy forecasting (Qu, et al., 2019). Both academia and industry practitioners have been using advanced techniques such as ML

ad DL algorithms to formulate a bankruptcy prediction model (Nwogugu, 2006; Nwogugu, 2008; Dellepiane, et al., 2015; Kadioglu, et al., 2017; Barboza, et al., 2017; Sujud, Hashem, 2017; Kou, et al., 2019; Devi, Radhika, 2018; Qu, et al., 2019). This literature review focused on two topics: ML and DL approaches.

### *2.1 ML Approaches*

Park and Han (2000) were one of the first researchers who developed a bankruptcy prediction model by using the k-nearest neighbour. Furthermore, Min and Lee (2005) are the first to use a support vector machine (SVM) with various kernels for building a bankruptcy prediction model. Boyacioglu et al. (2009) developed a bankruptcy prediction model for 65 Turkish banks by using SVM and multivariate statistical methods. The accuracy of the prediction model developed using SVM is superior to those of other models. A study built nine bankruptcy prediction models by using ML techniques, such as logistic regression (LR), SVM, K-nearest neighbour, and linear discriminant analysis, for US banks during the financial crisis and found that the accuracy of the model developed using SVM was higher than that of models developed using other ML techniques (Serrano-Cinca and Gutiérrez-Nieto, 2013). Chiamonte et al. (2015) formulated bankruptcy prediction models for 3242 European banks and showed that the neural network yielded more favourable results than did other techniques. The findings of the aforementioned studies indicated that SVM is the most suitable ML technique for developing bankruptcy prediction models (Bell, 1997; Olmeda, Fernández, 1997; Ahn, et al., 2000; Boyacioglu et al., 2009; Serrano, Gutiérrez, 2013; Chiamonte et al., 2015; Le, Viviani, 2018; Uthayakumar, et al., 2018; Alaka, 2018). Recently, researchers have used the SMOTE method to transform data in a balanced form for developing the most suitable prediction model and devised a technique to quantify the financial stress of firms under some constraints (Shrivastav, Ramudu, 2020; Shrivastava et al., 2020).

### *2.2 DL Approaches*

Although DL emerged almost two decades ago, it is now widely used in both academic research and industrial applications because of its ability to manage highly nonlinear data. DL has extensive applications in image recognition (Pak, Kim, 2017; Traore et al., 2018) voice recognition (Satt, et al., 2017; Khalil et al., 2019; Zhao, et al., 2019), and natural language processing (Deng, Liu, 2018; Kamath, et al., 2019). Some researchers have even used DL to solve issues encountered in the fields of finance and management science.

ANN and recurrent neural networks (RNN) are the two most common DL methods used for predicting stock price fluctuations (Fischer, Krauss, 2018). Convolutional neural network (CNN) is another crucial DL technique; however, this method has not been used for developing a bankruptcy prediction model (Qu, et al., 2019). Hosaka (2019) used CNN for the first time to analyze the bankruptcy of firms, the financial statement, and financial ratios of Japanese listed companies and convert the result into grayscale images. A theoretical framework for bankruptcy prediction was suggested by Hosaka (2019), and this framework

has dominated other predictive models including those developed using advanced ML techniques (Qu, et al., 2019).

DL models for bankruptcy prediction were introduced by Mai et al. (2019), particularly the neural network in which the model has more than one hidden layer. Mai et al. (2019) selected crucial features from the textual data of more than 9000 US public companies for bankruptcy prediction. The textual data collected from public news and the annual reports of these companies combined with the classical financial information of companies, such as financial ratios, yielded more suitable and efficient predictive models compared with those developed using standalone data. These findings and insights provide newer outlooks and motivations for research in this area. Some studies on bankruptcy prediction have even been conducted in the Indian context (Dhakar, et al., 2020; Smiti, Soui, 2020; Alexandropoulos, et al., 2019).

Most studies on bankruptcy prediction have primarily focused on countries that have a large number of bankrupt firms, especially banks. However, in a country such as India, surviving banks have far outnumbered failed banks, resulting in the issue of imbalanced classes. No study on this topic has yet used DL methods (Altman, 1968; Sinkey, 1975; Martin, 1977; Ohlson, 1980; Altman et al., 1994; Ahn et al., 2000; Wang et al., 2014; Chiaramonte et al., 2015; Le, Viviani, 2018; Uthayakumar et al., 2018; Shrivastav, 2019 and many more). In this study, we use an analytics-based methodological approach wherein we initially extract the most significant bankruptcy-related features, transform data from an imbalanced to a balanced form, choose suitable DL and ML techniques, and use them to develop the best predictive model.

### **3. Data Description, Descriptive Statistics, and Methodology**

We collect data for both failed and survived public<sup>5</sup> and private sector banks<sup>6</sup> in India from January 2000 to December 2018. We consider a bank to be a ‘failure’ when it meets one of the following conditions: merger or acquisition, bankruptcy, dissolution, and negative assets (Shrivastava, Ramudu, 2020). These conditions for failed or survived banks were verified by Altman (Altman, 1968; Altman, et al., 2017).

Data are collected for a total of 59 banks, of which 17 and 42 are failed and surviving banks, respectively. The target feature in the dataset has two classes, namely survival, and failure, and the proportion of classes is 0.97. The dataset has 618 instances with 26 financial and nonfinancial features depicted in table 1 below. Because collected data contains a mix of crucial and redundant features, we use a well-known feature selection technique called ‘Relief’ to formulate the bankruptcy model. ‘Relief’ is a nonparametric technique widely used for feature selection because of its simplicity and prevents the overfitting of the prediction model (Subsection 3.1.).

---

<sup>5</sup> Public Sector Banks (PSBs) are a major type of bank in India where a majority stake (i.e., more than 50%) is held by a government.

<sup>6</sup> India has banks where the majority of shares or equity are not held by the government but private shareholders.

**Table 1. Features and their descriptions**

Features	Description of Features
Status	Binary representation: 1 for failed banks and 0 for surviving banks
Total Assets	Current assets + advances + investment + fixed assets + others
Equity	Total capital – reserves and surplus
Total Liabilities	Net loans – reserves for impaired loans.
Deposits	Demand + saving + term deposits
Profit after tax	Operating profits + other incomes
Total Capital	Equity + reserves and surplus
Reserves and Funds	The reserve fund is a savings account or other highly liquid asset set apart by banks to meet any future costs
Return on assets	Net profit/total assets
Net Income	Posttax profit
Net Interest Revenue	Gross interest and dividend income minus total interest expense
Other Operating Income	Any other sustainable income that is related to a company's core business
Overheads	Personnel expenses and other operating expenses
Z-score	$(\text{Return on assets (ROA)} + \text{equity/asset})/\sigma$ (return on assets)
Loan Loss Reserves/Loans	Signifies how much funds have been put apart for potential losses.
Equity/Assets	Evaluates the amount of security the bank enjoys by its equity
Equity/Net Loans	Measures the equity insulation available to take up losses on the loan manuscript
Equity/Deposits	Estimates the amount of everlasting funding relative to undersized funding.
Equity/Liabilities	Identified as the capitalisation ratio and is the inverse of the leverage ratio.
Net Interest Margin	Net interest income expressed as a percentage of earning assets
Cost/Income	Estimates the costs of managing the bank, the main element of salaries, as a proportion of income produced before provisions.
Net Loans/Assets	Proportion of resources coupled up in loans
Growth of Real GDP	Gross domestic product at market price
Inflation	Logarithmic change of the GDP deflator year wise
C3/All	Percentage of total assets held by the big three banks of total assets of the banking industry
C5/All	Percentage of total assets held by the big five banks of total assets of the banking industry

In this study, extreme values (outliers) are winsorized upon 1% and 99% for surviving banks, whereas failed banks may represent some financial stress in the case of extreme values. The target feature in this study is the bank's status, namely survived or failed. The status of banks is used as a categorical variable where 0 represents surviving banks and 1 represents failed banks. Furthermore, significant features used for bankruptcy prediction are the statement, balance sheet, financial ratios, and country-specific variables. The dataset used in this study includes approximately 92% of Indian banks.

The descriptive statistics of all features included in the dataset are listed in Table 2. All the features in the dataset except those presented as a percentage or ratio are in millions. As shown in Table 3, the standard deviation values of most of the features are high, indicating a large variation in the bank's financial profile. The basic statistical measure of financial profiles for survived and failed banks are listed in Table 3. The t-test values for mean differences in the different features of banks with different profiles are provided in Table 3.

**Table 2. Descriptive statistics for private and public sector Indian banks over the period 2000-2017**

Bank-specific variables		Mean	Max	Min	Std. Dev.	N
Status	Survived (0) or failed (1)	0.03	1	0	0.16	838
Size	Total Assets	0.62	1	0	0.49	825
Bank type	Public sector banks as 1 and private sector banks as 0	0.64	1	0	0.48	838
Profit after tax	Operating profits $\pm$ other incomes	8268	145,496	-60,892	19,077	823
Total assets	Current assets+ advances + investment + fixed assets + others	1,185,955	27,059,663	0.5	2,239,587	823
Total capital	Equity + reserves and surplus	4371	45,739	0.5	5646.94	822
Deposits	Demand + saving + term deposits	952,140	20,447,514	866	1,725,551	814
Loans and advances	Loans and advances	705,731	15,710,784	763	1,381,711	821
Return on assets	Net profit/total assets	0.85	4.46	-6.5	0.81	794

**Table 3. Descriptive analytics of the dataset**

Feature	Minimum	Median	Mean	Maximum	Standard Deviation
Total Assets	0.0	472732.0	1184797.0	27059700.0	2238500.0
Equity	-9900.0	26900.0	78300.0	1883000.0	160050.0
Total Liabilities	700.0	401700.0	952200.0	20448000.0	1715000.0
Deposits	866.0	401609.0	952140.0	20447514.0	1725551.0
Profit after tax	-60892.1	3349.9	8269.0	145496.4	19065.0
Total Capital	0.0	472726.0	1184516.0	27059663.0	2238607.0
Reserves and Funds	-34971.0	23374.0	73863.0	1874887.0	158809.0
Return on assets	-6.5	0.9	0.9	4.5	0.8
Net Income	0.0	54039.0	132996.0	2700874.0	245053.0
Net Interest Revenue	-14064.0	11839.0	29703.0	625481.0	58665.0
Other Operating Income	79.5	40600.3	103652.6	2075392.8	187286.0
Growth Overheads	34.3	23469.5	61682.3	1139568.9	105937.0
Z-score	-3.3	2.0	2.3	11.5	2.1
Loan Loss Reserves/Loans	0.0	0.0	0.0	0.5	0.04
Equity/Assets	-50.6	0.1	0.0	1.0	1.8
Equity/Net Loans	-0.1	0.1	0.2	11.4	0.5
Equity/Deposits	-0.1	0.1	0.1	11.7	0.43
Equity/Liabilities	-1.0	0.1	0.1	19.9	0.8
Net Interest Margin	0.0	0.0	0.0	0.7	0.0
Cost/Income	0.9	1.6	1.6	22.8	0.8
Net Loans/Assets	0.0	0.6	0.5	0.7	0.1
GDP growth	0.0	0.1	0.1	0.2	0.0
Inflation CPI	2.2	6.3	6.9	15.0	3.2
C3/All	0.0	0.3	0.3	0.3	0.1
C5/All	0.0	0.4	0.3	0.4	0.1

**Table 4. The t-test values for mean differences in different features**

Features	Survived Banks	Failed Banks
Profit After Tax	10020	2100***
Total Assets	1435000	300500**
Return on Net worth	0.89	0.57***
Equity	95800	15900***
Total Liabilities	1339750	284550***
Total Provision	35392	6790***
Loans	935100	246900***
Net Interest Revenue	35900	8000 ***
Other operating income	125050	27887***
Growth overheads	74300	17260***
Loan Loss Reserves/Loans	0.04	0.03
Equity/Assets	-0.02	0.07
Equity/ Net loans	0.14	0.22
Equity/Deposits	0.08	0.15
Equity/Liabilities	0.1	0.12
Net Loans/Assets	0.51	0.5
Net Interest Margin	0.04	0.04
Cost/Income	1.58	1.69
Z-score	2.25	2.13**
Inflation CPI	7	6.6
C3 All	0.27	0.23***
C5 All	0.26	0.20***
GDP growth	0.11	0.13

\*\*\*, \*\*, and \* represents statistical significance at 1%, 5%, and 10%, respectively.

The financial profiles of surviving and failed banks are presented in column I and column II in Table 3. As shown in columns I and II, surviving banks are financially healthier than failed banks. The net income and equity are 10020 and 95800, respectively, for surviving banks and 2100 and 15900, respectively, for failed banks. The equity/assets for surviving banks are 0.07, whereas those for failed banks are -0.01. Overall, the financial health of surviving banks is superior to that of failed banks.

### 3.1 Two-Step Feature Selection

Kira and Rendall (1992) formulated an instance-based feature selection method called 'Relief'. This technique preserves the balance between the computational complexity and accuracy of ML and DL methods. The Relief technique allocates weights to independent features that indicate the significance of this feature with regard to target features. The maximum and minimum values of weights allocated to independent variables by the Relief algorithm are +1 and -1, where +1 shows the most crucial variables and -1 represents the most redundant variables.

The Relief algorithm is a nonparametric technique for feature selection and indicates the significance of features based on the contribution of other features. The Relief algorithm does not possess any assumptions regarding the distribution of independent variables or the size of the sample. The features that have positive weights are considered to be significant,

whereas those with negative weights are considered to be redundant and thus discarded from the model.

Computer pseudo-codes for the “Relief” technique are indicated below:

**Initial Requirement:** First, we use features of each record where ‘0’ represents the class for surviving banks and ‘1’ represents the class for failed banks. ‘R’ coding is used to implement the Relief algorithm in this study, where ‘I’ represents the number of records in the training dataset, ‘V’ denotes the number of features in each record of the training dataset, and ‘T’ represents randomly selected training records from the ‘I’ records of the training dataset. ‘A’ represents the randomly selected feature for randomly selected training records.

The dummy code for the ‘Relief’ technique is indicated as follows:

Assume that the weight of each feature is zero,  $W [A] = 0$ .

For  $i = 1$  to  $T$  do

Select a random target instance, e.g.  $L_i$

For this randomly selected instance, check the closest hit ‘H’ and closest miss ‘I’.

For  $A = 1$  to  $V$  do

Weight  $[A] = \text{Weight } [A] - \text{diff } [A, L_i, H]/T + \text{diff } [A, L_i, I]/T$

Finish (second loop)


Finish (First loop)

Return ( $W [A]$ ).

The algorithm chooses records from training data (e.g. ‘ $L_i$ ’) without replacement. For chosen instances from the training dataset, the weights of all variables are updated based on differences observed between target and neighbour instances. This process continues, and in each round, the distance of the ‘target’ instances from all other instances is calculated. Furthermore, this method selects the two closest neighbour instances from the same class (0 or 1), termed as the closest hit (‘H’), and the closest neighbour with a different class, termed as the closest miss (‘M’). The weights in each round are updated based on the closest hit or miss.

If it is the closest hit or features differ for the same class (0 or 1), then the weight decreases by  $1/N$ , and if it is a closed miss ‘M’, then the weight increases by  $1/N$ . This process continues until all the features of all instances are finished by the loop. An example of the ‘Relief’ technique is as follows:


Target Instance ( $L_i$ )	PQPQPQPQPQPQPQ	0
Closest Hit (H)	PQPQPQPQPQPQPQ	0





In this example, because of the mismatch of a feature where the instances are from the same class (red colour), a weight  $-1/T$  is assigned to the feature.

Target Instance (Li)	PQPQPQPQPQPQPQ	0
Closest Hit (H)	PQPQPQPPQPQPQ	1



In this example, because of the mismatch of features where instances are from different classes (red colour), a weight  $1/N$  is allocated to the feature. This process follows the last instances, and this method for estimating weights is valid for discrete features only.

The diff. function computes the difference in the value of feature 'A' with two instances  $I_1$  and  $I_2$ , where  $I_1 = L_i$  and  $I_2$  is either 'H' or 'M' during weight updates. The diff. function of a discrete feature is as follows:

$$\text{diff.}(A, I_1, I_2) = \begin{cases} 0 & \text{if } \text{value}(A, I_1) = \text{Value}(A, I_2) \\ 1 & \text{otherwise} \end{cases}$$

The diff. function of a continuous feature is as follows:

$$\text{diff.}(A, I_1, I_2) = \frac{|\text{value}(A, I_1) - \text{value}(A, I_2)|}{\max(A) - \min(A)}$$

The weights for feature 'A' are calculated for all instances. The weights are normalized so that their value is between 0 and 1. The weights calculated using the 'Relief' algorithm are then fed into various ML and DL methods for developing bankruptcy prediction models.

### 3.2 Imbalanced Class and SMOTE

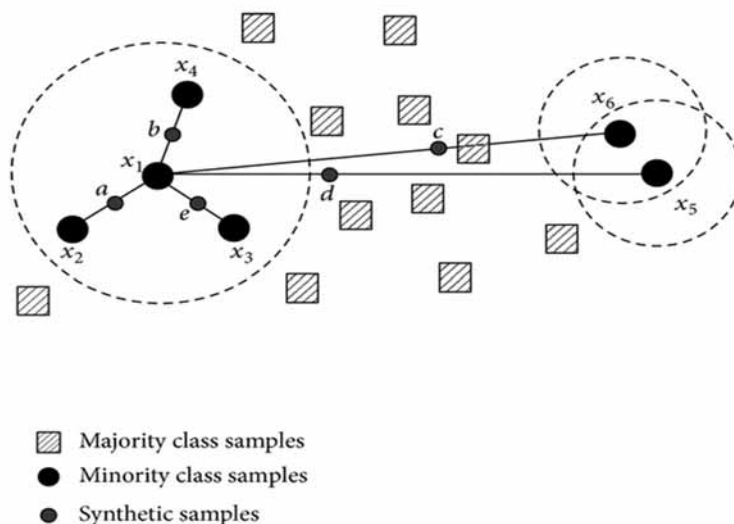
Imbalanced data are a crucial issue in ML and DL in which one class contains more instances than the other class. Undersampling, oversampling, and SMOTE are prominent techniques used to solve the problem of imbalanced data (Chawla et al., 2002). Both undersampling and oversampling replicate minority classes for balancing them in data, whereas SMOTE overcomes the imbalances of classes by creating dummy instances. SMOTE is a powerful and widely used method that generates dummy instances from minority instances. SMOTE generates a new minority class instance by interpolation of the nearest minority class instance randomly as a pictorial representation shown in figure 1.

First, for each minority class instance  $\mathcal{X}$ , one gets its k-nearest neighbours from other minority class instance. Second, select one minority class instance  $\bar{\mathcal{X}}$  among neighbours. Finally, create a synthetic instance  $\mathcal{X}_{new}$  by interpolating from  $\mathcal{X}$  and  $\bar{\mathcal{X}}$  as follows:

$$x_{\text{new}} = x + \text{rand}(0,1) \times (\bar{x} - x) \quad (1)$$

Here,  $\text{rand}(0, 1)$  creates a random number lying between 0 and 1.

**Figure 1. Synthetic data generation method using SMOTE**



SMOTE interpolates a new minority class instance from two minority class instances.

### 3.3 Logistic regression

Logistic regression (Kumar, U.D., 2017) is a supervised ML algorithm used to predict classes from an input feature. It provides the probability of a class by using the logit function. The logistics regression model is given below:

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (2)$$

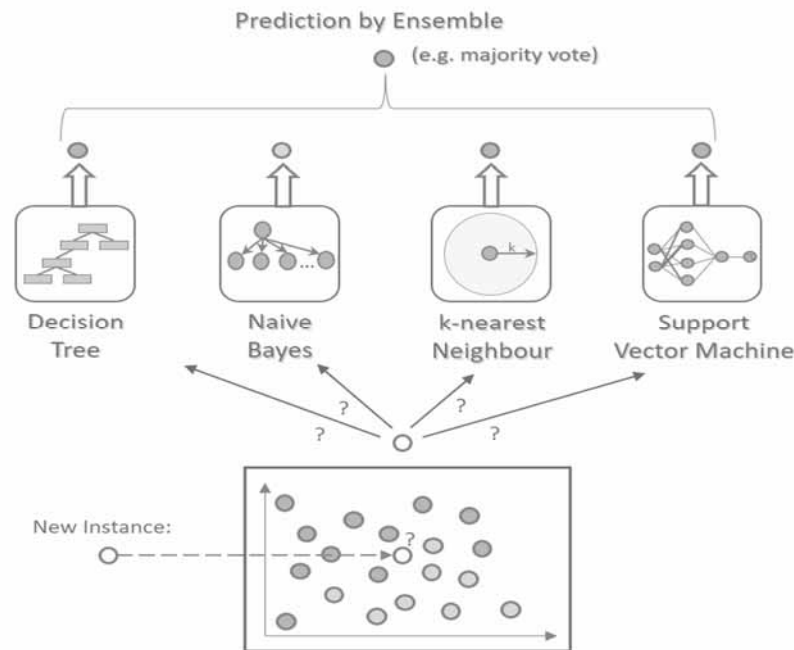
If  $p$  denotes the likelihood of success, then  $1 - p$  would be the likelihood of failure, especially in the case of binary class instances (failure and nonfailure). The  $x_0, x_1, \dots, x_n$  are the features of a logistics model and  $\beta_0, \beta_1, \dots, \beta_n$  are the coefficient estimates of features.

### 3.4 Ensemble Learning

Ensemble learning is a powerful method to increase the performance of a predictive model. Ensemble learning is a group learning method that provides higher accuracy and model

stability. This technique uses various ML algorithms to predict an accurate class. Classification is performed through majority voting, whereas regression is performed using the averaging method as depicted in figure-2.

**Figure 2. Ensemble Learning method**

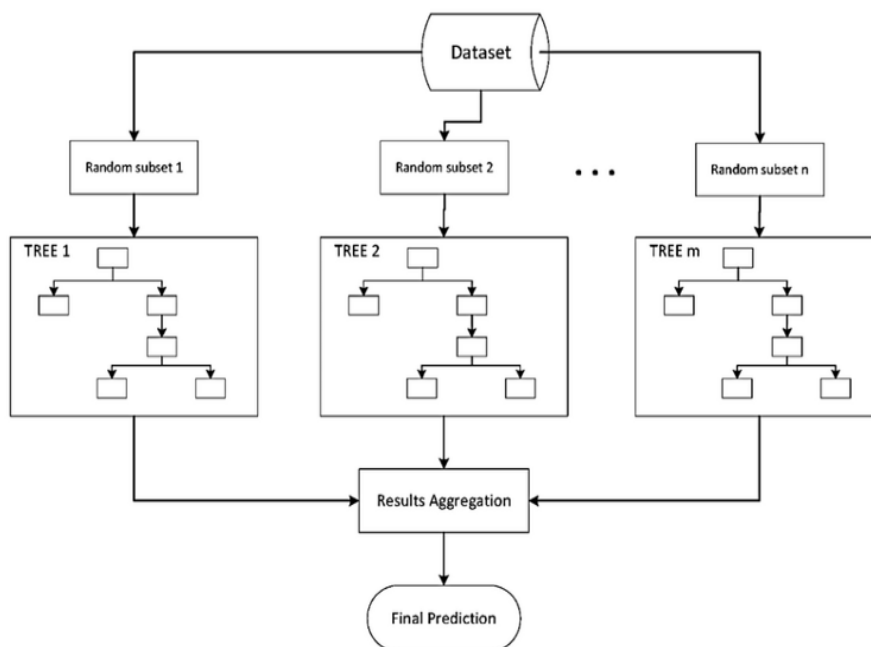


The two types of ensemble methods are bagging and boosting. Variance and overfitting can be reduced using ensemble techniques, thus reducing the bias of the predictive model.

#### 3.4.1. Bagging Technique: Random forest

Bagging is an ensemble technique used to improve the accuracy and stability of a model. In the bagging approach, the same learning algorithm is trained with the subsets of a dataset that are randomly picked from the training dataset. We select the subsets of the training dataset into bags randomly and then train the learning model of each bag (Figini et al., 2016). The final prediction is performed by combining the results of all model results. In this study, we use the random forest, a widely used bagging technique based on decision tree models. Random forest is particularly robust and allows for the presence of outliers and noise in the training set (Yeh et al., 2014) and the pictorial representation of the ensemble technique is given in Figure 3.

**Figure 3. Bagging procedure for algorithm learning**



The steps of a random forest algorithm (Yeh et al., 2014) are as follows:

1. Create random subsets of the parent dataset that are composed of an arbitrary number of observations and different features.
2. Each subset from step 1 produces a decision tree.
3. For each observation, the forest uses a large number of votes. The class with the most votes is chosen as the preferred classification of the element.

The random forest identifies the importance of each variable in classification results; therefore, it provides not only the classification of observations but also information regarding the significance of features for the separation of classes (Maione et al., 2016).

#### 3.4.2. Boosting Technique: Adaboost

Boosting is another ensemble method that combines weak learners to create a strong learner to make more accurate predictions. Boosting begins with a weak classifier that is prepared using training data. A classifier learning algorithm is considered to be weak when small changes in data induce large changes in the predictive model. In the next iteration, the new classifier focuses on or places more weight on those cases that were incorrectly classified in the last round.

AdaBoost is a successful and efficient method for classification (Kim, Upneja, 2014). Initially, Adaboost assigns weights to all  $k$  observations  $1/k$ . Thus, the first sample is uniformly generated from initial observations. After the training set  $X_i$  is extracted from  $X$ , a classifier  $Y_i$  is trained on  $X_i$ . The error rate is calculated considering the number of observations of the training set. The new weight for each observation is based on the effectiveness of the classifier  $Y_i$ . If the error rate is higher than a random guess, then the test set is discarded, and another set is generated using original weights (initially  $1/k$ ). If the error rate is satisfactory, the weights of the observation are updated according to the importance of the classifier. These new weights are then used to generate another sample from initial observations. The boosting technique involves the following steps (Heo, Yang, 2014):

1. The distribution of weights  $w_1(i) = 1/k$  is created, where  $i = 1, 2, \dots, k$ , and  $w_t$  is the iterative weighting ( $t = 1, \dots, T$ ),

$$w_{t+1}(i) = \frac{w_t(i) e^{\alpha_t(2I(y_i \neq h_t) - 1)}}{\sum_{i=1}^k w_t(i) e^{\alpha_t(2I(y_i \neq h_t))}}, \text{ where } h_t = \arg \max |0.5 - \epsilon_t| \text{ is the error such that}$$

$\epsilon_t = \sum_{i=1}^k w^t(i) |I(y_i \neq h_t(x_i)) - 0.5|$ ,  $I = 1$  when the measure is accurately computed; otherwise, it is 0.

2. In each cycle,  $\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$  is recalculated. The process completes when  $|0.5 - \epsilon_t| \neq \delta$ , where  $\delta$  is constant.

3.  $Y(x)$  is evaluated for the complete boost by  $Y(x) = \text{sign} \sum_{t=1}^T \alpha_t h_t(x)$ .

### 3.5 DL approach (ANN)

ANN is a deep learning technique that can be used to determine the pattern of nonlinear data. ANN is based on input variables that communicate to one or more hidden layers with a combination of neurons and predict the output class (survival of failed banks in this case). The idea behind the ANN method is to simulate the human brain where neurons communicate with others with the help of signals (layers) (Shanmuganathan, 2016). The output of the ANN is based on input, weights, and bias term  $b_i$  as follows:

$$h_i = f^{(i)} \left( b_i^{(i)} + \sum_{j=1}^n w_{ij} x_j \right)$$

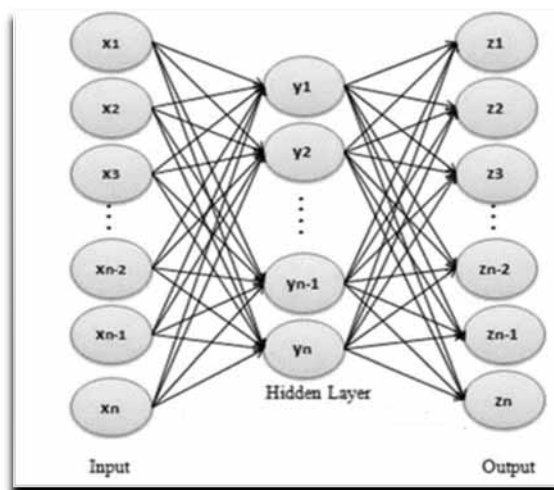
Here,  $x_i$  represents input data and  $w_{ij}$  is the weight of layers from the  $j^{\text{th}}$  input neuron to the  $i^{\text{th}}$  hidden neuron. In the ANN model, the first layer is the input layer that is equal to the number of input features in the model represented by  $x_i, i=1, \dots, n$ . Each input has different

weights  $w$  based on its relationship with the output class (survival and failure). The output of the neuron is calculated using the following formula:

$$z_i = f^{(2)} \left( b_i^{(2)} + \sum_{j=1}^{n_k} v_{ij} h_j \right)$$

where  $n_k$  is the number of hidden neurons, and  $v_{ij}$  denotes the weight connecting the hidden neuron  $i$  to the output neuron  $j$ . At the start, networks are initialized using random weights. Subsequently, the values of weights are iteratively adjusted to reduce the loss function. ANN has been criticized for its black-box nature and lengthy training process in the development of an optimal model. The raw structure of ANN is shown in Figure 4.

**Figure 4. Basic structure of the ANN model**



#### 4. Empirical Results

The data used in this study has 618 instances in which the class proportions (survival and failure) are 2.4% and 97.6%, respectively. The parent dataset is separated into two portions, called train and test, accounting for 80% and 20% of data, respectively. First, we formulate the logistics prediction model for bankruptcy by using the training dataset and validate it by using the test dataset. The predictive model provides a precision of 0/0 with 0.5 as the threshold value. The prediction model has a high number of false negatives with a low recall value. The F value of the prediction model is undefined at 0/0, and the overall precision and accuracy of the model are low. The low accuracy of the model can be attributed to the model being biased for the majority class, thus being unable to understand the pattern for the minority class. Therefore, transforming the dataset into a balanced form is essential before formulating the prediction model. The SMOTE technique is used to transform the dataset into a balanced form. In this case, after using the SMOTE technique, the ratios of minority

and majority classes become nearly equal with 1180 instances and 26 features (Table 1). Next, the Relief technique is applied to the dataset for significant variables related to the failure of the bank. The feature whose weights are more than 0 (Table 5), as indicated by the Relief algorithm, are considered to be significant features and fed into a different ML and DL model to formulate the bankruptcy prediction model. The weights calculated using the Relief algorithm are given in Table 4.

**Table 5. List of features and their weight estimated using the Relief algorithm**

Features Name	Relief Score
Total Assets	00
Equity	-0.09
Total Liabilities	0.08
Deposits	-0.012
Profit after tax	0.14
Total Capital	0.02
Reserves and Funds	0.09
Return on assets	0.15
Net Income	00
Net Interest Revenue	-0.12
Other Operating Income	0.09
Overheads	-0.1
Z-score	0.23
Loan Loss Reserves/Loans	-0.1
Equity/Assets	0.10
Equity/Net Loans	0.12
Equity/Deposits	-0.1
Equity/Liabilities	0.12
Net Interest Margin	00
Cost/Income	0.1
Net Loans/Assets	0.02
Growth of Real GDP	00
Inflation	00
C3/All	0.08
C5/All	0.09

The balanced data are divided into training and testing datasets, accounting for 80% and 20% of data, respectively. Because of the use of the SMOTE algorithm, the risk of bias and the overfitting of the model is high. To prevent these issues, random forest and Aaboost algorithms are used. Notably, the models can make wrong predictions when it is validated using the test dataset. The bankruptcy prediction model can predict a bank to be failing or surviving when the true status of banks may differ. A total of 4 conditions may occur as given in table 5 below:

1. Failed banks are falsely classified as surviving, whereas failed banks are correctly classified as failing.
2. Surviving banks are rightly classified as surviving, whereas surviving banks are falsely classified as failing.

**Table 6. Confusion Matrix: True versus Forecasted Results**

Forecasted Results		True Results	
		Positive	Negative
	Positive	True Positive	False Positive
	Negative	False Negative	True Negative

$$\text{Sensitivity}(1 - \text{Type II error}) = \frac{\text{true negatives}}{\text{true negatives} + \text{false positive}}$$

&

$$\text{Specificity}(1 - \text{Type I error}) = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

As discussed above, in two cases, errors occur. The first case is when failed banks are classified as surviving, and the second case is when surviving banks are classified as failing. These two types of wrong classification of banks are related to Type I and Type II errors.

Forecasting accuracy and Type I or II errors are calculated using the test dataset to compare various predictive models in this study (Lin et al., 2012). However, the total accuracy of the model is not an appropriate measure to compare various predictive models because the Type II error is more sensitive compared with Type I error in this study. The Type I error indicates the number of surviving banks that have been incorrectly classified as failed banks. By contrast, the Type II error indicates the number of failed banks that the model incorrectly classified as surviving banks. The Type II error is more acute for banks because if the predictive model makes wrong decisions that are highly likely for bankruptcy, it creates a challenging problem for banks as time passes. Overall, a prediction model with low Type II error and high accuracy is considered the best prediction model in this study. Therefore, the prediction model that can provide the highest accuracy and lowest Type II error rate can be regarded as the best predictive model. We use ML and DL techniques such as logistic regression, random forest, AdaBoost, and ANN to formulate bankruptcy predictive models and validate them using training and testing datasets, respectively.

The Type II error is 64.34%, and the accuracy is 68.65% in the bankruptcy prediction model developed using logistic regression as given in Table 7.

**Table 7. True vs Forecasted**

Forecasted		True	
		1	0
	1	41	5
	0	74	132

The type II error is 71.8% and the accuracy is 58.26% for the bankruptcy prediction model developed using the random forest as given in Table 8.



Table 8. True vs Forecasted

		True	
		1	0
Forecasted	1	48	4
	0	67	133

The Type II error is 1.73% and the accuracy is 98.8% for the bankruptcy prediction model developed using AdaBoost as given in table-9 below:

Table 9. True vs Forecasted

		True	
		1	0
Forecasted	1	113	1
	0	2	136

The Type II error is 0.86% and the accuracy is 99% for the bankruptcy prediction model developed using ANN as given in table-10 below:

Table 10. True vs Forecasted – Artificial Neural Network

		True	
		1	0
Forecasted	1	115	1
	0	1	135

Table 11. Error of the model on the test dataset

Techniques	Type-II error
Logistics Regression	64.34%
Random Forest	58.25%
AdaBoost	1.74%
Artificial Neural Network	.87%

The high accuracy and low Type-II error rate are statistical measures used to compare ML or DL models. The accuracy and Type-II error rates of all bankruptcy prediction models in Table 11 indicate that ANN is the most favourable bankruptcy prediction model, although none of the formulated bankruptcy models in the study has a 0 Type-II error. One of the likely reasons may be that some banks are financially healthy but acquisitions or mergers occurred due to government policies or to reduce the operational cost. The second reason can be that feature selection techniques have eliminated some of the financial features of firms from modelling even if they may be important for bankruptcy. For example, SBI Commercial and Intl. Bank has been forecasted as a surviving bank, although we consider it as a failed bank in the original dataset. This bank was merged with the SBI by the government to minimize operational costs and not due to the financial crisis. Typically, these are some scenarios that result in Type-II errors in models. Therefore, based on the trade-off among complexity, accuracy, and Type II error of the bankruptcy prediction model, ANN has the highest accuracy and is the most favourable model.

## 5. Conclusions and Implications of the Study

In this study, a systematic framework is developed for analyzing a bank's financial stress and to formulate an efficient and generalized bankruptcy model. In this study, data are collected from the Prowess Database, a publically available dataset for Indian banks that contains data from 2000 to 2018 with several missing values. We develop bankruptcy prediction models by using logistics, random forest, AdaBoost, and ANN and perform a comparison based on their accuracy. Finally, based on Type-I error and the accuracy of the model, ANN is found to be the most favourable prediction model. The possible reason is that ANN can identify a highly nonlinear pattern in the dataset compared with other techniques. The proposed method provides a holistic approach, starting from selecting a list of significant features for bankruptcy prediction by using the 'Relief' algorithm, transforming the dataset into a balanced form through SMOTE, and selecting appropriate ML techniques that can predict bank failures. This model can be useful for decision-makers who can obtain a future warning regarding firms before they undergo insolvency.

## References

- Aha, D. W., Kibler, D., Albert, M. K. (1991). Instance-based learning algorithms. – *Machine learning*, 6(1), pp. 37-66.
- Ahn, B. S., Cho, S. S., Kim, C. Y. (2000). The integrated methodology of rough set theory and artificial neural network for business failure prediction. – *Expert systems with applications*, 18(2), pp. 65-74.
- Alaka, H. A., Oyedele, L. O., Owolabi, H. A., Kumar, V., Ajayi, S. O., Akinade, O. O., Bilal, M. (2018). Systematic review of bankruptcy prediction models: Towards a framework for tool selection. – *Expert Systems with Applications*, 94, pp. 164-184.
- Alexandropoulos, S. A. N., Aridas, C. K., Kotsiantis, S. B., Vrahatis, M. N. (2019 May). A deep dense neural network for bankruptcy prediction. – In: *International Conference on Engineering Applications of Neural Networks*. Springer, Cham, pp. 435-444.
- Altman, E. I. (1968). Financial ratios, discriminant analysis and the prediction of corporate bankruptcy. – *The journal of finance*, 23(4), pp. 589-609.
- Altman, E. I., Iwanicz-Drozdowska, M., Laitinen, E. K., Suvas, A. (2020). A race for long horizon bankruptcy prediction. – *Applied Economics*, pp. 1-20.
- Altman, E. I., Iwanicz-Drozdowska, M., Laitinen, E. K., Suvas, A. (2017). Financial distress prediction in an international context: A review and empirical analysis of Altman's Z-score model. – *Journal of International Financial Management & Accounting*, 28(2), pp. 131-171.
- Altman, E. I., Marco, G., Varetto, F. (1994). Corporate distress diagnosis: Comparisons using linear discriminant analysis and neural networks (the Italian experience). – *Journal of banking & finance*, 18(3), pp. 505-529.
- Appiah, K. O., Chizema, A., Arthur, J. (2015). Predicting corporate failure: a systematic literature review of methodological issues. – *International Journal of Law and Management*.
- Barboza, F., Kimura, H., Altman, E. (2017). Machine learning models and bankruptcy prediction. – *Expert Systems with Applications*, 83, pp. 405-417.
- Bell, C. M. (1997). *Ritual: Perspectives and dimensions*. Oxford University Press on Demand.
- Boyacioglu, M. A., Kara, Y., Baykan, Ö. K. (2009). Predicting bank financial failures using neural networks, support vector machines and multivariate statistical methods: A comparative analysis in the sample of savings deposit insurance fund (SDIF) transferred banks in Turkey. – *Expert Systems with Applications*, 36(2), pp. 3355-3366.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., Kegelmeyer, W. P. (2002). SMOTE: synthetic minority over-sampling technique. – *Journal of artificial intelligence research*, 16, pp. 321-357
- Chiaromonte, L., Poli, F., Oriani, M. E. (2015). Are cooperative banks a lever for promoting bank stability? Evidence from the recent financial crisis in OECD countries. – *European Financial Management*, 21(3), pp. 491-523.

- Dellepiane, U., Di Marcantonio, M., Laghi, E., Renzi, S. (2015). Bankruptcy prediction using support vector machines and feature selection during the recent financial crisis. – *International Journal of Economics and Finance*, 7(8), pp. 182-195.
- Deng, L., Liu, Y. (eds.). (2018). *Deep learning in natural language processing*. Springer.
- Devi, S. S., Radhika, Y. (2018). A survey on machine learning and statistical techniques in bankruptcy prediction. – *International Journal of Machine Learning and Computing*, 8(2), pp. 133-139.
- Dhakar, P., Srivastava, S., Tiwari, K. (2020). DLBR: Bankruptcy Prediction using Deep Learning-A Case Study on Indian Firms. – *Asian Journal of Research in Banking and Finance*, 10(3), pp. 1-10.
- Fernández, A., Garcia, S., Herrera, F., Chawla, N. V. (2018). SMOTE for learning from imbalanced data: progress and challenges, marking the 15-year anniversary. – *Journal of artificial intelligence research*, 61, pp. 863-905.
- Figini, S., Savona, R., Vezzoli, M. (2016). Corporate default prediction model averaging: a normative linear pooling approach. – *Intelligent Systems in Accounting, Finance and Management*, 23(1-2), pp. 6-20.
- Fischer, T., Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. – *European Journal of Operational Research*, 270(2), pp. 654-669.
- Hosaka, T. (2019). Bankruptcy prediction using imaged financial ratios and convolutional neural networks. – *Expert systems with applications*, 117, pp. 287-299.
- Kadioglu, E., Telceken, N., Ocal, N. (2017). Effect of the asset quality on the bank profitability. – *International Journal of Economics and Finance*, 9(7), pp. 60-68.
- Kamath, U., Liu, J., Whitaker, J. (2019). *Deep learning for nlp and speech recognition (Vol. 84)*. Springer.
- Khalil, R. A., Jones, E., Babar, M. I., Jan, T., Zafar, M. H., Alhussain, T. (2019). Speech emotion recognition using deep learning techniques: A review. – *IEEE Access*, 7, pp. 117327-117345.
- Kim, S. Y., Upneja, A. (2014). Predicting restaurant financial distress using decision tree and AdaBoosted decision tree models. – *Economic Modelling*, 36, pp. 354-362.
- Kira, K., Rendell, L. A. (1992). A practical approach to feature selection. – In: *Machine Learning Proceedings 1992*. Morgan Kaufmann, pp. 249-256.
- Kou, G., Chao, X., Peng, Y., Alsaadi, F. E., Herrera-Viedma, E. (2019). Machine learning methods for systemic risk analysis in financial sectors. – *Technological and Economic Development of Economy*, 25(5), pp. 716-742.
- Kumar, U. D. (2017). *Business analytics: The science of data-driven decision making*. Wiley India.
- Le, H. H., Viviani, J. L. (2018). Predicting bank failure: An improvement by implementing a machine-learning approach to classical financial ratios. – *Research in International Business and Finance*, 44, pp. 16-25.
- Liang, D., Lu, C. C., Tsai, C. F., Shih, G. A. (2016). Financial ratios and corporate governance indicators in bankruptcy prediction: A comprehensive study. – *European Journal of Operational Research*, 252(2), pp. 561-572.
- Lin, J. W., Chen, C. W., Peng, C. Y. (2012). Potential hazard analysis and risk assessment of debris flow by fuzzy modelling. – *Natural hazards*, 64(1), pp. 273-282.
- Mai, F., Tian, S., Lee, C., Ma, L. (2019). Deep learning models for bankruptcy prediction using textual disclosures. – *European journal of operational research*, 274(2), pp. 743-758.
- Maione, C., Batista, B. L., Campiglia, A. D., Barbosa Jr. F., Barbosa, R. M. (2016). Classification of geographic origin of rice by data mining and inductively coupled plasma mass spectrometry. – *Computers and Electronics in Agriculture*, 121, pp. 101-107.
- Makinen, M., Solanko, L. (2018). Determinants of Bank Closures: Do Levels or Changes of CAMEL Variables Matter?. – *Russian Journal of Money and Finance*, 77(2), pp. 3-21.
- Martin, D. (1977). Early warning of bank failure: A logit regression approach. – *Journal of banking & finance*, 1(3), pp. 249-276.
- Min, J. H., Lee, Y. C. (2005). Bankruptcy prediction using support vector machine with optimal choice of kernel function parameters. – *Expert systems with applications*, 28(4), pp. 603-614.
- Nwogugu, M. (2006). Decision-making, risk and corporate governance: New dynamic models/algorithms and optimization for bankruptcy decisions. – *Applied mathematics and computation*, 179(1), pp. 386-401.
- Nwogugu, M. (2008). Illegality Of Securitization, Bankruptcy Issues And Theories Of Securitization. – *Journal of International Banking Law & Regulation*, 23(7), pp. 363-375.
- Ohlson, J. A. (1980). Financial ratios and the probabilistic prediction of bankruptcy. – *Journal of accounting research*, pp. 109-131.
- Olmeda, I., Fernández, E. (1997). Hybrid classifiers for financial multicriteria decision making: The case of bankruptcy prediction. – *Computational Economics*, 10(4), pp. 317-335.
- Pak, M., Kim, S. (2017). August. A review of deep learning in image recognition. – In: *2017 4<sup>th</sup> international conference on computer applications and information processing technology (CAIPT)*, IEEE, pp. 1-3.

- Pappas, V., Ongena, S., Izzeldin, M., Fuertes, A. M. (2017). A survival analysis of Islamic and conventional banks. – *Journal of Financial Services Research*, 51(2), pp. 221-256.
- Qu, Y., Quan, P., Lei, M., Shi, Y. (2019). Review of bankruptcy prediction using machine learning and deep learning techniques. – *Procedia Computer Science*, 162, pp. 895-899.
- Satt, A., Rozenberg, S., Hoory, R. (2017). August. Efficient Emotion Recognition from Speech Using Deep Learning on Spectrograms. – In: *Interspeech*, pp. 1089-1093.
- Serrano-Cinca, C., Gutiérrez-Nieto, B. (2013). A decision support system for financial and social investment. – *Applied Economics*, 45(28), pp. 4060-4070.
- Shanmuganathan, S. (2016). Artificial neural network modelling: An introduction. – In: *Artificial neural network modelling*, Springer, Cham, pp. 1-14.
- Shrivastav, S. K. (2019). Measuring the Determinants for the Survival of Indian Banks Using Machine Learning Approach. – *FIIIB Business Review*, 8(1), pp. 32-38.
- Shrivastav, S. K., Ramudu, P. J. (2020). Bankruptcy Prediction and Stress Quantification Using Support Vector Machine: Evidence from Indian Banks. – *Risks*, 8(2), p. 52.
- Shrivastava, S., Jeyanthi, P. M., Singh, S. (2020). Failure prediction of Indian Banks using SMOTE, Lasso regression, bagging and boosting. – *Cogent Economics & Finance*, 8(1), p. 1729569.
- Sinkey Jr., J. F. (1975). A multivariate statistical analysis of the characteristics of problem banks. – *The Journal of Finance*, 30(1), pp. 21-36.
- Sujud, H., Hashem, B. (2017). Effect of bank innovations on profitability and return on assets (ROA) of commercial banks in Lebanon. – *International Journal of Economics and Finance*, 9(4), pp. 35-50.
- Traore, B. B., Kamsu-Foguem, B., Tangara, F. (2018). Deep convolution neural network for image recognition. – *Ecological Informatics*, 48, pp. 257-268.
- Uthayakumar, J., Metawa, N., Shankar, K., Lakshmanprabu, S. K. (2018). Intelligent hybrid model for financial crisis prediction using machine learning techniques. – *Information Systems and e-Business Management*, pp. 1-29.
- Yeh, C. C., Chi, D. J., Lin, Y. R. (2014). Going-concern prediction using hybrid random forests and rough set approach. – *Information Sciences*, 254, pp. 98-110.
- Yeh, Q. J. (1996). The application of data envelopment analysis in conjunction with financial ratios for bank performance evaluation. – *Journal of the Operational Research Society*, 47(8), pp. 980-988.