

# CLASSIFICATION OF A BANK DATA SET ON VARIOUS DATA MINING PLATFORMS

<sup>1</sup>NAKKA MANEESHA, SMT.LAKSHMI PRAVEENA<sup>2</sup>

<sup>1</sup>MCA STUDENT, DEPARTMENT OF MCA, VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, GUNTUR, ANDHRA PRADESH.

<sup>2</sup>ASSISTANT PROFESSOR, DEPARTMENT OF MCA, VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, GUNTUR, ANDHRA PRADESH.

## Abstract

The process of extracting meaningful rules from big and complex data is called data mining. Data mining has an increasing popularity in every field today. Data units are established in customer-oriented industries such as marketing, finance and telecommunication to work on the customer churn and acquisition, in particular. Among the data mining methods, classification algorithms are used in studies conducted for customer acquisition to predict the potential customers of the company in question in the related industry. In this study, bank marketing data set in UCI Machine Learning Data Set was used by creating models with the same classification algorithms in different data mining programs. Accuracy, precision and fmeasure criteria were used to test performances of the classification models. When creating the classification models, the test and training data sets were randomly divided by the holdout method to evaluate the performance of the data set.

## Keywords

mining, Classification, Decision, Data, Training, Industries, Communications technology.

## I. INTRODUCTION

Today, data mining is utilized in the arrangement of problems in numerous fields, for example, wellbeing, finance and instruction. Data mining considers are being carried out in the field of wellbeing for finding of the infection, in customer-oriented industries, for example, media transmission, insurance and banking to work on customer churn and customer acquisition. Right now, forecasting study was carried out to see whether the battle of a bank results in new customer

acquisition. Another purpose of this investigation is to see the results of a similar classification algorithms in different data mining programs. Gotten results were appeared in tables in the results segment. There are numerous classification algorithms ceaselessly produced for various applications in the literature on bank marketing data set. Bach et al. [2] have identified customers who responded emphatically to the battles by performing customer division with a variety of methods, for example, artificial neural networks. Sumathi and Sivanandam have utilized data mining methods of financial organizations and other foundations to discover interrelationships between data [3]. Keramati et al. have utilized decision trees, artificial neural networks, k-nearest neighbors and support vector machines, among the machine learning algorithms, to predict existing customers who might prefer contending banks utilizing a media transmission organization data situated in Iran. They have identified the algorithm that gives the best result by comparing the algorithms utilized in the examination [4].

Today, data mining is utilized in the arrangement of problems in numerous fields, for example, wellbeing, finance and training. Data mining contemplates are being carried out in the field of wellbeing for conclusion of the malady, in customer-oriented industries, for example, media transmission, insurance and banking to work on customer churn and customer acquisition. Right now, forecasting study was carried out to see whether the crusade of a bank results in new customer acquisition

Marketing is a procedure of recognizing the goal consumers to purchase or make an arrangement with a product by means of fitting frameworks. It presently promotes the process to purchase the

merchandise or service and even helps with planning the necessary for the product and persuade customers to purchase it. The overall purpose is to improve the selling of merchandise and enterprises for the industry, marketing, and commercial organizations. It additionally suits to preserve the status of the business [1].

## II. RELATED WORK

This section addresses the data mining and classification algorithms and data mining programs used throughout the study.

### A. Data Mining

Data mining is the process of extracting meaningful and structures information in the complex data sets. During this procedure, data mining methods such as classification, clustering and association rules are used. Data mining methods are used to analyze, categorize, summarize and determine the relationships using different dimensions of data [5]. These methods are divided into two groups as predictive or descriptive methods [7].

### Existing System

Data mining is used in the solution of problems in many fields such as health, finance and education. Data mining studies are being carried out in the field of health for diagnosis of the disease, in customer-oriented industries such as telecommunication, insurance and banking to work on customer churn and customer acquisition. In this research, a forecasting study was carried out to see whether the campaign of a bank results in new customer acquisition. Another purpose of this study is to see the results of the same classification algorithms in different data mining programs.

### Modules

#### ➤ The Applicant

The applicant is register and upload the value document to the bank employee for getting the business loan. The documents are identity document, project detail, asset document and collateral documents.

#### ➤ Document verification

The bank employee of clerk is first check the completion of requirement after he redirected to the staff. The staff is verifying the asset and collateral document and check the worth ness to retrieve the loan amount and credit score. After the application is moved to the manager. The manager analysis the event and detect the fraud by using HMM and he decided to provide loan.

#### ➤ Data Mining

Data mining is the process of extracting meaningful and structures information in the complex data sets. During this procedure, data mining methods such as classification, clustering and association rules are used. Data mining methods are used to analyze, categorize, summarize and determine the relationships using different dimensions of data

#### ➤ Classification Algorithms

Bank marketing data set in UCI Machine Learning Data Set was used. Models were created using classification algorithms on this data set. Classification algorithms used in the study are the k-nearest neighbor (k-nn), Naive Bayes (NB), and C4.5 decision tree.

### B. Classification Algorithms Used in the Study

In this study, bank marketing data set in UCI Machine Learning Data Set [1] was used. Models were created using classification algorithms on this data set. Classification algorithms used in the study are the k-nearest neighbor (k-nn), Naive Bayes (NB), and C4.5 decision tree. The classification algorithms used are addressed in this section.

#### 1) k-nearest neighbor algorithm (k-nn)

It is one of the most basic algorithms of sample-based learning algorithms. In this algorithm learning process is performed with the data in training set. The

new samples are classified according to the similarity within the samples in the training set [8]. The k-nearest neighbor algorithm finds k samples that are closest to the unknown data by looking at the pattern space to find which class the unknown data belongs to. Distance is calculated by distance calculation methods such as Euclidean and Manhattan, and distance between neighbors is found. Unknown data are assigned to the class value that most closely resembles the nearest neighbors [9].

### 2) Naive Bayes algorithm

The Naive Bayes algorithm is named after the English mathematician Thomas Bayes. Bayesian algorithms are among the statistical classification techniques and are based on the statistical Bayesian theorem. Bayes classifier is a predictive model, easier to apply. Naive Bayes is a classification algorithm that shows the relationship between the independent variables and the target variable [10].

Let  $X = \{x_1, x_2, x_3, \dots, x_n\}$  is the sample set, and  $c_1, c_2, c_3, \dots, c_n$  is the class set.

$$P(X|C_i) = \frac{P(X|C_i)P(C_i)}{P(X)}$$

Probabilities are computed as seen in Eq. (1) for the sample to be classified. The data sample with the highest probability, calculated for each class, belongs to that class [11].

### 3) C4.5 Decision Tree Algorithm

C4.5 algorithm has been developed by Ross Quinlan. The Gain Ratio is used in the C4.5 decision tree. C4.5 algorithm can work with either categorical or numerical attributes. Decision Trees generated by C4.5 can be used for classification; therefore, C4.5 is generally called a statistical classifier [12].

### Proposed System

There are many classification algorithms continuously developed for various applications in the literature on bank marketing data set. It has identified customers who responded positively to the campaigns by performing customer segmentation with a variety of methods such as artificial neural networks, k-nearest neighbors and support vector

machines, among the machine learning algorithms, to predict existing customers who would prefer competing banks using a telecommunication company data located in Iran.

### C. Data Mining Programs

Numerous programs have been developed to implement data mining applications. Commercial programs such as SAS and open source programs such as RapidMiner (YALE), Waikato Environment for Knowledge Analysis (WEKA), R, Konstanz Information Miner (KNIME) can be given as examples of data mining programs developed [13]. In this section, the data mining programs used throughout the study are described briefly.

#### 1) Knime

Konstanz Information Miner (Knime) is a data mining program developed by the Konstanz University data science team [14]. Knime can import data of various file extensions (such as .txt, .arff, .csv) [15].

#### 2) RapidMiner (Yale)

It is a program developed by Ralf Klinkenberg, Ingo Mierswa and Simon Fischer in Artificial Intelligence Unit of Dortmund University of Technology. The Yale program has been developed at Yale University [16]. Yale has been reintroduced in 2007 with the RapidMiner name [17]. It works with 22 different file formats. It supports many databases such as Oracle, MS SQL Server, MySQL, IBM DB2 and text files [14]. It can run on MS Windows, Linux, and Mac OS X operating systems.

#### 3) Weka

Waikato Environment for Knowledge Analysis (Weka) is an open source data mining program developed using Java in Waikato University under the GNU general public license [18]. It accesses the SQL database using Java Database Connectivity (JDBC) [19]. It includes all the data mining and machine learning algorithms. It works on the .arff (Attribute Relationship File Format) file format specially designed for WEKA.

#### 4) R

It is a computer program developed for statistical calculation as well as being a programming language on its own. It has thousands of modules. With these packages, numerous operations can be performed such as data mining and data visualization. The R

language, which has been developed by Ross Ihaka and Robert Gentleman in the New Zealand University of Auckland, is continually evolving due to the increased number of packages programmed in accordance with the needs [20]. It has been developed open source as a alternative to the S software [13].

III. APPLICATION

In this study, an application was carried out with classification algorithms of data mining methods in order to predict customer acquisition using the bank marketing data set in the UCI database.

A. Data Set

In this study, bank marketing data set in UCI Machine Learning Data Set [1] was used by establishing models with the same classification algorithms in different data mining programs. The bank marketing data set contains 17 attributes and 45211 customer records. Table 1 shows the data type and description of the attributes.

B. Model Performance Evaluation Criteria

Evaluation of the model created by classification algorithms is carried out by various methods. One of these methods is the confusion matrix [21]. The actual values and the values predicted by the classification algorithm are shown in Table 1. Performance evaluation criteria of classification algorithms are shown in Table 1 below [21-23].

TABLE I. CONFUSION MATRIX

		Prediction	
		True	False
Actual	True	TT	TF
	False	FT	FF

Accuracy and error value of the model generated by the classification algorithms according to Table 2 are given by Eq. (2) and Eq. (3), respectively [22].

$$Accuracy = \frac{TT + FF}{TT + TF + FT + FF} \quad (2)$$

$$Error = 1 - Accuracy$$

Precision and sensitivity values of the model generated by the classification algorithms according to Table 2 are given by Eq. (4) and Eq. (5), respectively [22].

$$Precision = \frac{DD}{DD + YD}$$

$$Sensitivity = \frac{TT}{TT + TF}$$

Specificity and F-measure values of the model generated by the classification algorithms according to Table 2 are given by Eq. (6) and Eq. (7), respectively [22].

$$Specificity = \frac{YY}{YY + YD} \quad (6)$$

$$F - measure = \frac{2 \times Sensitivity \times Precision}{Sensitivity + precision} \quad (7)$$

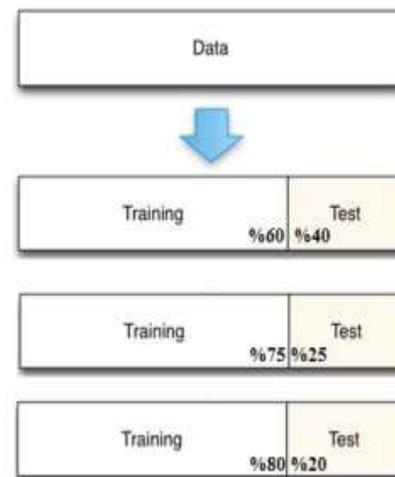


Fig. 1. Test and training set separation with the hold out method.

Models are created using classification algorithms to make estimations on the data set. To see the performances of the classification models, the classification models are divided into training and test data. Various methods have been developed for this splitting process. Among these methods, the holdout method was used in this study. In the hold

out separation, the test and the training data sets are divided once with a specific ratio. Figure 1 shows the flow of this method.

### Conclusion

The performances of different data mining programs were examined by establishing models with classification algorithms. Different results were obtained in the four programs used. However, the algorithm that gives the best result in all programs was the decision tree algorithm. This result suggests that decision tree method gives better performance regardless of the program used. Further studies are also needed to support this result by working with data other than the bank data set. This is a subject of another research to be further investigated in the future.

### References

- [1] S. Moro, P. Cortez, and P. Rita, "A data-driven approach to predict the success of bank telemarketing," *Decision Support Systems*, vol. 62, s. 22-31, 2014
- [2] M.P. Bach, S. Juković, K. Dumiči, and N. Šarlija, "Business client segmentation in banking using self-organizing maps," *South East European Journal of Economics and Business*, vol. 8, no. 2, s. 32-41, 2013.
- [3] S. Sumathi, S. Sivanandam, Introduction to Data Mining and Its Applications, *Springer Science & Business Media*, 2006.
- [4] A. Keramati, R. Jafari-Marandi, M. Aliannejadi, I. Ahmadian, M. Mozaffari, and U. Abbasi, "Improved churn prediction in telecommunication industry using data techniques," *Applied Soft Computing*, vol. 24, s. 994-1012, 2014.
- [5] M.S. Başarslan, F. Kayaalp, "Customer churn analysis with classification algorithms in telecommunication sector. ICAT'17, Istanbul, Turkey, 2017
- [6] J. Han, M. Kanber, Data Mining: Concepts and Techniques, *Morgan Kaufmann*, 2006
- [7] S. Akyokuş, "Veri Madenciliği Yöntemlerine Genel Bakış," TBD Veri Madenciliği Günü sunumu, Doğuş Üniversitesi, 2006
- [8] T. Mitchell, Machine Learning, *McGraw Hill*, New York, 1997.
- [9] P. Harrington, Machine Learning In Action,

- Manning*, New York 2012.
- [10] H. Arslan, "Sakarya üniversitesi web sitesi erişim kayıtlarının web madenciliği ile analizi," Yüksek lisans tezi, Elektronik-Bilgisayar Eğitimi, Sakarya Üniversitesi, Sakarya, Türkiye, 2008.
- [11] J. Han, M. Kamber, Data Mining: Concepts and Techniques, Morgan Kaufmann Publisher, San Francisco, USA, 2000.
- [12] I.H. Witten, E. Frank, Mark Data Mining, Morgan Kaufmann, Elsevier, San Francisco, 2005.
- [13] M.S. Başarslan, F. Kayaalp, "Telekomünikasyon Sektöründe Müşteri Kayıp Analizi," Bilgisayar Mühendisliği Ana Bilim Dalı, Düzce Üniversitesi, 2017.
- [14] KNIME, Bilgisayar Programı, Konstanz Üniversitesi, Zürih Teknopark, 2004.
- [15] T.T. Bilgin, "Veri akışı diyagramları tabanlı veri madenciliği araçları ve yazılım geliştirme ortamları," Akademik Bilişim'09 Konferansı, Şanlıurfa, Türkiye, 2009.
- [16] YALE, Bilgisayar Programı, Yale Üniversitesi, 2001.
- [17] RAPIDMINER, Bilgisayar Programı, Dortmund Teknoloji Üniversitesi Yapay Zeka Birimi, 2006.
- [18] M. Dener, M. dörterler, ve A. Orman, "Açık kaynak kodlu veri madenciliği programları: Weka'da örnek uygulama," Akademik Bilişim'09 Konferansı, Şanlıurfa, Türkiye, 2009.
- [19] K. Dahiya, S. Bhatia, "Customer churn analysis in telecom industry," Reliability, Infocom Technologies and Optimization (ICRITO) (Trends and Future Directions), 2015 4th International Conference on IEEE, Noida, India, s. 1-6, 2015.
- [20] A. Demirci. (2015, 28 Eylül). Data Driven Kavramı [Online]. Erişim: <http://devveri.com/kategori/haberler>.
- [21] N. Japkowicz, "Performance evaluation for learning algorithms," International Conference on Machine Learning, Edinburg, Scotland, 2012.
- [22] M. Clark, "An Introduction to machine learning with Applications in R," Lecture Notes, University of Notre Dame, 2015.
- [23] P. Flach, "The many faces of ROC analysis in machine learning," Lecture Notes, University of Bristol, 2004.