



Empirical Study on Hybrid technique to predict customer's intention

Pallabi Baruah¹, Dr. Bhairab Sarma²

¹Assistant Professor, Department of Computer Science, Asian Institute of Management and Technology, Guwahati, Assam

²Associate Professor, Department of Computer Science and Electronics, University of Science & Technology, Meghalaya

Abstract: Customer churn is one of the most important metrics for a growing business to evaluate. While it's not the happiest measure, it's a number that can give a service oriented company, the hard truth about its customer retention. The churn rate, also known as the rate of attrition or customer churn, is the rate at which customers stop doing business with an entity. It is most commonly expressed as the percentage of service subscribers who discontinue their subscriptions within a given time period. It is also the rate at which employees leave their jobs within a certain period. For a company to expand its clientele, its growth rate (measured by the number of new customers) must exceed its churn rate. Purposes behind customer churn might be the disappointment with the nature of administration, high costs, ugly plans, no comprehension of the administration design, awful help, and so on. This paper gives emphasis on a hybrid data mining techniques which involves Decision Tree and Logistic Regression which is used to create the prediction model to predict the intention of the customers in service oriented organization.

Keywords: Churn, prediction, sectors, hybrid, logistic regression, decision tree.

INTRODUCTION

Customer churn is an agonizing reality that affects all businesses at some point. In any industry, it refers to the shift of customers from one organization to another. Major cause of revenue loss in this sector is because of customer churns. It is important for a lasting and sustainable business growth to understand what has caused previously loyal customers and users to abandon ship and find a new provider. Moreover, the client stops contract without the point of changing to a contender. Explanations behind this are changes in the circumstance that makes unimaginable for the client from additionally requiring the administration, e.g. money related issues, prompting inconceivability of installment or change of the topographical area of the client to a place where the organization isn't accessible or the administration is inaccessible. Now and then the organization additionally stops or pulls back the agreement because of organization approach reason. A company can compare its new subscribers versus its loss of subscribers to determine both its churn rate and growth rate to see if there was overall growth or loss in a specific time period. While the churn rate tracks lost customers, the growth rate tracks new customers. Prediction techniques gives the expectation about clients who liable to stir sooner which means to recognize such churners and to finish some positive exercises to restrain the upset effect. The Customer Relationship Model (CRM) office to avert endorsers who are probably going to stir in future by taking the required maintenance approaches to draw in the feasible churners and to hold them. Along these lines, the potential loss of the organization could be dodged. The contribution for this issue incorporates the information on past requires every portable supporter, together with all individual and business data that is kept up by the specialist organization. Likewise, for the preparation stage, marks are given as a rundown of churners. After the model is set up with most critical precision, the model must have the ability to envision the once-over of churners from the real dataset which excludes any unsettle stamp. In the point of view of learning disclosure process, this issue is classified as prescient mining or prescient demonstrating.

Types of Customer Churns

- 1 Active churner (Volunteer):** those clients who need to stop the agreement and move to the following supplier.
- 2 Passive churner (Non-Volunteer):** At the point when an organization suspends administration to a client.
- 3 Rotational churner (Silent):** Those clients who cease the agreement without the earlier information of the two gatherings (client and friends), where each gathering (e.g. client or organization) may abruptly end the agreement with no warning.

The initial two sorts of stirs can be anticipated effortlessly with the assistance of customary methodologies as far as the Boolean class esteem; in any case, the third kind of agitate is hard to foresee since there is a plausibility of clients who



may beat sooner rather than later for assortment of reasons that are either not known or hard to anticipate. It ought to be the objective of the chief to diminish the stir proportion, for existing clients are the most significant resources for organizations when contrasted with getting new ones. Client beat conduct affects the organization's execution which are condensed as tail: (i) a negative effect on the general execution of the organization, (ii) a potential reason for low deals on the grounds that new/here and now clients purchase less administrations, (iii) encourages contenders to increase disappointed clients with business promotion(s), (iv) prompts income misfortunes, (v) puts negative effect on long haul clients, (vi) builds vulnerability which lessens the proportion of conceivable new clients, (vii) drawing in new clients is more costly than holding existing and (viii) hazard to organization's picture in the focus recognize computerized techniques that can help organizations in the mind boggling assignment of anticipating client beating. If a client maintenance is incorporated, it can build benefits. This features the incomparable significance of empowering client steadfastness and keeping away from client stirs. Keeping up client dedication and beat have turned out to be essential articles of exchange media transmission enterprises, with the goal that proactive maintenance battles can be sent in an offer to hold them. As maintenance crusades are expensive and tedious, watchful arranging is required. As it were, this cost can be diminished by utilizing information mining methods that perform activities in light of found learning. The prompt necessity of the market is to have frameworks that can perform exact 1. Identification of steadfast clients (with the goal that organizations can offer more administrations to hold them) 2. Prediction of churners to ensure that solitary the customers who are needing to switch their authority centers are being engaged for upkeep 3. commend activities that progressions clients from an undesired status, (for example, churners) to a coveted one, (for example, steadfast) while expanding a target work which is the normal net benefit. This researcher proposes to give emphasis on churn prediction and the use of data mining techniques to detect the intention of the customers in the industry. Data mining is a process that finds a small set of precious nuggets from a great deal of raw material. Data Mining makes use of various algorithms to perform a variety of tasks. The research sample data is examined by these algorithms and a model is determined that is close to solving of the problem of formulating a prediction model to detect customer churn

DATA MINING

Data Mining is the process of discovering interesting knowledge from large amounts of data stored either in databases, data warehouses, or other information repositories. It is a method of extraction and analysis of patterns, relationships and information from huge databases. It is primarily concerned with discovering patterns and anomalies within datasets, but it is not related to the extraction of the data itself. The mining process is known as knowledge discovery in databases (KDD). KDD is the non-trivial process of recognizing suitable, new, potentially useful and reasonable patterns in data. The issues involve a discrete valued target variable and the main objective is to determine a subscriber as a potential churning or a potential non-churner. Thus, the KDD function is described as a classification issue. The initial step in predictive modelling is the acquisition and preparation of data. In the data mining process, four subtasks are involved. These are classification, clustering, regression and association rule learning. In addition, in this context, data mining techniques can be classified into the two types of hypothesis and discovery-oriented. Verification methods manage the estimation of a hypothesis designed by an external source. Statistical methods like goodness-of-fit test, t-test of means and analysis of variance fall into this type. These techniques are linked with data mining rather than their discovery-oriented counterparts as data mining issues. The issues involve choosing a hypothesis rather than testing, although discovery methods are employed to recognize patterns in data.

IMPORTANCE OF THE PRESENT WORK

Many companies have one main concern which is related to the customer retention. Customers are likely to churn from one service provider to another. Churn is a word which means change of service provider due to better services and rate which is offered by the different telecommunication, insurance and any service related companies. These companies are thinking of productive ways through which they could identify the customers who have high probability of changing so that they can send proactive customer retention campaign. Data mining techniques are being used to predict the customer churn based on customer behaviour, data and patterns which are derived. Consumer behaviour is the study of individual, or group about their process of selecting and using the product, services, ideas or experiences to satisfy needs. It involves ideas, services and tangible products. The customer behaviour is analysed to making the marketing strategies and public policy. It is very important to understand the psychology of the customer while he/she purchases a product. The stored data contains the information of about the spending behaviour of customer, how much they buy, which day at what time he/she does the shopping, and what they buy most often, in that locality etc. New techniques of data mining are used to understand the Customer Relationship Management(CRM) and with different strategies can be implemented for different sets of customers in different areas. Thereby, it is an efficient way to help the telecommunication companies in order to detect the customers behaviour and take action accordingly.



REVIEW OF RELATED WORK

Data mining is a process that discovers the knowledge or hidden pattern from large databases. It is one of the core processes of Knowledge Discovery in databases(KDD)[9] Data mining techniques are applied in telecom database for various purposes. The data generated by the telecom companies are broadly grouped into three types. They are Customer Data(Demography), Network data and Bill data[23].For churn prediction, statistical methods had been used for many years but Data mining techniques has earned more popularity in telecommunication sector. Advance techniques in DM are available for determination of customer churn [18]. There are many developments and designs so far as Datamining techniques are concerned. They are classified based on the type of database, the knowledge to be discovered or the techniques used. Based on database, Data mining system can be classified based on the type of database it is designed for (relational, transaction, multimedia, legacy, web database). Based on techniques, it may be categorized as autonomous knowledge mining, data driven mining, query driven mining, query driven mining, and interactive Data mining technique. Alternatively, it can be classified into underlying mining approach such as generalization based mining, pattern based mining, statistical or mathematical base mining and integrated approached. Based on knowledge, Data mining Systems can discover various types of knowledge, including association, classification, clustering, prediction, including association, classification, clustering, prediction, sequential patterns and decision trees [9] . The Researcher is willing to further the understanding of grounded theories in Prediction model using Data Mining technique in customer churn determination in various companies.

In the telecom industry, there is a huge revenue loss due to the behavior of the customers. Such customers create an undesired and unnecessary financial burden on the company. This financial burden results in to negative impact on the company and ultimately may lead to sickness of the company, Data mining techniques are used for three main purpose in this sector so as to analyze customer loyalty. Firstly, it is churn prediction which is prediction of customers who are at risk of leaving the company due to various reasons. Secondly, Insolvency prediction which is the increasing due bills in any service oriented company. Thirdly, Fraud detection which deals with fraudulent users and their usage patterns. This helps the CRM department to prevent subscribers who are likely to churn in future by taking the required retention policies to attract the likely churners and to retain them. [23].

[4] This paper is to segment airline customers into four groups, set different churn rules to evaluate churn rate and analyze customer churn propensity based on logistic model. With the help of these strategies, the airlines can take positive and effective measures to reduce the company's operating costs and enhance the company's core competencies.

[10] Another research work proposed the Neural Network based Ensemble classifiers which are performing best for the given two datasets of telecom industry and the two datasets used different types of attributes for the churn prediction so we are sure that the same model can be applied to any dataset for acquiring the best prediction result and can save the loyal customer of a company before churn.

[6]A researcher in his study tested 5 different classification methods with a dataset consisting of 57 attributes. Experiments were carried out several times using comparisons between different classes and Support Vector Machine (SVM) with a comparison of 50:50 Class sampling data is the best method for predicting churn customers at a private bank in Indonesia. The results of this model can be utilized by company who will apply strategic action to prevent customer churn.

[11] proposed multiple regressions analysis to predict the customers churn in the telecommunications industry based on recommended features in which the results have shown that the performance of multiple regressions for predicting customer churn is acceptably good. [20] used hybrid supervised and unsupervised techniques in order to deploy to achieve improved churn prediction instead of single classifier resulting in low efficiency. In this work, the orange data set is preprocessed using effective data cleaning methods. After cleaning process, the clustering approach is carried out using unsupervised techniques such as K-Means and weighted K-Means algorithms. The various clusters obtained from the above method are divided into training and testing sets using hold-out method. And each of the cluster training and testing data sets is handled by various supervised data mining techniques. The hybrid mechanism efficiency is measured in terms of sensitivity, specificity, and accuracy.

[1] in their research work propounded that the use of the bank's database imposed some limitations on the present study. implemented the CRISP methodology for predicting customer churn in electronic banking services. The aim of the present study is to identify the features of churners from electronic banking services. [24] implemented and tested well in this study, some work can be improved further. First, in terms of time complexity, FSDTE model is a little time-consuming. To evaluate the time complexities of different models, we compared the average time consumed by ten independent runs of each model under the same conditions

[18] proposed METALAB as the data mining tool which supports many algorithms and Back propagation algorithm was used for the research. [22] tried to overcome the churn prediction problem, using machine learning algorithms and data mining tools. One of the popular tools in the field is Weka which is an open source software for data mining, developed by the University of Waikato in New Zealand. ([23] used the two well-known models, Decision tree and Artificial Neural Network to build the churn prediction model. It has been observed that decision tree model surpasses the neural network



model in the prediction of churn and it is also easy to construct. This research has no steps to analyzed to include retention policies but is related to a prediction model to detect customer churn.

[2] attempted to make use of rough set theory as one-class and multi-class classifiers and investigated the trade-off in the selection of an effective classification model for forecasting customer churn in telecommunication industry. A series of experiments were performed and explored the performance of four different rule generation algorithms. It is also investigated that the true churn prediction rate has also improved from 86% in MCC to 96% in OCC. [8]in their research paper, it is observed that decision tree model surpasses the neural network model in the prediction of churn and it is also easy to construct.

The Telecommunication, Insurance and Banking industry produces enormous quantities of data and is bedeviled with a vast array of imperfect customer information that decision makers need to deal with. Customers regularly port their mobile numbers from one Telecommunication provider to the other thereby making the companies lose large amount of data. A simple predictive model for the prediction of potential churn customers using primary data from customers has been of great importance in the sector.

[7] It is proposed that Decision Tree and Logistic Regression in their research paper for building the churn prediction model. Discussions of the various prediction models and also comparisons of the quality measures of prediction models like regression analysis, decision trees were made. [17] in the study suggests a statistical survival analysis tool to predict customer churn based on comparison between decision trees and logistic regression The proposed model suggests that data mining techniques can be a promising solution for the customer churn management. Using this model, the telecom companies can predict in advance which customers are at risk of leaving, and can target those customers consequently saving a lot of revenues namely the ones which is used for replacing the lost customers and also the ones that are wasted for retaining already loyal customers.

[16] in his research work analyzed a large customer dataset for churn by using Decision Tree Classification Technique. After implementing all the possible variants of decision tree in SPSS it was observed that Exhaustive CHAID technique proved to be more efficient and accurate than others to predict the customers who are likely to churn in nearby future.

[5] clustering is suitable and effective in clustering the data set into two categories together with giving high percentage of non-churners over churners at the cluster-0. As regards the association rule angle Aperiore and FPGrowth, they have been applied to specify the behavior of each customer and that is supposed to help the business organization to focus on certain services or attributes to stop customer churn from happening.

[14] in this research conducts a real world study on customer churn prediction and proposes the use of boosting to enhance a customer churn prediction model. Logistic regression is used in this research as a basis learner, and a churn prediction model is built on each cluster, respectively. The result is compared with a single logistic regression model. Experimental evaluation reveals that boosting also provides a good separation of churn data; thus, boosting is suggested for churn prediction analysis. [3]he projected model is capable of predicting subscriber churn pattern well in prior. We have proposed, hybrid Classification techniques that have shown their superiorities over single technique. Our approach is to compare 2 kinds of hybrid methods by classification with classification and clustering with classification hybrid methods. This paper gives the implementation details of decision tree and logistic regression to predict the customer churn in telecom industry. We implement Hybrid Decision Tree and Logistic Regression Classifier. These techniques provide better results, but it takes lot of time for execution. To tackle this problem; we proposed hybrid Fuzzy unordered rule induction algorithm (FURIA) with Fuzzy C-Means Clustering for predicting customer churn

Based on the study of some research work, researcher is enthusiastic to explore and further do research work related to the customers' behavior and a formulate a prediction model to find out the list of customers who are at high risk to churn in companies.

PROPOSED METHODOLOGY

In the proposed system, Python programming is used to build the model for prediction of customer churn. It is a language and environment for statistical computing and graphics. Python provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, ...) and graphical techniques, and is highly extensible.

Data Preprocessing and preparation

The data is downloaded from IBM Sample Data Sets. Each row represents a customer, each column contains the customer's attributes. In data mining problems, data acquisition and preparation is one of the important phase which is time consuming as it involves data collection, integration and cleaning. Data needs to be arranged, aggregated and characterized. Aggregation leads to the generation of new data in addition to the existing ones. It is cleaned by removing any ambiguities, errors and also fields with null values are excluded. The entire dataset is checked and it involves a lot of analysis before it is ready for being processed.

The model will have three parts –



1. Display performance analysis – which gives a view of the results obtained by applying the two classification prediction methods - The logistic regression and decision tree on the available dataset.
2. Testing – to extract the list of customers which have a high probability to churn from the input, given that the attributes of the input data are same as the available dataset used for training,
3. Training and testing – which develops a model along with generating a churn list if any other type of dataset is provided.

In the display performance analysis, the logistic regression and decision tree is used.

Decision tree

Decision tree is one way to display an algorithm that only contains conditional control statements. They are commonly used in operations research, specifically in decision analysis, to help identify a strategy most likely to reach a goal, but are also a popular tool in machine learning. [15] CART (Classification and Regression Trees) technique is more suited for continuous dependent variable and categorical predictor variable. CART splits the feature space recursively into non-overlapping regions. In order to predict the value of dependent categorical variable, a classification tree is generated. CART incorporates both testing with a test data set and cross-validation to assess the goodness of fit more accurately. CART can use the same variables more than once in different parts of the tree. This capability can uncover complex interdependencies between sets of variables. CART can be used in conjunction with other prediction methods to select the input set of variables.

Some important variables in the dataset are selected for decision tree visualization and confusion matrix table is being created based on the selected variables and thus predictions are being made.

Logistic regression

Logistic regression predicts the probability of an outcome that can only have two values (i.e. a dichotomy). The prediction is based on the use of one or several predictors (numerical and categorical). A linear regression is not appropriate for predicting the value of a binary variable for two reasons: A linear regression will predict values outside the acceptable range (e.g. predicting probabilities outside the range 0 to 1) . Since the dichotomous experiments can only have one of two possible values for each experiment, the residuals will not be normally distributed about the predicted line.

On the other hand, a logistic regression produces a logistic curve, which is limited to values between 0 and 1. Logistic regression is similar to a linear regression, but the curve is constructed using the natural logarithm of the “odds” of the target variable, rather than the probability. Moreover, the predictors do not have to be normally distributed or have equal variance in each group.

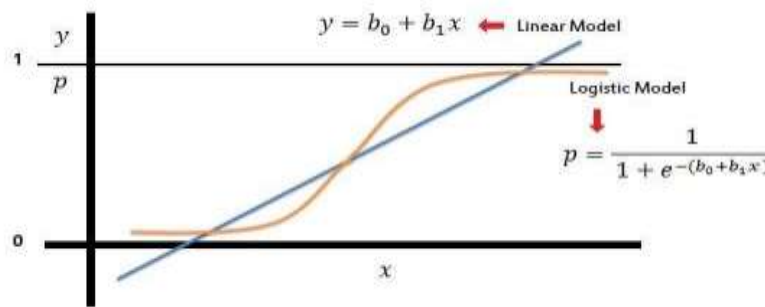


Figure – 1

In the logistic regression the constant (b_0) moves the curve left and right and the slope (b_1) defines the steepness of the curve (Figure-1). By simple transformation, the logistic regression equation can be written in terms of an odds ratio.

$$\frac{p}{1-p} = \exp(b_0 + b_1x) \tag{1}$$

Finally, taking the natural log of both sides, we can write the equation in terms of log-odds (logit) which is a linear function of the predictors. The coefficient (b_1) is the amount the logit (log-odds) changes with a one unit change in x .

$$\ln\left(\frac{p}{1-p}\right) = b_0 + b_1x \tag{2}$$

As mentioned before, logistic regression can handle any number of numerical and/or categorical variables. The output of p is normally a numeric value between 0 and 1 which gives the probability of the outcome. That is, if the output is 0.6



then there is a 60% probability of getting 1 which is predicted to be a positive case and if the output is below 0.5 then there is a probability of getting 0 and it is predicted to be a negative case. The churn prediction model is built on logistic regression and decision tree algorithms and on selecting the correct combination of attributes, training and test datasets can provide accurate outcomes for prediction of customer.

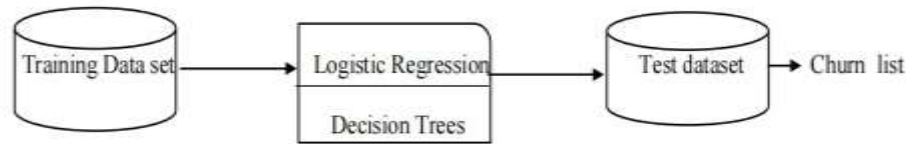


Figure-2. The Churn Prediction Model

The data is splitted into training and testing sets. The data is fitted on to the Logistic Regression model and analysis is done as variables are being added.

CONCLUSION

The experimental result show that the planned hybrid technique yields accurate result. The proposed model gives a factual survival examination apparatus to predict churn in light of correlation between decision trees and logistic regression. Choosing the privilege mix of properties and settling the correct edge qualities may deliver more precise results. The selection of attributes is verified to increase the performance of DT.

REFERENCES

1. Abbas Keramati1, H. G.. Developing a prediction model for customer churn from electronic banking services using data mining. Keramati et al. Financial Innovation (2016).
2. Adnan Amin, C. K. Customer Churn Prediction in Telecommunication Industry: With and without Counter Example. A. Gelbukh et al. (Eds.): MICAI 2014, Part II, LNAI 8857, pp. 206–218, 2014. © Springer International Publishing Switzerland 2014.
3. Atul Sunil Choudhari1, D. M. Predictive To Prescriptive Analysis For Customer Churn in Telecom Industry Using Hybrid Data Mining Techniques. IEEE. (2018).
4. Hong-Yu. Research on Customer Churn Prediction Using Logistic Regression Model. © Springer Nature Switzerland AG 2019.
5. Ibrahim M.M.Mitkees, A. P. Customer Churn Prediction Model using Data Mining Techniques. IEEE. (2017).
6. Ketut Gde Manik Karvana, S. Y. Customer Churn Analysis and Prediction Using Data Mining Models in Banking Industry. IEEE. (2019).
7. Kiran Dahiya, S. B. Customer Churn Analysis in Telecom Industry. IEEE. (2015)
8. Balasubramaniam, M. Churn Prediction in Mobile Telecom System using Data Mining Techniques. International Journal of Scientific and Research Publications. (2018)
9. Maninderjeet Kaur, P. A comparative study of techniques to predict customer churn in telecommunication industry. International Research Journal of Engineering and Technology (IRJET). (2017)
10. Mehpara Saghir, 2. B. Churn Prediction using Neural Network based Individual and Ensemble Models. Proceedings of 2019 16th International Bhurban Conference on Applied Sciences & Technology (IBCAST). (2019).
11. Mohd Khalid Aang, wM. N. . Data Mining for Churn Prediction: Multiple Regressions Approach . EL/DTA/UNESST 2012, CCIS 352, pp. 318–324. (2012)
12. Suneel ramachandra joshi; Sharadhi joshi. "Continuous improvement strategy- lifeblood of organizations". International Research Journal on Advanced Science Hub, 2, Special Issue ICAET 11S, 2020, 42-47. doi: 10.47392/irjash.2020.231
13. Salini Suresh; Suneetha V; Niharika Sinha; Sabyasachi Prusty; Sriranga H.A. "Machine Learning: An Intuitive Approach In Healthcare". International Research Journal on Advanced Science Hub, 2, 7, 2020, 67-74. doi: 10.47392/irjash.2020.67
14. Kousalya M.; Lakshi E.; Mukesh Kanna R.K.; Pravin G.; Prabhu T.. "Automatic Restaurant Food Ordering Menu Card". International Research Journal on Advanced Science Hub, 3, Special Issue ICARD-2021 3S, 2021, 7-12. doi: 10.47392/irjash.2021.052
15. Vikas R Gangadhar; Ajim Shaikh. "Cloud Technology and Return on Investment (ROI)". International Research Journal on Advanced Science Hub, 3, Special Issue ICEST 1S, 2021, 73-79. doi: 10.47392/irjash.2021.023
16. Priyanka Kumari. "A Study of Customer Preference and Attitude towards Online Shopping In Bihar". International Research Journal on Advanced Science Hub, 3, Special Issue ICEST 1S, 2021, 7-11. doi: 10.47392/irjash.2021.012
17. Mümin Yıldız, S. A. Customer Churn Prediction in Telecommunication. IEEE. (2015).
18. Nabareseh, I. S.. Predictive analytics: a data mining technique in customer churn management for decision making. Tomas Bata University in Zlín. (2017)
19. Ning Lu, H. L. A Customer Churn Prediction Model in Telecom Industry Using Boosting. IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS, VOL. 10, NO. 2., (2014, May).
20. Nisha Saini, M. D.. Churn Prediction in Telecommunication Industry using Decision Tree. International Journal of Engineering Research and Technology. (2017)
21. Preeti K. Dalvi, S. K. Analysis of Customer Churn Prediction in Telecom Industry using Decision Trees and Logistic Regression. Symposium on Colossal Data Analysis and Networking (CDAN). Pune: IEEE. (2016).
22. Rahul J. Jadhav, U. T. Churn Prediction in Telecommunication Using Data Mining Technology. International Journal of Advanced Computer Science and Applications.(2011).
23. Saini, N. Churn Prediction in Telecommunication Industry using Decision Tree. International Journal of Engineering Research & Technology (IJERT), ISSN: 2278-0181. (2017).



24. Sivasankar, J. V. Improved Churn Prediction Based on Supervised and Unsupervised Hybrid Data Mining System. © Springer Nature Singapore Pte Ltd. 2018.
25. THYAGARAJAN, K. Generic Model for Customer Relationship Management Based on Back Propagation Algorithm. THIRUCHIRAPPALLI. (2014).
26. Utku Yabas, H. C. . Customer Churn Prediction for Telecom Services. IEEE 36th International Conference on Computer Software and Applications. (2012)
27. V. Umayaparvathi, K. I. Applications of Data Mining Techniques in Telecom Churn Prediction. International Journal of Computer Applications. (2017).
28. Wang, J. X. Feature-selection-based dynamic transfer ensemble model for customer churn prediction. Received: 20 November 2012 / Revised: 11 November 2013 / Accepted: 04 December 2013 © Springer Verlag London 2014.