

**An Intelligent System for Churn Prediction and Customer Retention:
The Case of a Telecommunications Company**

by

Parth Sarangi

A thesis submitted in partial fulfillment of the requirements for the
degree of Master of Engineering in
Information Management

Examination Committee: Dr. Vatcharaporn Esichaikul (Chairperson)
Dr. Matthew N. Dailey
Prof. Sumanta Guha

Nationality: Indian
Previous Degree: Bachelor of Technology in Electronics and Communication
National Institute of Technology Srinagar, India

Scholarship Donor: AIT Fellowship

Asian Institute of Technology
School of Engineering and Technology
Thailand
April 2018

Acknowledgments

I would like to thank all my family members for showing great support towards my dream of pursuing a Master's degree. I express my sincere gratitude to my advisor Dr. Vatcharaporn Esichaikul for showing faith in me, with encouragements and guidance. Without her support I would not have been able to conceive and accomplish this thesis.

I also express my sincere thanks to my Professors Dr Matthew N. Dailey and Dr. Sumanta Guha for guiding me in my academics and for helping me understand the fundamental concepts of machine learning and web development.

I would like to extend my thanks to my seniors Mr. Manish Taneja, and Mr. Siraj Muhammad for their invaluable and selfless guidance.

I wish to express humble thanks to the staff of my Institution for their support and assistance in my academic pursuit.

Lastly, I want to extend my gratitude to Asian Institute of Technology (AIT) for providing a fellowship and enabling me for pursuing my degree.

Above all I am grateful to the omnipresent spirit of the universe which guides me in achieving this milestones.

Abstract

The telecommunications industry is very competitive in most of the developed and developing countries. A few companies operate to provide numerous services to a huge consumer base. Rapid technological advancements in the ICT sector has increased the availability and affordability of mobile telephony devices. With an increasing adoption of mobile telephony devices there is a greater increase in the demand for mobility services. Services such as phone calls, sms, internet etc., have increased over the past decade. With the growing demand and growing consumer base, services providers have reduced prices of these services. Profitability of companies is decided by the number of consumers. In these highly competitive market, customer satisfaction and retention is of high importance.

In this study an intelligent churn prediction and customer retention (ICPCR) system is developed. An open source dataset of call detail records (CDR), with 5000 records and 21 features is selected for purpose of system design and model development. Eight data mining methods are employed to generate prediction models. Prediction models such as decision tree, random forest, support vector machines, neural networks and naive bayes are compared for performance and evaluated based on accuracy, sensitivity, specificity and positive prediction metrics. Based on the performance decision tree is selected as the prediction engine in implementation of ICPCR. The customer retention system is designed on decision table rules and upselling techniques are suggested for every customer predicted as churner.

This study is designed to facilitate the customer service representative with a mechanism to visualize customer call detail data. In addition, it also enables them to predict the churning status of current customers, and also presents them with an option to control churning by suggesting marketing solutions or product benefits. This study presents a web based approach with KPI dashboards, charts and maps, OALP for roll-up and drill-down analysis, prediction and finally retention strategies for controlling customer churn.

Table of Contents

Chapter	Title	Page
	Title Page	i
	Acknowledgments	ii
	Abstract	iii
	Table of Contents	iv
	List of Figures	vi
	List of Tables	vii
1	Introduction	1
	1.1 Overview	1
	1.2 Problem Statement	2
	1.3 Objectives	3
	1.4 Limitations and Scope	3
	1.5 Thesis Outline	3
2	Literature Review	4
	2.1 Customer Churn & Retention	4
	2.2 OLAP & Datawarehouse	4
	2.3 Data Mining	5
	2.4 Model Evaluation Metrics	13
	2.5 Review of Selected Research Papers	15
	2.6 Summary of Selected Research Studies	16
3	Methodology	21
	3.1 Research Methodology	21
	3.2 Data Preprocessing and Datawarehouse Development	22
	3.3 Development and Evaluation of the Prediction Models	24
	3.4 To do	24
	3.5 System Development & Evaluation	26
4	Data preprocessing & Data warehouse	28
	4.1 Data collection	28
	4.2 Meta-data evaluation	28
	4.3 Data cleaning	30
	4.4 Data-warehouse design	30
	4.5 ETL process	35
5	Model Development and Evaluation	36
	5.1 Prediction models selected	36
	5.2 Decision tree : rpart's classification tree	37
	5.3 Decision tree : ctree's conditional inference	40
	5.4 Random forest	43

5.5	Decision tree C5.0	44
5.6	Decision tree C5.0 boosted	47
5.7	Naive bayes	49
5.8	Support vector machine	50
5.9	Neural network	51
6	Results and Discussion	53
6.1	Model performance comparisons	53
6.2	Discussion	58
7	System Development and Evaluation	59
7.1	ICPCR web service	59
7.2	Dashboard	60
7.3	Data explorer	61
7.4	OLAP	62
7.5	Visualizations	63
7.6	Map	65
7.7	Prediction functionality	66
7.8	Retention functionality	69
8	Conclusion and Recommendations	70
8.1	Conclusion	70
8.2	Recommendations	71
	References	72

List of Figures

Figure	Title	Page
2.1	OLAP Solution - Apache Kylin	5
2.2	Mammal classification problem	6
2.3	A sample neural network	7
2.4	Kohonen SOM	8
2.5	Select the Right Mining Technique	9
2.6	Another approach to select the Data Mining. Reprinted from Scikit	10
2.7	PreditionIO Engine interaction with Apps and Prediction Engine	11
2.8	R Shiny architecture	12
2.9	Confusion Matrix	14
3.1	Research Methodology	22
3.2	Data preprocessing	22
3.3	OLAP Star Schema	23
3.4	The Intelligent Churn Prediction Architecture	26
4.1	Star schema of churn data-warehouse	34
5.1	Rpart decision tree	39
5.2	Ctree - conditional inference tree	42
5.3	C 5.0 decision tree	46
5.4	C 5.0 decision tree	48
5.5	Neural network 2 hidden layers	52
6.1	Model ROC curves	54
7.1	Dashboard - Data explorer	60
7.2	Dashboard - Key performance index	61
7.3	Dashboard - field analysis	61
7.4	OLAP - functionality helps to investigate data via drill-down and roll-up	62
7.5	Dashboard - Bar plot of Total, Churner and non-Churner customer counts region wise	63
7.6	Dashboard - Bar plot of Churner customers to Non-churners State wise	64
7.7	Dashboard - Map to visually show the density of customers state wise	65
7.8	Map visualization of churners and non-churners state wise	65
7.9	Dashboard - Upload of file for prediction	66
7.10	Dashboard - View uploaded file contents	67
7.11	Dashboard - View prediction results	68
7.12	Dashboard - View retention strategies	69

List of Tables

Table	Title	Page
1.1	Approx. subscriber counts(in millions) of select companies in Indian telecom industry.	2
2.1	Previous literature reviewradioButtons("disp", "Display", choices = c(Head = "head", All = "all"), selected = "head")	16
4.1	Meta data description	29
4.2	Dimensional table - t_day fields	31
4.3	Dimensional table - t_eve fields	31
4.4	Dimensional table - t_night fields	31
4.5	Dimensional table - t_state fields	31
4.6	Dimensional tables - t_pid fields	32
4.7	Dimensional table - t_acct_len fields	32
4.8	Dimensional table - t_vm fields	32
4.9	Dimensional table - t_intl fields	32
4.10	Fact table fields	33
5.1	Confusion matrix rpart decision tree	37
5.2	Accuracy, Sensitivity, Specificity - Rpart tree	38
5.3	Confusion matrix ctree decision tree	41
5.4	Accuracy, Sensitivity, Specificity - Ctree tree	41
5.5	Confusion matrix random forest	43
5.6	Accuracy, Sensitivity, Specificity - Random forest	43
5.7	Confusion matrix c5.0	45
5.8	Accuracy, Sensitivity, Specificity - C5.0 tree	45
5.9	Confusion matrix c5.0 Boosted	47
5.10	Accuracy, Sensitivity, Specificity - C5.0 Boosted tree	47
5.11	Confusion matrix naive bayes	49
5.12	Accuracy, Sensitivity, Specificity - Naive Bayes	49
5.13	Confusion matrix support vector machines	50
5.14	Accuracy, Sensitivity, Specificity - Naive Bayes	50
5.15	Confusion matrix neural networks	51
5.16	Accuracy, Sensitivity, Specificity - Naive Bayes	51
6.1	Performance comparisons for prediction models	53
6.2	Decision table for classification	55

Chapter 1

Introduction

1.1 Overview

In the past two decades, many nations are witnessing growth in telephonic services due to availability of affordable cellular devices and increasing fidelity of mobile services. Major policies of deregulation by governments have encouraged private corporations to invest funds and support invention of improved technologies. Telecommunications infrastructure and services are the major contributors to the economic prosperity of any country (Cronin et al., 1993). The telecom industry is largely customer service oriented with goals of loyalty, retention and satisfaction (Gerpott et al., 2001). The major source of revenue is from direct selling of cellular and Internet services. The companies involved in delivering the services have invested in expensive infrastructures and software systems.

Over the period of time most telecom organizations provide almost the same service and similar value proposition to the customer. Companies experience high customer defection when competitors bring in new offers, services and technologies. Incumbent telecom operators face consumer churning on a regular basis.

Profitable telecom companies generally have a large customer base and their databases hold a wealth of information. It has become imperative that company leaders need to look into their own subscriber base and study the trends that can reveal customer behavior. The biggest asset for companies in the services domain is the customer (Poel & Lariviere, 2004). Thus companies are resorting to data mining techniques and tools to predict customer churn prediction (Berson et al., 1999). From previous data mining techniques it is inferred that it is more profitable to retain and service existing users than to bring in new subscribers (Reinartz & Kumar, 2003). A small effort to retain customers results in major contributions.

Reports published by TRAI shows the mobile phone subscriptions for each telecom operator in India (TRAI - Telecom Regulatory Authority of India, n.d.). Data accessed from the reports is tabulated as below:

Operators	Customer Count in Aug 2016	Increase or Decrease in Period				Customer Count in Dec 2016
		Aug - Sep	Sep – Oct	Oct – Nov	Nov – Dec	
Airtel	257	2	2	1	2	265
Vodafone	200	0.5	1	0.8	1.8	204
Tata Indicom	58	- 1	- 1	- 1	- 1.6	52
Reliance Jio	0	15	19	16	20	72

Table 1.1: Approx. subscriber counts(in millions) of select companies in Indian telecom industry..

It can be deduced from this report that “Tata Indicom” is continuously loosing customers and “Airtel” & “Vodafone” are adding new subscribers at relatively the same rate as they did before. Whereas “Reliance Jio”, a new entrant is experiencing an extraordinary influx of customers so much so that it almost crossed the numbers held by Tata Indicom in Aug 2016.

This thesis presents an intelligent system which predicts customer churn, helps managers and decision makers to identify the valuable proportion for customer retention strategies. The thesis proposes a system supplemented by a data warehouse on the back-end and a visualizations dashboard as the front-end for decision makers. The predictive model is devised after comparison of prediction performance between Decision tree, Support vector machine and neural networks. The proposal is to build a single system as opposed to using separate softwares for prediction, data manipulation and displaying performance indicators.

1.2 Problem Statement

The telecommunication industry’s income is based primarily on the sale of services to customers. A company’s income can dwindle severely if the mindset of its customers changes. As of this decade we have witnessed a growth of smart-phones and so the need to consume data has increased. Ever so often rivals advertise customer centric plans. Internet service providers are trying to woo customers with free, limited, high speed, unlimited, day only, night only and various other the Internet data campaigns. In the recent history, in Indian telecommunications market, the incumbent operators like Airtel, BSNL, Vodafone, Idea Cellular lost plenty of customers to a new entrant, Reliance Jio. Jio launched its services September 5th 2016. It has been reported that Jio has signed about 72 million customers for its paid services that were free in the past. (Reuters, 2017). This shows the loyalty factor among the customers staying with Reliance Jio. Thus identification of the correct customer segment and understanding their current and future needs is a proactive decision that needs to be taken by company’s management. If leaders are tardy and resist change, they could leave their customers dry and sulky. This would obviously result in customer defection and ultimately loss in revenue.

1.3 Objectives

The overall objective of the thesis is to develop an intelligent system for churn prediction and customer retention (ICPCR).

The specific objectives of the thesis are to:

1. Design models and evaluate their churn prediction performance.
2. Build the system of intelligent churn prediction and customer retention system.
3. Evaluate the system for reliable performance.

1.4 Limitations and Scope

There are many models available for churn prediction. The scope of this thesis is to build the system based on three data mining predictive models viz., Decision tree, Support Vector Machines, Artificial Neural Network tentatively.

1.5 Thesis Outline

The organization of this dissertation is as follows:

- In Chapter 2, the literature review is explored.
- In Chapter 3, the methodology of system development is proposed.
- In Chapter 4 stage of data pre-processing is discussed.
- In Chapter 5, model development is discussd.
- In Chapter 6, results are discussed.
- In Chapter 7, system development is explained and screenshots of the functionalities are shown.
- In Chapter 8, conclusion of the thesis is presented with future recommendations.

Chapter 2

Literature Review

This chapter introduces concepts, technologies, techniques, consulted papers and articles pertaining to the core concepts of Customer Churn & Retention, OLAP & Datawarehouse, Data mining, Model evaluation metrics, Review of selected papers and Summary of selected papers.

2.1 Customer Churn & Retention

Customers are the most volatile asset of a services based company. Many frequently churn in search of better services. Customers are frivolous and those with prepaid or prepay plans are most unfaithful. Companies are generally in profit if they are able to retain customers and it pays off to almost six times (Bhattacharya, 1998). Customers spending longer durations with a company are not easily churned and would not be affected by marketing strategies of rival companies. These customers are valuable to the company and generate profit in revenue. Research studies have shown that long standing customers would be engaged in influencing newer customers to buy into a contract with their service provider (Mizerski, 1982).

The ARPU of a stable customer is high compared to that of a churning customer. Thus marketing managers are focusing on advertising competitive products to retain customers from churning. The loss of capital due to a defecting customer is higher than the cost of retention. As per Forbes, Nov 11, 2013, earnings can swing positively by about 10 % if customers are successfully retained.

2.2 OLAP & Datawarehouse

Systems and companies are ever expanding. They are collecting data at unprecedented rates. Managing data becomes easier with the implementation of Data-warehouse. In many cases the database of a company is segregated into different schema's. Segregation of schema's helps to avoid necessary access privileges and grants confusion. It also helps to maintain the organizational level of segregation in the database, ie., the HR department tables will be unaccessible to an accounts official and vice versa. But company leaders and decision makers should be accessing specific key counts and aggregations from all of their departments. A collection of tables sourcing data from their individual units.

OLAP - Online Analytical Processing is an extension of Data-warehouse technology (Han, 1997). Olap consists of four main processes viz., Drill-down, Roll-up, slice and dice. Multi-demensional data can be fetched by OLAP from the Datawarehouse, and the unit of this is called the OLAP cube. There are two types of OLAP - MOLAP & ROLAP. MOLAP Multidimensional OLAP is a solution used widely.

One very famous open-source OLAP solution is the Kylin™ (Kylin, n.d.). Shown in Figure 2.1 is the architecture of the product.

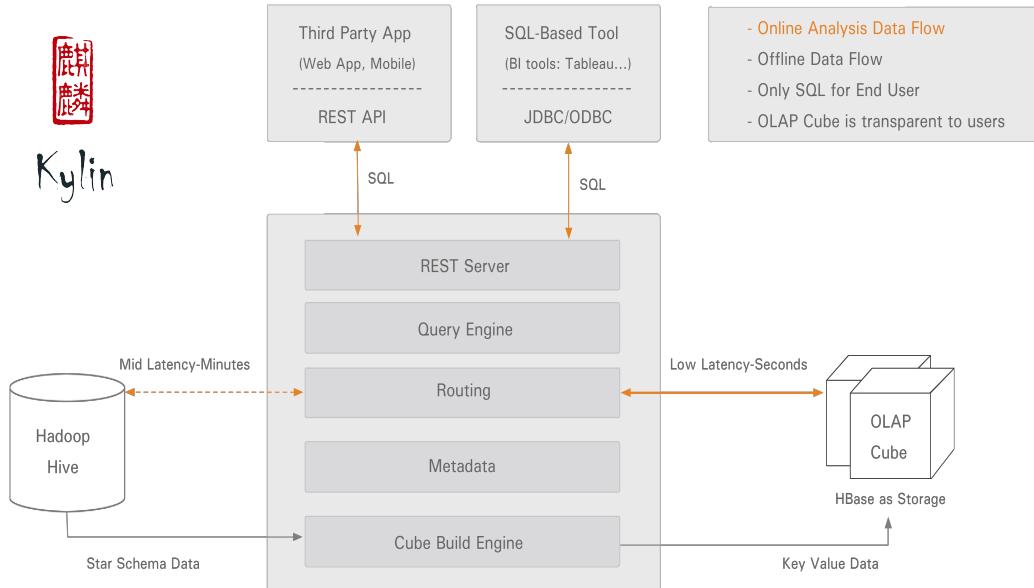


Figure 2.1: OLAP Solution - Apache Kylin.

2.3 Data Mining

Data mining is the process of extracting useful trend and patterns from structured and unstructured sources of data. Sometimes many academicians refer to it as KDD (Knowledge Discovery in Databases). John Naisbett (author of famous 'Megatrends') said "We are drowning in information but starved for knowledge." There are various techniques to perform data mining and these can be broadly classified into two categories Supervised Learning, Un-Supervised Learning. A very common terminology used in the data science field is of machine learning and it also used instead of data mining.

2.3.1 Supervised Learning

This part of the data mining consists of classification and regression algorithms. Control and dependent variables of the given data are known entities. The use of these algorithms is to predict the outcome given past data. These algorithms have to be trained with a set of data and then they have to be tested. After reaching certain acceptable level of accuracy, these algorithms are used for prediction.

Below are some of the Supervised learning techniques :

- Linear regression : The prediction of dependent variable is done given the value of known variable. There is only 1 dependent variable. For example, $y = \beta_0 + \beta_1 x + \varepsilon$
 y = dependent variable, x = independent variable
- Multiple regression : is an extension of the linear regression but has more number of independent variables.
- Nonlinear regression : there are two variables but they are related in a curvilinear fashion i.e., not governed by the straight line equations.
- Logistic regression : A regression based modeling technique, which is better than linear regression when more variables are considered. Output variable is categorical in nature.
- Decision tree : This is a classification algorithm which when plotted resembles an upside down tree structure. Given that a set of data has many attributes and there is a need to classify them, a decision tree is very suitable method to do so. There are many types of decision trees like the ID3, CART, C4.5 and C5.0. In Figure 2.2 a simple DT for mammal classification model is shown. A decision tree can be designed using **Hunt's Algorithm**.

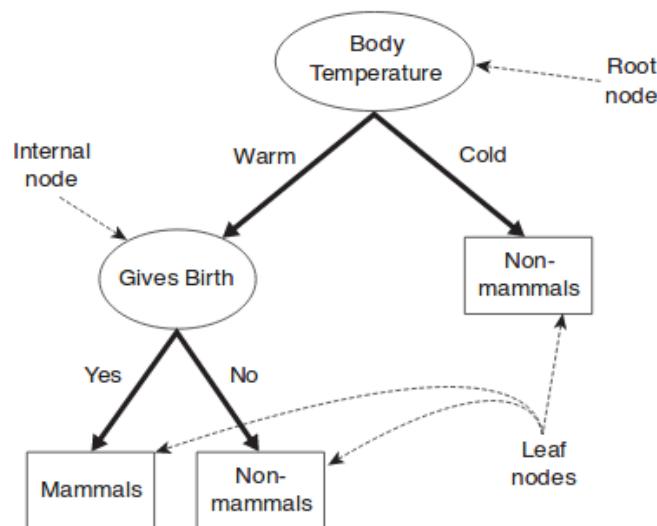


Figure 2.2: Mammal classification problem.

- Random Forest : This technique can be used for both classification or regression type problems. A random forest is combination of many decision trees. In some cases random forest is sometimes very accurate.
- Support Vector Machine : This is a classifier technique where the data is segregated by generating hyperplanes. If there are n-features in the data then there have to be n-hyperplanes. The best classification is the hyperplane which clearly separates the data points.
- k-Nearest Neighbors : A learning algorithm that classifies the data into clusters nearest to them. The euclidean distance or manhattan distance could be some of the methods to find the nearest cluster. It is sometimes considered a lazy learning algorithm.

- Naive Bayes : This is an classification rule working on the probabilistic Bayes theorem.

$$P(H|X) = P(X|H)P(H)/P(X).$$
- Artificial neural networks : Neural networks are classification methods modeled after neurons (karpathy@cs.stanford.edu, n.d.). There are many layers with nodes Figure 2.3. There are many types of neural networks viz., Feed Forward NN, Radial Bias function, Recurrent NN, Backpropagation NN, Perceptron etc. Neural networks are very fast learners.

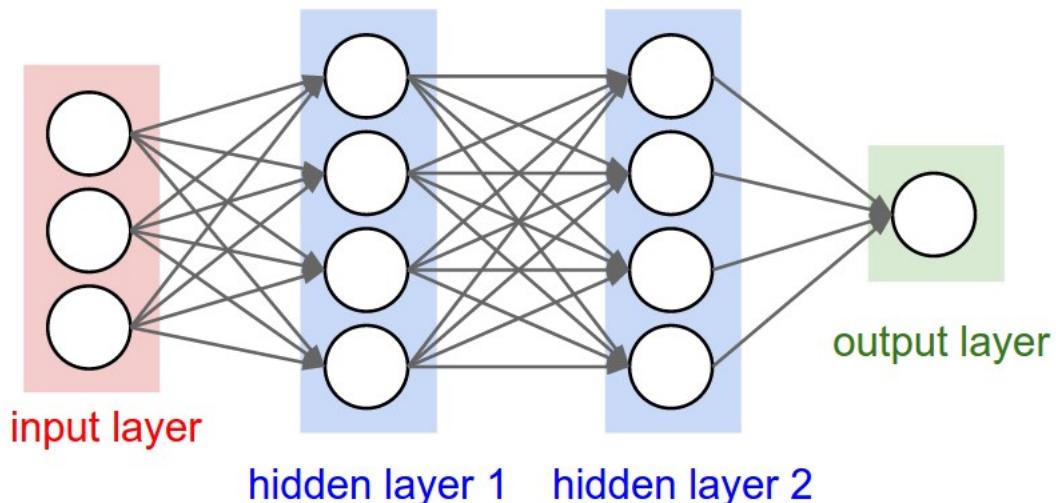


Figure 2.3: A sample neural network.

2.3.2 Un-Supervised Learning

The clustering and association techniques in data mining are grouped into Un-Supervised learning. The output variables are not known. Below are some Unsupervised class of algorithms :

- K-means clustering : it is a means of clustering a set of data points with some k centroids. For each data point the distance is calculated and the nearest centroid is chosen and data point is associated with that cluster. After every iteration of cluster formation a new centroid is calculated and the distance of the data points are taken. The clusters are reformed and the iteration is performed till no data point movement happens.
- Apriori clustering : Here in the A priori algorithm is used to create the clusters. A priori is used for frequent item set mining states that sets of items are frequent if the items themselves are frequent.
- Hierarchical clustering : This is a clustering method in which large clusters are further segregated into smaller clusters. This is the Divisive type of HC. In the Agglomerative type of HC, the nearby clusters are joined to form larger clusters. A Dendrogram is used to graphically represent the clusters.

- Hidden Markov models : These are used to analyze or predict time series problems in fields of speech, language, medicine, and robotics. Core of the technique is formed on the foundations of Bayes Network. In a markov chain a future state depends only on the current state. It is called Hidden because only certain measurements can be seen of the states, not the states itself. Particle filter and Kalman filter are HMM's.
- Self organizing maps (SOM) : This is a type of neural network. Types are of Vector Quantizer or Kohonen SOM. In Figure 2.4 is an illustration of an Kohonen SOM.

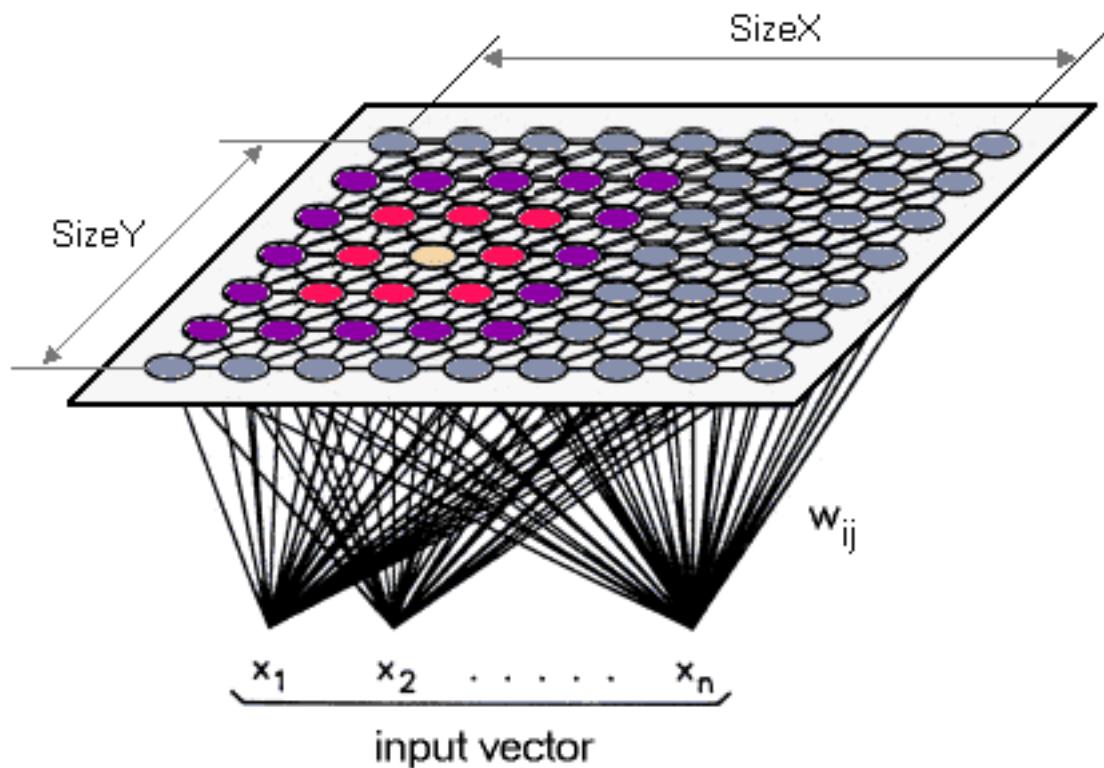


Figure 2.4: Kohonen SOM.

2.3.3 Selecting the Right technique

It is of utmost importance that a data scientist select the important mining technique. Of all the process involved in the knowledge discovery process, selection of algorithm is quite difficult. Figure 2.5, from “Choosing the Right Data Mining Technique: Classification of Methods and Intelligent Recommendation” (Gibert et al., 2010) shows the approach which could be taken to select between the various models available for data mining.

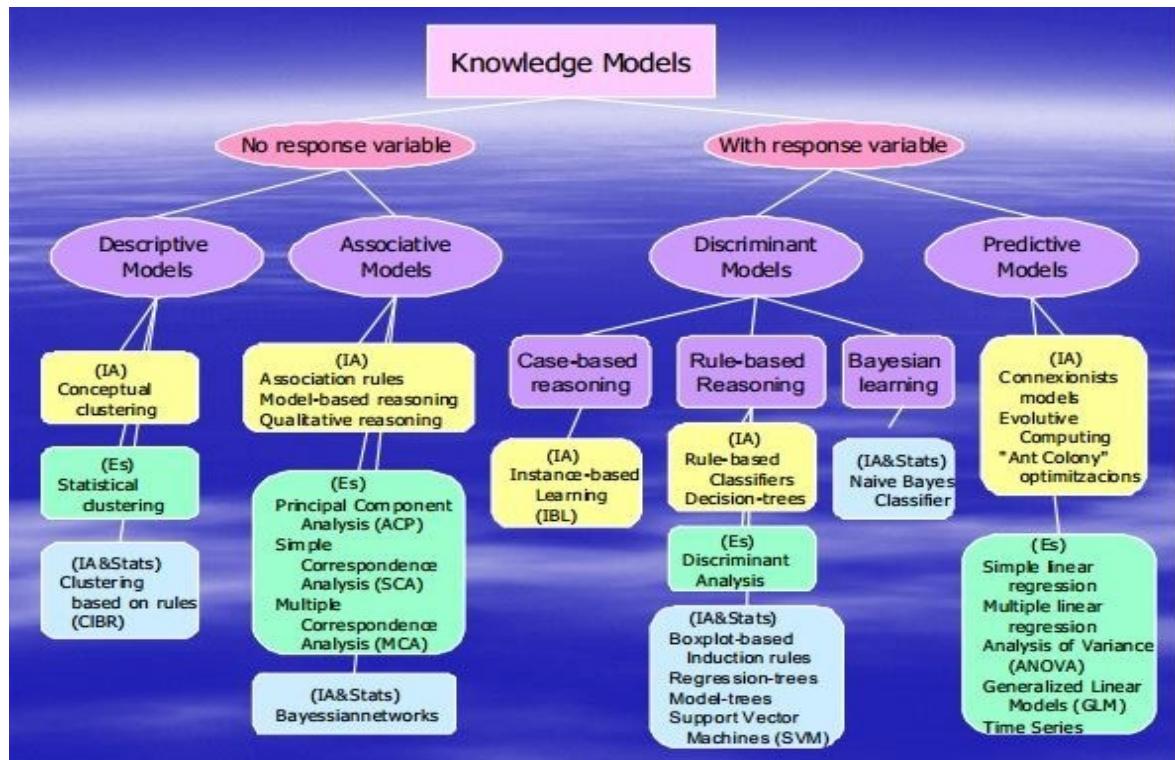


Figure 2.5: Select the Right Mining Technique .

In addition to the above there is another approach, shown in Figure 2.6 suggested by the very popular scikit (machine learning library) of python for data mining (Scikit, n.d.).

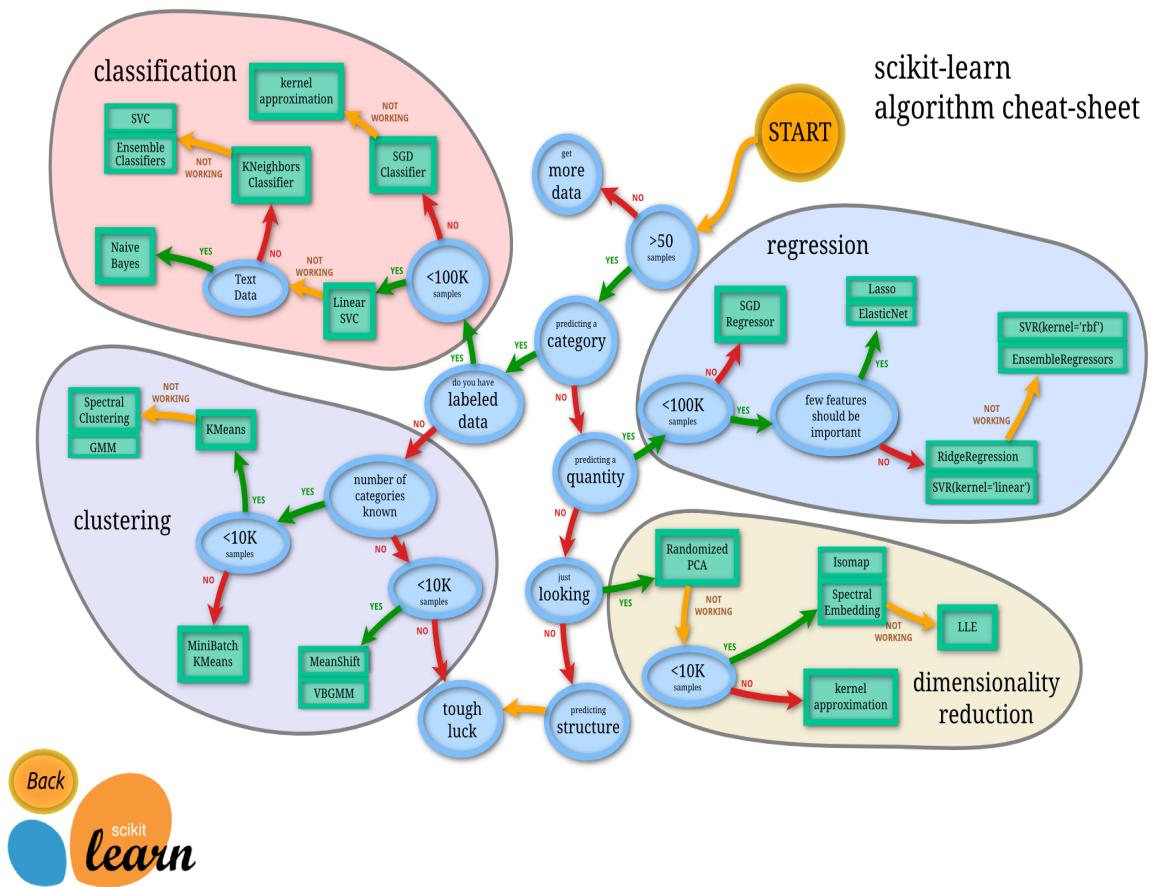


Figure 2.6: Another approach to select the Data Mining. Reprinted from Scikit.

2.3.4 Softwares, Libraries and Servers

Data mining techniques have been implemented into modules by a number of generous contributors. There are some very famous solutions that an academician can utilize for implementing predictive analytics. The algorithms and techniques of neural networks, clustering, classification and associations are available as solutions and API's. Following is a list in no particular order:

- **Softwares :** The following softwares are open source and available for data analytics by downloading to the desktop.
 - Weka : This is an open source software containing data mining algorithms. This can be used as a software or called by users own Java code.
 - Knime : An open source tool for data mining, comprises of many functions from data cleaning to pattern analysis.
 - Rapidminer : Also another popular tool for machine learning with plenty of algorithms for analysis.

- Libraries : These are available for use as toolbox and academic can program own solution.
 - Tensorflow
 - mlpack
 - H2O
 - Mlib
 - Scikit
- Servers : The following servers have built in modules that can be accessed via web applications and can be modeled to process real time analytics instead of one of processing as with above solutions
 - **DeepDetect** : is an open source deep learning server implemented in C++. It can be supported with back end machine learning applications with TensorFlow XGBoost and Caffe. Model assessment is built in the framework.
 - **Apache Prediction IO** : This a open source stack for academicians to deploy machine learning. The stack has an Event Server that can be used to query from a web application and respond in real time. The Event server co-ordinates with the Engine to respond to API inputs and respond with predicted outcomes Figure 2.7 (PredictionIO, n.d.). PredictionIO provides various templates for varied mining algorithms. Classification templates like Decision trees, Logistic Regression, NLP are available for use.

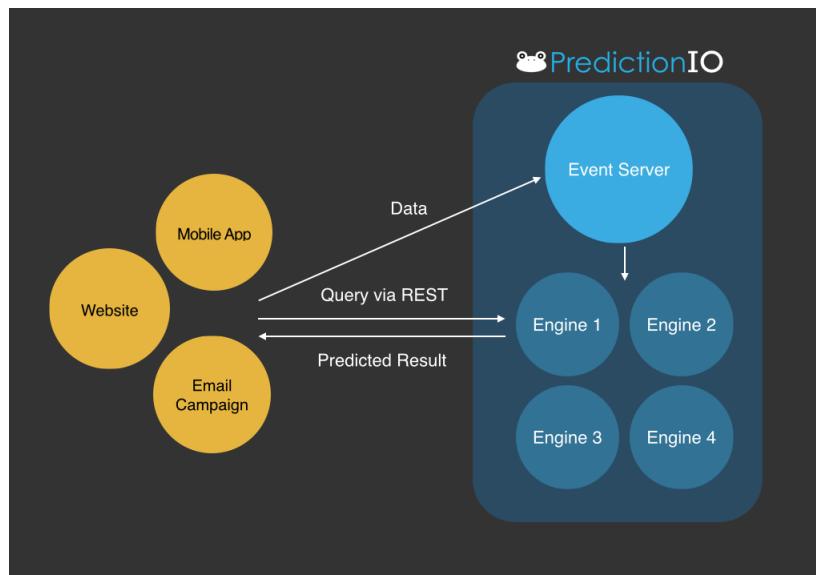


Figure 2.7: PredictionIO Engine interaction with Apps and Prediction Engine.

- **Shiny** : This is an R package and allows for easy to build web applications. It is made of two parts UI script and server script. In Figure it can be seen how Shiny can be implemented to exploit the data mining capabilities of R. Shown in Figure 2.8 how multiple users can access shiny R applications (Rstudio, n.d.).

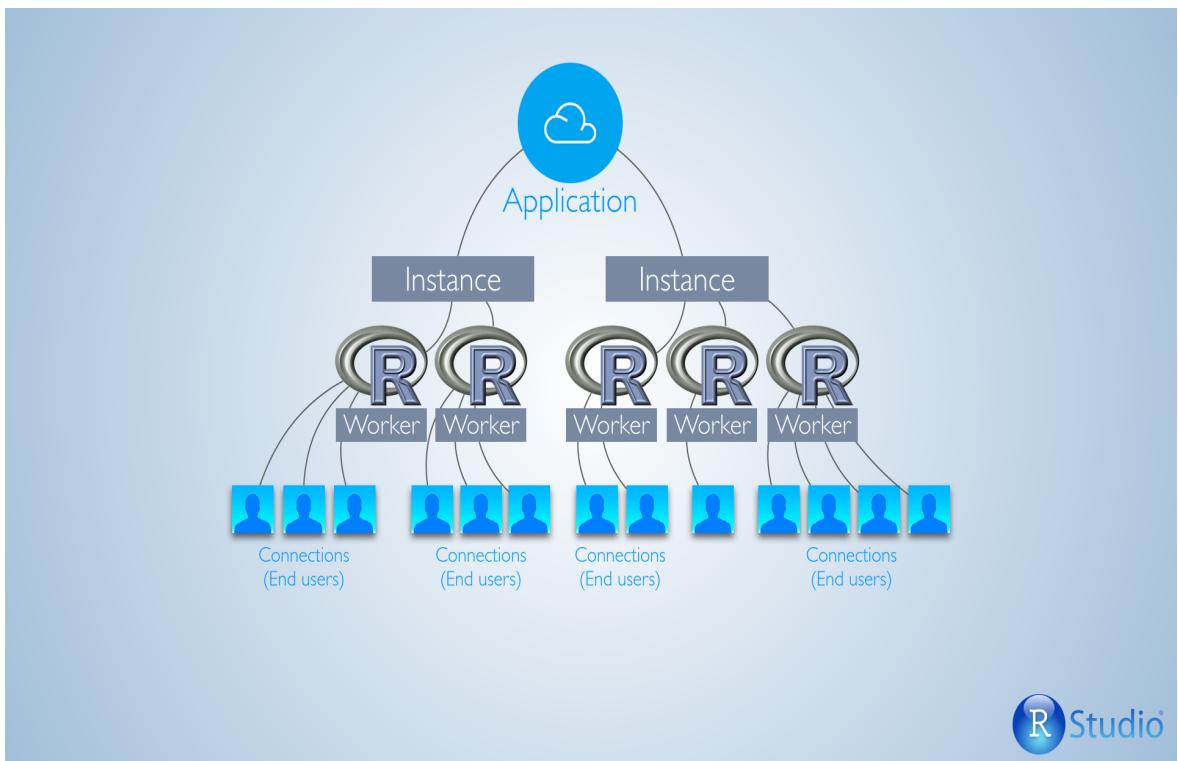


Figure 2.8: R Shiny architecture.

2.4 Model Evaluation Metrics

Model development is an important process, but evaluations of the model to ascertain its performance is as much an important procedure. The dataset is partitioned suitably and the testing set is not in the view of the model during training. There are however two methods of evaluations.

- Holdout technique
- K-fold Cross validation technique
- Leave one out CV
- Bootstrap method
- Sensitivity & Specificity

2.4.1 Holdout technique

This method is chosen for evaluation if the dataset is large enough. The data is segregated into three parts viz., Training, Validation and Test sets.

- Training dataset : It is some part of the dataset used for training the models. Predictive models are necessarily trained before actual prediction can be performed eg., Decision Trees, Random forest, Neural network need to be trained.
- Validation dataset : This is a subset of the data used to validate the output after model training. It helps to optimize the models performance. It is not mandatory to have validation sets for certain prediction models.
- Test dataset : Also a part of the whole dataset, it helps to

2.4.2 k-fold cross validation technique

This method of evaluations is chosen if the dataset is small and limited. The data is partitioned into k equal sized sets with an unbiased process. The model is built k times, with every $K-1$ data sets selected as training set, leaving out 1 set to be used as test set. A round robin process is followed to select the testing set in every iteration.

2.4.3 Sensitivity & Specificity

For calculating the performance of the model, a confusion matrix is plotted. The matrix is a cross table between predicted values and the actual values Figure 2.9. There are generally four types of

values that can be calculated from the matrix and those are as follows :

- TP - true positives : The predictor predicts “True” for actual true value of data.
- TN - true negatives : The predictor predicts “False” for actual false values of data.
- FP - false positives : The predictor predicts “False” for actual true value of data.
- FN - false negatives : The predictor predicts “True” for actual false values of data.

Sensitivity : the ratio of the count of the True Positives to the total count of events. This is also called the **Recall**.

$$Sensitivity(\text{or} Recall) = \frac{TP}{TP + FN}$$

Specificity : the ratio of the count of the True Negatives to the total count of non-events.

$$Specificity = \frac{TN}{FP + TN}$$

In addition to the above, True Positive value is called the **Precision**.

Form the values of *Precision* and *Recall* another statistical measurement called F-score can be derived.

$$F = 2 \times \frac{precision \times recall}{precision + recall}$$

		Prediction	
		Positive	Negative
Actual	Positive	TP	FN
	Negative	FP	TN

Figure 2.9: Confusion Matrix.

2.5 Review of Selected Research Papers

In the paper titled “Modeling & Simulation of a Predictive Customer Churn Model for Telecommunication Industry” the authors emulated a neuro fuzzy inference system to study the customer churn in the telecom industry (O et al., 2015). They modeled membership functions for the attributes of the dataset. Then they employed search algorithm for feature selection of the variables that indicate churn. Thereafter they model fuzzy equations to relate the dependent variables to the independent variables. This fuzzy system is trained to tune the Adaptive neuro fuzzy system based on the Sugeno FIS. The call detail records of 5000 subscribers was used to model this FIS. The dataset has 21 attributes but here they selected 9. Then the variables were modeled into three categories. For performance evaluation they calculated the Precision rate and the recall rate. After the testing it was found that accuracy was 95.8% , precision 80.86%, recall 92.7%.

A research study “A Hybrid Churn Prediction Model in Mobile Telecommunication Industry ” (Olle & Cai, 2014) presents a combination of LR and VP method. The academics used two algorithms of supervised learning viz., Logistic regression and Voted perceptron. They then combined the two into a Hybrid model for classification in WEKA. The obtained the data from an Asian telcom operator, records of around 2000 customers and 23 attributes.

From the results it was observed that hybrid model performed better than each of them individually.

In the study “A comparison of machine learning techniques for customer churn prediction” by (Vafeiadis et al., 2015) the researchers present a well meted out comparison between the normal model functions and their corresponding boosted models. The performance criteria was based on the F-score. They had used a series of simulations based on the Monte Carlo method. The models selected for analysis were Back-Propagation algorithm , Support Vector Machines, Decision Trees, Naive Bayes and Logistic Regression. The data was obtained from the publicly available churn dataset hosted at UCI Machine learning repository. The 100-fold cross validation technique was used to reduce bias. Ratio of training to testing set is about 2 : 3. A type of the most common boosting algorithm Adaboost, *Adaboost.M1* with DT and BPN as weak classifier was used.

The R programming was used for modeling the simulation experiment. Two steps were followed : Step 1 - tested classifiers run with data and performance of F-score measured. Step 2 - boosting algorithm was applied and performance F-score measured. 100 Monte carlo realizations were generated for cross validation of results. Monte carlo is synthesis of datasets that resemble the actual data. It was derieved from the results that two prediction models performed the best. 2 layer BPN with 15 hidden nodes and Decision tree classifier. An accuracy of 94% and F-measure around 77%. The SVM scored lower followed by Naive Bayes and Logit Regression at last. After application of the Boosting algo, SVM reported the best accuracy of 97% and Fmeasure over 84%.

2.6 Summary of Selected Research Studies

Here some of the past relevant literature in the domain of churn prediction and the results are discussed in Table 2.1.

Table 2.1: Previous literature reviewradioButtons("disp", "Display", choices = c(Head = "head", All = "all"), selected = "head").

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
1	Modeling & Simulation of a Predictive Customer Churn Model for Telecommunication Industry (O et al., 2015)	Adaptive neuro fuzzy inference system for prediction emulation of customer churn Neural network + fuzzy logic.	Data : 5000 subscribers CDR – call detail record with 21 variables. Partitioned into 5 sets each containing 1000 records. Method : Number of predictor variables taken is 9. Target variable is Churn with value Y or N. Membership function for each variable.	Found that 3 variables are very important. Total no of minute calls, no of customer service calls, no of repaired calls. Fuzzy churn model Precison 80.86% recall 92.7% and predicted accuracy 95.8%.	None suggested
2	A Hybrid Churn Prediction Model in Mobile Telecommunication Industry (Olle & Cai, 2014)	A model combined with VotedPerceptron and Logisti Regression is performance compared to the models of VP and LR as individual predictors.	Data : 2000 customers CDR from an Asian telecom company with 23 attributes. Method : A hybrid model of VP and LR was used. WEKA tool was used to model.	The hybrid model performs better than the models prediction accuracy seperately.	None suggested

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
3	A comparison of machine learning techniques for customer churn prediction (Vafeiadis et al., 2015)	The normal model functions were performed compared to their corresponding boosted models.	Data : publicly hosted churn dataset at UCI machine learning repository. Method : Machine learning techniques of Back-Propagation algorithm , Support Vector Machines, Decision Trees, Naive Bayes and Logistic Regression were used. The boosting algorithm Adaboost.M1 a type of Adaboost was used. R programming was used for modeling the system.	2 prediction models performed the best : 2-layer BPN with 15 hidden nodes and Decision tree classifier. SVM scored lower followed by Naive Bayes and Logit Regression at last. After application of the Boosting algo, SVM reported the best accuracy of 97% and Fmeasure over 84%.	None suggested
4	Turning telecommunications call details to churn prediction: a data mining approach (Wei & Chiu, 2002)	The company experiences a high monthly churn rate of 1.5 – 2Neural network requires a long time due to it's iterative nature. Highly skewed class distribution between churners and non-churners.	Data : Telecom company of Taiwan. Contractual and call details of subscribers Oct 2000 – Jan 2001. 9100000 records. Method : Multi classifier class combiner, Decision tree C4.5	Churn prediction is relatively high within 1 month duration. Multi classifier performs better than single classifier.	To include more variables from logs and complaints. Evaluation of empirical stats between customers from different geographic locations. Integration with data-warehouse for constantly learning behavior of customer. Research with other industry data from credit card to Internet service providers.

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
5	Applying Fuzzy Data Mining to Telecom Churn Management (Liao & Chueh, 2011).	To determine the most effective marketing strategies of customer retention, by analyzing the responses of customers.	Data : Taiwan telecom company, retention activity & response data for customer contract expiry between June and Junly 2008 Method : ID3 decision tree for classification.	Using fuzzy set the customer retention shows that marketing via telemarketing is more effective compared with Direct mailing. Also fuzzy marketing technique is better than direct mailing marketing for customers with higher bill amounts.	Fuzzy data mining techniques to analyze the past records of results of various marketing activities to establish a marketing mode.
6	Customer churn prediction using improved balanced random forests (Xie et al., 2009).	a novel learning method, called improved balanced random forests (IBRF), and demonstrate its application to churn prediction	Data : Chinese bank data. 1524 [762 train, 762 test]. Method : IBRF = Balanced random forest + weighted random forest. Introduce 2 interval variables ‘m – middle pt’ & ‘d – length of interval’. apply IBRF to a set of churn data in a bank as test the performance of our proposed method, we run several comparative experiments comparison of results from IBRF and other standard methods, namely artificial neural network (ANN), decision tree (DT), and CWC-SVM (Scholkopf, Platt, Shawe, Smola, & Williamson,	Accuracy rate follows this pattern $IBRF > CWC - SVM > ANN > DT$, Top-decile Lift varies as this $IBRF > CWV - SVM > DT > ANN$. IBRF offers great potential compared to traditional approaches due to its scalability, and faster training and running speeds.	Experimenting with some other weak learners in random forests. Improving effectiveness and generalization ability.

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
7	Churn prediction in subscription services: An application of support vector machines while comparing two parameter-selection techniques (Coussement & Poel, 2008)	Churn prediction using SVM. Benchmarked to Logit regression and random forest.	Data : Belgian newspaper publishing company. Training set 45000, Test set 45000 Method : Use of random forest software and SVM-toolbox. SVM compared to Logit regression & random forest. Grid search using 5-fold cross-validation	SVM trained on balanced distribution, outperforms logit regression when parameter selection applied. Random forest surpass SVM. Academincs and practitionerx don't need to rely on traditional Logit reg, SVM with parameter selection technique and random forest offer better alternative	No complete working meta-theory to choose kernel function and SVM parameters. Thus deriving a procedure to select proper kernel function and SVM parameter.
8	Customer churn prediction by Hybrid neural networks (Tsai & Lu, 2009)	Very few studies for hybrid data mining approach for prediction.	Data : CRM dataset from American telephone company, July 2001 to Jan 2002 51,306 subscribers. Method : 2 methods developed and compared for performance. M1 – SOM + ANN clustering + classification is used. M2 – ANN + ANN 2 classifiers are used. 5 fold cross validation, each set of the 5 are tested 5 times. Baseline is 20 ANN's	Baseline ANN models had prediction accuracy of 88% performance : $ANN + ANN > singleANN$ 3 * 3 SOM is best among 2 * 2 , 3 * 3, 4 * 4 and 5 * 5 clustering Performance of the hybrid models is : $ANN + ANN > SOM + ANN > ANN$	Need to explore dimensionality reduction or Feature selection of data preprocessing. Application of SVM or genetic algorithms. Explore other domains for churn prediction.

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
9	Predicting customer retention and profitability by using random forest and regression forest (Larivière & Poel, 2005)	The paper discusses more than one variable of retention and profit outcome.	Data : 100,000 Belgian finance company. Divided into 2 random parts, one for estimation other for evaluation. Method : Authors used random forest for regression to predict profitability, next purchase and defection decision. Benchmarked to linear regression model.	Random forest are better than logit and linear regression.	None suggested.
10	Churn prediction using comprehensible support vector machine: An analytical CRM application	The paper discusses more than one variable of retention and profit outcome.	Data : 100,000 Belgian finance company. Divided into 2 random parts, one for estimation other for evaluation. Method : Authors used random forest for regression to predict profitability, next purchase and defection decision. Benchmarked to linear regression model.	Random forest are better than logit and linear regression.	None suggested.
11	Churn prediction for high-value players in casual social games	Paper presents churn prediction of players of social games and the business impact of retaining high valued players.	Data : dataset of high value users of games - Diamond dash and Monster World, for 2 days. Method : The researchers trained and predicted neural networks, logistic regression, decision tree and support vector machine. Radial basis function for support vector machine was used with 10-fold cross validation. For business impact of churning the researchers designed A/B test.	Single neural network with tuned learning rate is better than other algorithms. A/B test reveals that sending free coins to high value customers does not affect churn rate.	None suggested.

Chapter 3

Methodology

In this chapter the methodology for implementing the ICPCR system is illustrated. Also the steps that would be followed are outlined.

3.1 Research Methodology

The following steps will be conducted also shown in Figure 3.1 :

Step 1: Data Preprocessing and Datawarehouse Development

- Data Collection
- Meta-data evaluation
- Data cleaning
- Datawarehouse design
- ETL process

Step 2: Development and Evaluation of the Prediction Models

- Select three churn prediction models
- Models to be trained and tested with the data
- Model Evaluation

Step 3: System Development & Evaluation

- Build the ICPCR system as a web application.
- Integration of Web app with OLAP and prediction model.
- Develop the Dashboards to display KPI's.
- Test the system.

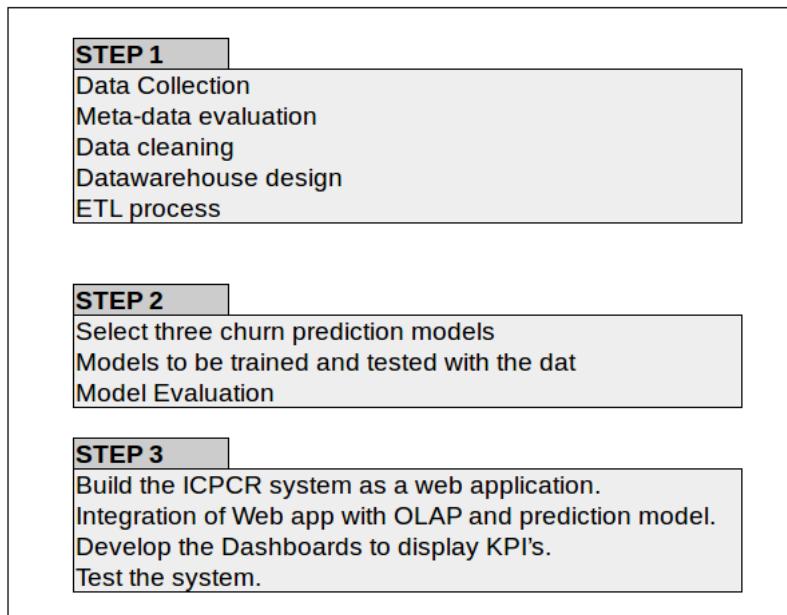


Figure 3.1: Research Methodology.

3.2 Data Preprocessing and Datawarehouse Development

3.2.1 Data preprocessing

Data will be collected from available open source sites. In this section a sequence of steps for data preparation are listed. In Figure 3.2 the process flow is shown.

1. Study of meta-data of the dataset. This study reveals the important attributes to be used for prediction.
2. Cleaning of un-useable data, either by replacing with suitable or by entirely removing it. Un-useable data is the one that may be invalid like null or special characters in numeric fields etc.
3. Extract the data and load into the database. This helps in querying the data faster with Structured Query Language.



Figure 3.2: Data preprocessing.

3.2.2 Datawarehouse development

Following steps will be followed for design of data warehouse:

The attributes generated from above step are summarized. This summary is used to design the OLAP cube. The OLAP will be used in generating reports and KPI's for the dashboard generation. The OLAP will be designed with the star schema. Figure 3.3 shows a typical implementation of the star schema (Tutorials Point, n.d.). A similar structure will be implemented for the study after the dimensions of the data are finalized.

Like for example the count of all the people between the age of 22 to 24 using prepaid service for the year 2013 could be one data whereas the count for 2014 would be another.

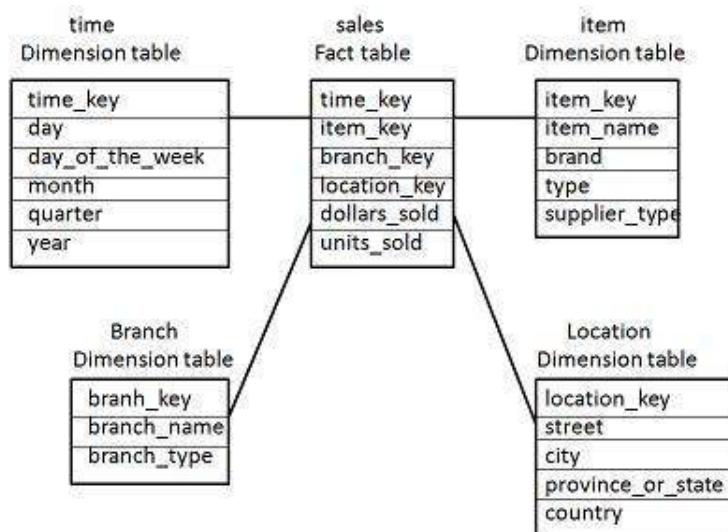


Figure 3.3: OLAP Star Schema.

After the Datawarehouse is designed, the tables have to be loaded with data. Thus the next step of ETL is done. Extract Transform and Load processing (ETL) This is a necessary step that would be required to properly extract data from the data file, transform the data types in order that they may be suitable for the database and finally loading to database.

3.3 Development and Evaluation of the Prediction Models

3.3.1 Model Design

In this section, the models are selected for churn prediction. Tentatively it is decided to select Decision tree, Support Sector Machine and ANN. The models will be trained with a training set and then the performance will be evaluated with the testing set. The proposal is to select either the machine learning library of MLlib under Apache or Scikit of Python or libraries under R. It will largely depend on the availability of the models in the libraries. In case a model is not available it will be sourced from another library. Also in addition it is proposed that a boosting algorithm like Adaboost would be used to measure change in prediction performance.

3.3.2 Model Evaluations

In order to judge the better performing model or rather the accuracy of predictability by the classification techniques, it is but necessary to perform an evaluation. The evaluations that are commonly performed by academicians are the k-Fold Cross Validation, Sensitivity & Specificity measurements (Larivière & Poel, 2005).

- K-Fold Cross Validation : It is proposed to per form this process to make the classification model more accurate. From previous literature it is learned that $k = 100$ is highly appropriate.
- Plotting of confusion matrix, as followed by other academicians and then deriving the Sensitivity, Specificity, Precision, Recall and F-score are the proposed evaluation techniques

3.4 To do

Adopt Rule based system because following :

reference: <http://www.dataskills.it/en/dalla-predictive-analytics-alla-prescrizione-di-una-soluzione>

As we know, predictions are not valid forever and therefore the clues the model provides should leave enough time to act or they would end up being of little use during the decision-making and implementation process. If the rules produced by the model are simple and easy to interpret, then the time necessary for the action to take place will be short, allowing decision makers to fit in into the time frame of the prediction's validity. We can now ask ourselves what are the algorithms that, in addition to generating a prediction, create rules that are immediately usable. These are three classes of algorithms:

Decision trees Fuzzy Rule-Based System Switching Neural Networks (Logic Learning Machine)

Other systems such as, for example, Neural Networks or Support Vector Machines, despite being effective from the point of view of accuracy, specificity, and sensitivity of the model, are black-box machines, i.e. they do not provide any insight into how they reach a certain prediction.

Decision trees are simple to use but do not have a particularly brilliant predictive performance yet they are widely used (perhaps because they are one of the few systems that can produce rules) and implemented in many data mining tools. Fuzzy systems are the best from the predictive viewpoint but are not very common. The last algorithm, Switching Neural Networks, is implemented only in one machine learning tool – Rulex (www.rulexinc.com) – and presents a very high predictive ability.

3.5 System Development & Evaluation

In this section the architecture of the ICPCR system is proposed. The application, shown in Figure 3.4, would be developed in a 3-tire format i.e, Database Layer, Application Layer, and Presentation Layer. The system is designed in two modes. One is the learning phase mode and the other is the Prediction phase mode. In the learning phase the system is fed data and the inference engine learns the trend. Testing and benchmarking along with weighting.

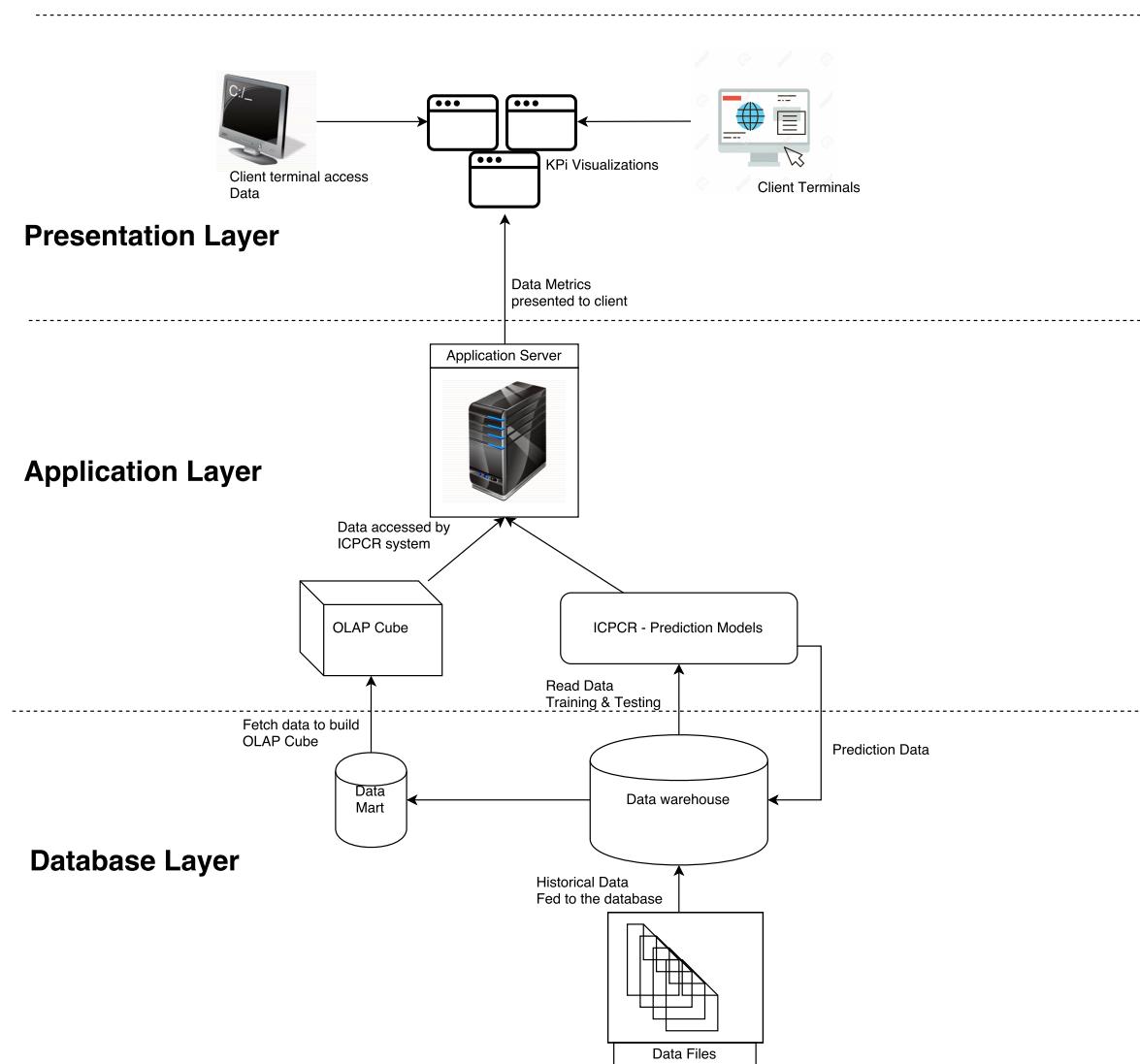


Figure 3.4: The Intelligent Churn Prediction Architecture.

3.5.1 Presentation Layer

In this thesis, the presentation layer is the section of the system which is accessible to the user or client. This is used to view the key values obtained from the OLAP and the mining results. There would be a display of metrics of the data.

1. It is proposed to deploy a suitable application to display a dashboard of KPI's.
2. The display of KPI's will be in graphs and charts format. The KPI's are taken from the OLAP cube.

3.5.2 Application Layer

This layer would be comprised of three parts.

1. Application server : This consists of the set of logic codes which will fetch the appropriate data for display in the front end. It may fetch the data directly from the tables or from the OLAP Cube, as is requested from the user.
2. Prediction model : This part is comprised of the predictive model to predict the outcome of data presented to it in the database. The model will go through a phase of training, testing, and prediction of churn value for new data. Also it is proposed that Prediction model be able to identify the variables which could be addressed for retaining the customer.
3. OLAP : This is the MOLAP implementation for building the Key metrics from the data. This part of the system would be responsible for the dashboard metrics display to the user.

3.5.3 Database Layer

This layer will be comprised of the data-warehouse tables. The OLAP calculation and the Model predictions will be updated whenever a set of new data is identified. The Olap cube feed tables will also be present here. A Star schema will be implemented for fetching of data for the various dimensions of the OLAP.

3.5.4 System Evaluation

The thesis proposes a system evaluation process to audit the performance. A set of test from latency in display and run will be calculated and improved before the process of deployment. This would ensure that system does not behave erratically under normal situations.

Chapter 4

Data preprocessing & Data warehouse

This chapter presents the progress in data preprocessing and data-warehouse development.

4.1 Data collection

The data was collected from the open source data available at SGI MLC++ website hosted at : <https://www.sgi.com/tech/mlc/db/>. The site has two files suitable for churn prediction. The data was donated to the public domain by Orange telecom.

- Data file “*churn.all*“ at location <https://www.sgi.com/tech/mlc/db/churn.all>.
- Meta-data file “*churn.names*“ for the data at location <https://www.sgi.com/tech/mlc/db/churn.names>

4.2 Meta-data evaluation

The meta-data is the set of field names which describe the purpose of the values contained in the field/column. In our call details records there are 21 fields. The description and data types are described below in Table 4.1.

Table 4.1: Meta data description.

Serial	Name of data field	Description	Units	Type
1	State	state's of USA	alphabetic string	discrete
2	Account Length	months of active usage	count in months	continuous
3	Area code	area code for phone	digits	continuous
4	Phone number	phone number	digits	discrete
5	Voice mail plan	subscribed to voice mail	binary value True/False	discrete
6	Number vmail messages	number of voice-mail messages	count	continuous
7	International plan	subscribed to international plan	binary values True/False	discrete
8	Total international minutes	total number of international calls	count in minutes	continuous
9	Total international calls	total number of international calls	count	continuous
10	Total international charge	total charge of international calls	dollars	continuous
11	Total day minutes	total minutes of day calls	count in minutes	continuous
12	Total day calls	total number of day calls	count	continuous
13	Total day charge	total charge of day calls	dollars	continuous
14	Total eve minutes	total minutes of evening calls	count in minutes	continuous
15	Total eve calls	total number of evening call	count	continuous
16	Total eve charge	total charge of evening calls	dollars	continuous
17	Total night minutes	total minutes of night call	count in minutes	continuous
18	Total night calls	total number of night calls	count	continuous
19	Total night charge	total charge of night calls	dollars	continuous
20	Number customer service calls	number of calls to customer service	count	continuous
21	Churn value	if customer churned or not	binary value True/False	discrete

4.3 Data cleaning

Data cleaning is a procedure that is very important to get rid of unwanted data values or checking for missing data values. The data is first loaded into MySQL database table **tab_churn**. A sample insert statement is shown below .

```
insert insert into tab_churn values
("OH", 84, 408, "375-9999", "yes", "no", 0, 299.4, 71,
 50.9, 61.9, 88, 5.26, 196.9, 89, 8.86, 6.6, 7, 1.78, 2,
 "False"),
("OK", 75, 415, "330-6626", "yes", "no", 0, 166.7, 113,
 28.34, 148.3, 122, 12.61, 186.9, 121, 8.41, 10.1, 3,
 2.73, 3, "False"),
("AL", 118, 510, "391-8027", "yes", "no", 0, 223.4, 98,
 37.98, 220.6, 101, 18.75, 203.9, 118, 9.18, 6.3, 6, 1.7,
 0, "False"),
("MA", 121, 510, "355-9993", "no", "yes", 24, 218.2, 88,
 37.09, 348.5, 108, 29.62, 212.6, 118, 9.57, 7.5, 7,
 2.03, 3, "False"),
("VT", 86, 415, "373-8058", "no", "yes", 34, 129.4, 102,
 22, 267.1, 104, 22.7, 154.8, 100, 6.97, 9.3, 16, 2.51,
 0, "False");
```

Now the data needs to be accessed into R environment for data analysis. This was done by using the MySQL connection parameters and pull the data into R and a data frame is created **churn_data**.

```
churn_data <- fetch(dbSendQuery(mydb,"select * from
tab_churn"),n=-1)
```

4.4 Data-warehouse design

The data-warehouse was developed in the MySQL environment. And the data was loaded from the tab_churn already loaded before. This helps in maintaining data integrity. To develop the Data-warehouse, fields of the churn dataset were modeled into dimensions and fact tables. The dimensional tables defined are listed as follows in tables 4.2 to 4.9:

Table 4.2: Dimensional table - t_day fields.

primary key	id_day
	total_day_minutes
	total_day_calls
	total_day_charges

Table 4.3: Dimensional table - t_eve fields.

primary key	id_eve
	total_eve_minutes
	total_eve_calls
	total_eve_charges

Table 4.4: Dimensional table - t_night fields.

primary key	id_night
	total_night_minutes
	total_night_calls
	total_night_charges

Table 4.5: Dimensional table - t_state fields.

primary key	id_state
	state
	state_full
	region

Table 4.6: Dimensional tables - *t_pid* fields.

primary key	id_pid
	area_code
	phone_number

Table 4.7: Dimensional table - *t_acct_len* fields.

primary key	id_acctlen
	account_length

Table 4.8: Dimensional table - *t_vm* fields.

primary key	id_vm
	voice_mail_plan
	number_vmail_messages

Table 4.9: Dimensional table - *t_intl* fields.

primary key	id_intl
	international_plan
	total_intl_minutes
	total_intl_calls
	total_intl_charge

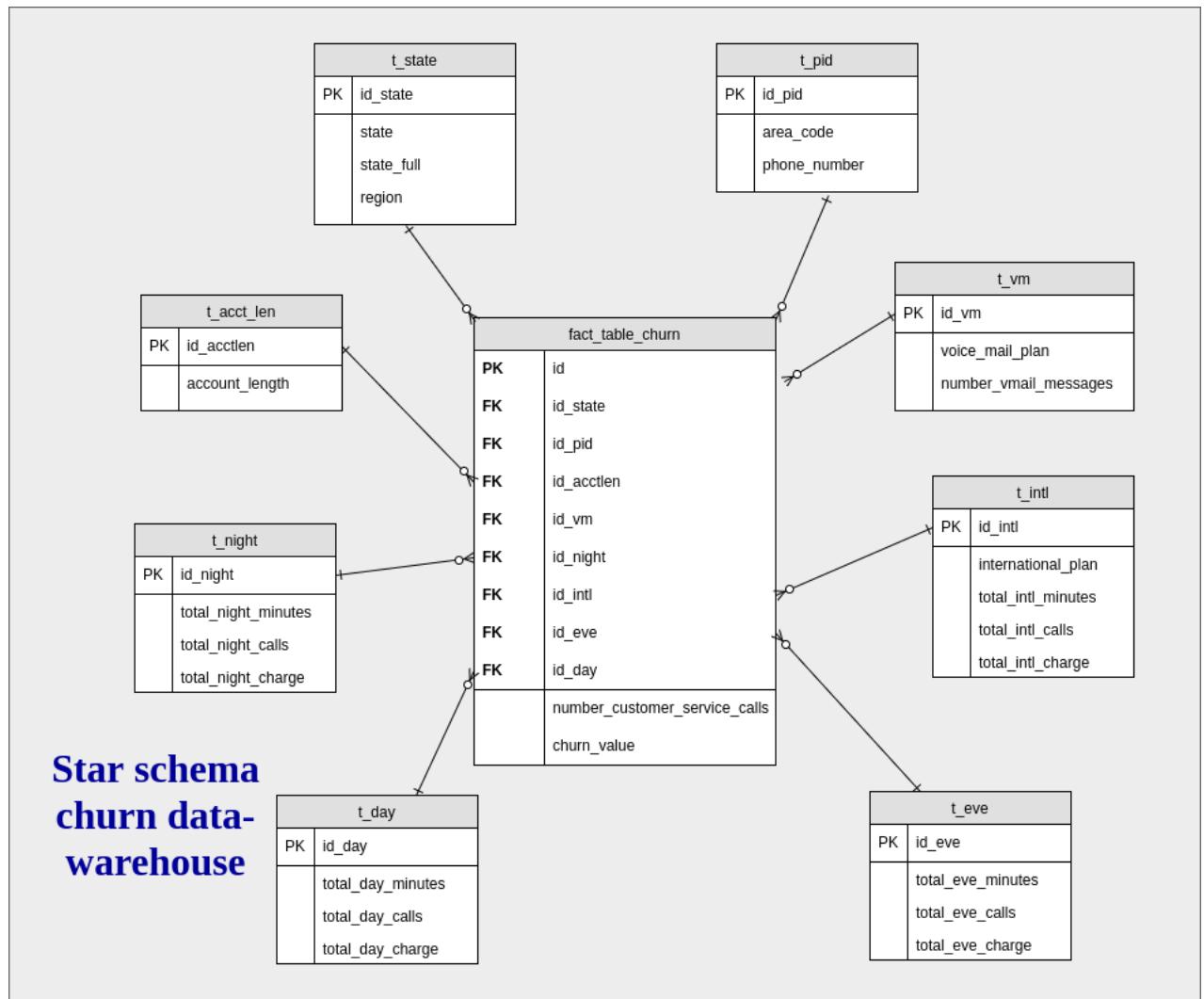
The fact table has foreign key attributes referencing the dimension tables defined above and contains numerical fields as follows :

Table 4.10: Fact table fields.

primary key	id
foreign key	id_night
foreign key	id_state
foreign key	id_day
foreign key	id_pid
foreign key	id_eve
foreign key	id_acctlen
foreign key	id_night
foreign key	id_vm
	churn_value
	number_customer_service_calls

The structure of the schema follows the star schema design recommended for data-warehouse. In Figure 4.1, the star schema of the system can be seen.

Figure 4.1: Star schema of churn data-warehouse.



4.5 ETL process

In the development of the system the extract transform and load of data was handled as a manual operation.

- **Extract** : in the process for extraction of data, the data is extracted from the csv files into the R session and MySQL database.
- **Transform** : in the process for transforming the data certain measures had to be taken before data is useful for machine processing. The data has to be checked for missing values and records containing NA's have to be removed.
- **Load** : The churn data is loaded into the data frames that are internal data structures of the R environment. This structure helps in processing the data as a table and thus reduces the dependency on the MySQL database. Data is loaded into data frames by using the `churn_data <- fetch(dbSendQuery(mydb, "select * from tab_churn"), n=-1)` command. In here the `churn` data frame is loaded by the `fetch` command. The data is pulled from the MySQL server.

Chapter 5

Model Development and Evaluation

This chapter presents the progress on development and evaluation metrics of prediction models. Also the data is to be transformed into two sets viz. training and testing dataset. The models were developed in the R environment and using the standard machine learning libraries available under open source license.

5.1 Prediction models selected

In the chapter the models trained thus far are :

- Decision tree ,
- Naive Bayes,
- Random Forest,
- Support Vector Machine.
- Neural Network
- C5.0, and
- C5.0 Boosted tree

All the models were trained with same training data set and tested with the same testing data set. Also the training was done on the same 18 features of the data. Below are the list of features :

1. state
2. account.length
3. international.plan
4. voice.mail.plan
5. number.vmail.messages
6. total.day.minutes
7. total.day.calls
8. total.day.charge

9. total.eve.minutes
10. total.eve.calls
11. total.eve.charge
12. total.night.minutes
13. total.night.calls
14. total.night.charge
15. total.intl.minutes
16. total.intl.calls
17. total.intl.charge
18. number.customer.service.calls

The features are converted into a formula to be used for training. The formula is as follows :

```
f <- as.formula(paste("churned ~ ", paste(n[-19], collapse = " + "), sep=""))  
> f  
churned ~ state + account.length + international.plan + voice.mail.plan +  
number.vmail.messages + total.day.minutes + total.day.calls +  
total.day.charge + total.eve.minutes + total.eve.calls +  
total.eve.charge + total.night.minutes + total.night.calls +  
total.night.charge + total.intl.minutes + total.intl.calls +  
total.intl.charge + number.customer.service.calls
```

5.2 Decision tree : rpart's classification tree

The library rpart of R was used for generating a tree. In the Figure 5.1, the generated decision tree is shown.

Table 5.1 shows confusion matrix of decision tree.

Table 5.1: Confusion matrix rpart decision tree.

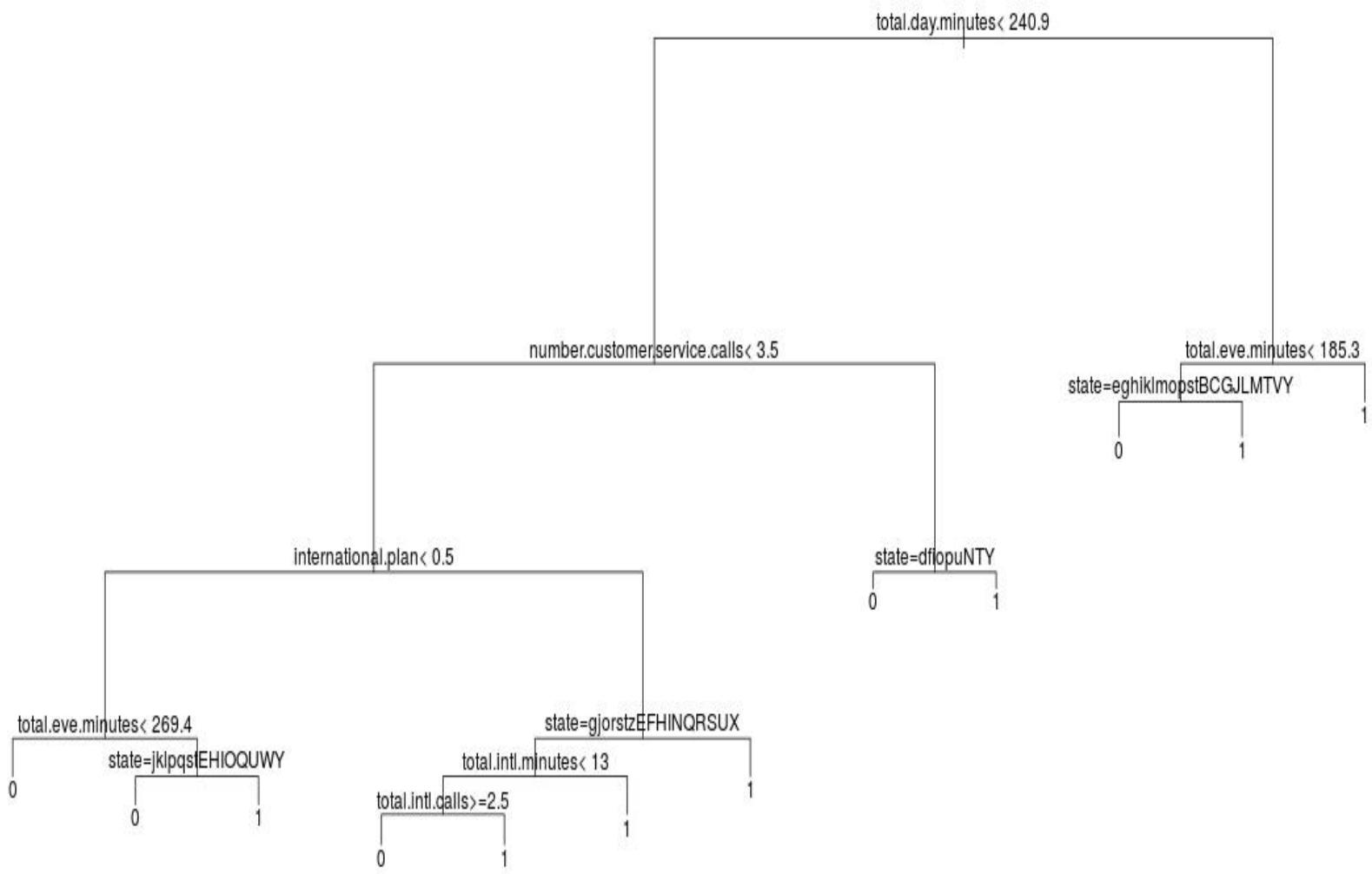
	Actual	
Prediction	False (0)	True (1)
False (0)	194	36
True (1)	29	165

Table 5.2 records the accuracy, sensitivity, and specificity.

Table 5.2: Accuracy, Sensitivity, Specificity - Rpart tree.

Accuracy	:	0.8467
Sensitivity	:	0.8700
Specificity	:	0.8209
Pos Pred Value	:	0.8435
Neg Pred Value	:	0.8505

Figure 5.1: Rpart decision tree.



5.3 Decision tree : ctree's conditional inference

The library ctree generates a conditional inference tree. By using this library we generated Conditional inference tree with 11 terminal nodes. The Figure 5.2 shows the decision tree model generated from ctree library.

The model generated is as follows:

```
1) international.plan <= 0; criterion = 1, statistic = 91.057
2) total.day.minutes <= 240.8; criterion = 1, statistic = 93.393
   3) number.customer.service.calls <= 3; criterion = 1, statistic = 132.291
      4) total.eve.minutes <= 269.3; criterion = 1, statistic = 44.098
         5)* weights = 448
      4) total.eve.minutes > 269.3
         6) total.day.minutes <= 207.1; criterion = 0.998, statistic = 32.855
            7)* weights = 25
         6) total.day.minutes > 207.1
            8)* weights = 32
   3) number.customer.service.calls > 3
      9) total.day.charge <= 33.44; criterion = 1, statistic = 59.779
         10) total.eve.minutes <= 233.2; criterion = 0.986, statistic = 46.025
            11)* weights = 90
         10) total.eve.minutes > 233.2
            12)* weights = 12
      9) total.day.charge > 33.44
         13)* weights = 14

2) total.day.minutes > 240.8
   14) voice.mail.plan <= 0; criterion = 1, statistic = 41.633
      15) total.eve.minutes <= 185.2; criterion = 1, statistic = 49.755
         16) total.day.charge <= 44.91; criterion = 0.993, statistic = 27.242
            17)* weights = 15
         16) total.day.charge > 44.91
            18)* weights = 36
      15) total.eve.minutes > 185.2
         19)* weights = 142
   14) voice.mail.plan > 0
      20)* weights = 13

1) international.plan > 0
   21)* weights = 163
```

Table 5.3 shows the statistics of confusion matrix.

Table 5.3: Confusion matrix ctree decision tree .

	Actual	
Prediction	False (0)	True (1)
False (0)	206	26
True (1)	17	175

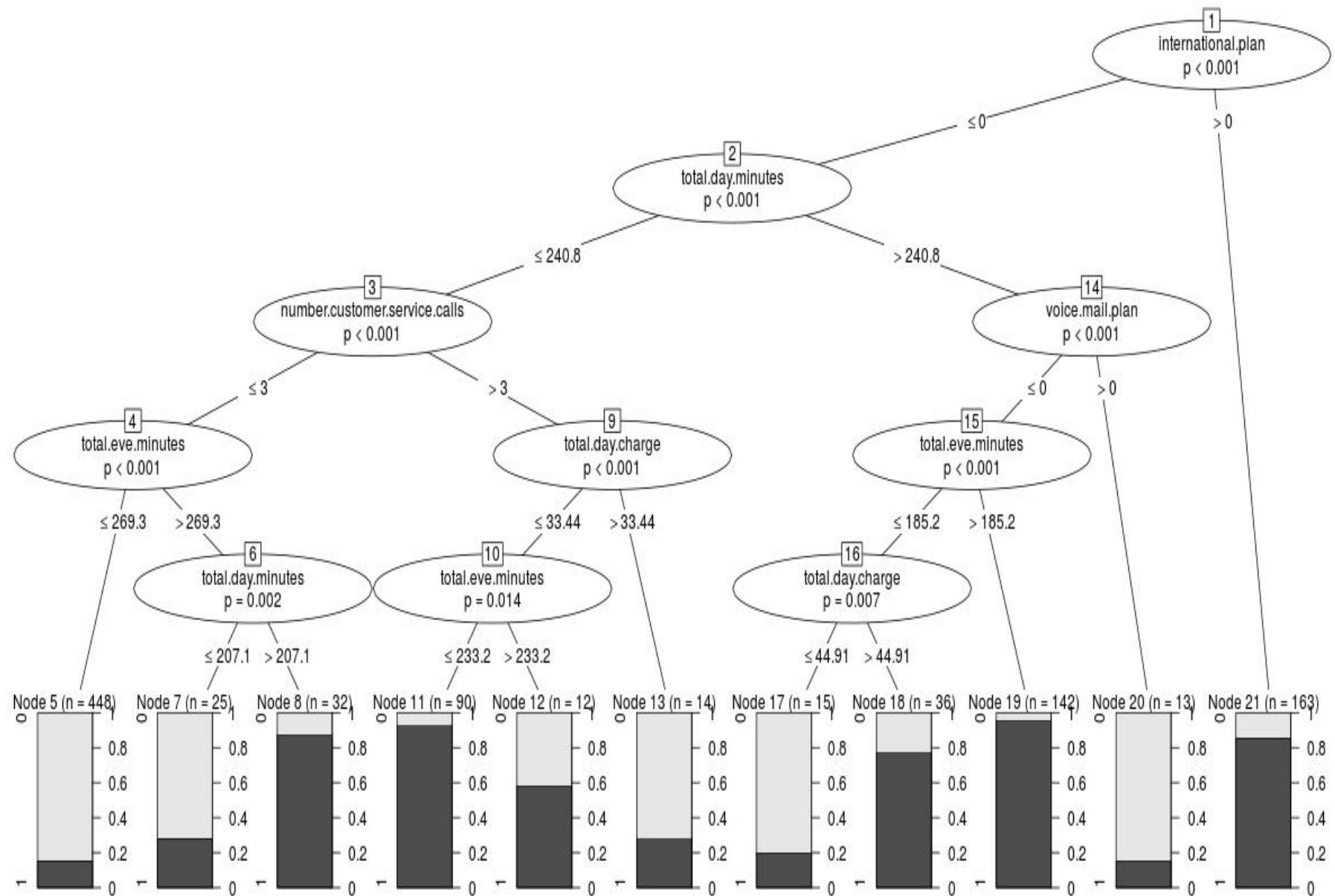
Table 5.4 records the accuracy, sensitivity, and specificity.

Table 5.4: Accuracy, Sensitivity, Specificity - Ctree tree.

Accuracy	:	0.8986
Sensitivity	:	0.9238
Specificity	:	0.8706
Pos Pred Value	:	0.8879
Neg Pred Value	:	0.9115

Figure 5.2: Ctree - conditional inference tree.

42



5.4 Random forest

This is a prediction model which can perform classification based on predictors. In model development random forest has not shown spectacular performance. Below are the

Table 5.5 shows the statistics of confusion matrix.

Table 5.5: Confusion matrix random forest.

	Actual	
Prediction	False (0)	True (1)
False (0)	215	79
True (1)	8	122

Table 5.6 records the accuracy, sensitivity, and specificity.

Table 5.6: Accuracy, Sensitivity, Specificity - Random forest.

Accuracy	:	0.7948
Sensitivity	:	0.9641
Specificity	:	0.6070
Pos Pred Value	:	0.7313
Neg Pred Value	:	0.9385

5.5 Decision tree C5.0

This is a C5.0 tree designed as a Rule based tree. Figure 5.3 shows the C5.0 decision tree model. The characteristics are as follows :

Rules :

Rule 1: (22/2, lift 1.8)
international.plan > 0
total.day.minutes <= 240.8
total.intl.minutes <= 13
total.intl.calls > 2
number.customer.service.calls <= 3
-> class 0 [0.875]

Rule 2: (448/68, lift 1.7)
international.plan <= 0
total.day.minutes <= 240.8
total.eve.charge <= 22.89
number.customer.service.calls <= 3
-> class 0 [0.847]

Rule 3: (62, lift 1.9)
international.plan > 0
total.intl.calls <= 2
-> class 1 [0.984]

Rule 4: (55, lift 1.9)
international.plan > 0
total.intl.minutes > 13
-> class 1 [0.982]

Rule 5: (239/39, lift 1.6)
total.day.minutes > 240.8
-> class 1 [0.834]

Rule 6: (152/25, lift 1.6)
number.customer.service.calls > 3
-> class 1 [0.831]

Rule 7: (96/27, lift 1.4)
total.eve.charge > 22.89
-> class 1 [0.714]

Table 5.7 shows confusion matrix of C5.0 tree.

Table 5.7: Confusion matrix c5.0 .

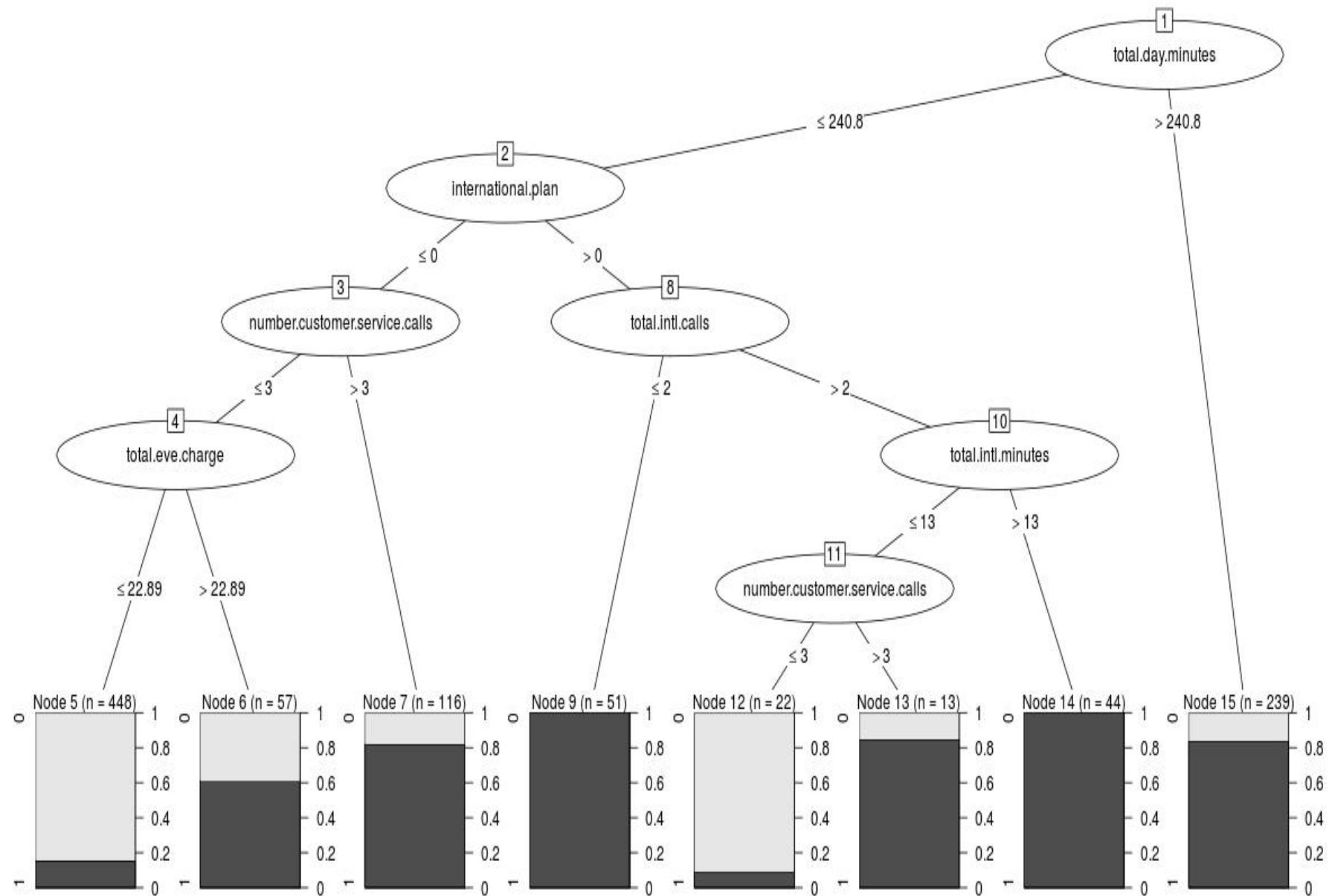
	Actual	
Prediction	False (0)	True (1)
False (0)	182	24
True (1)	41	177

Table 5.8 records the accuracy, sensitivity, and specificity.

Table 5.8: Accuracy, Sensitivity, Specificity - C5.0 tree.

Accuracy	:	0.8467
Sensitivity	:	0.8161
Specificity	:	0.8806
Pos Pred Value	:	0.8835
Neg Pred Value	:	0.8119

Figure 5.3: C 5.0 decision tree.



5.6 Decision tree C5.0 boosted

This is the boosted version of the C5.0 decision tree. The tree was grown by performing learning cycles 80 times and that resulted in pruning and error control. Figure 5.4 shows the C5.0 boosted decision tree model.

Table 5.9 shows the statistics of confusion matrix.

Table 5.9: Confusion matrix c5.0 Boosted.

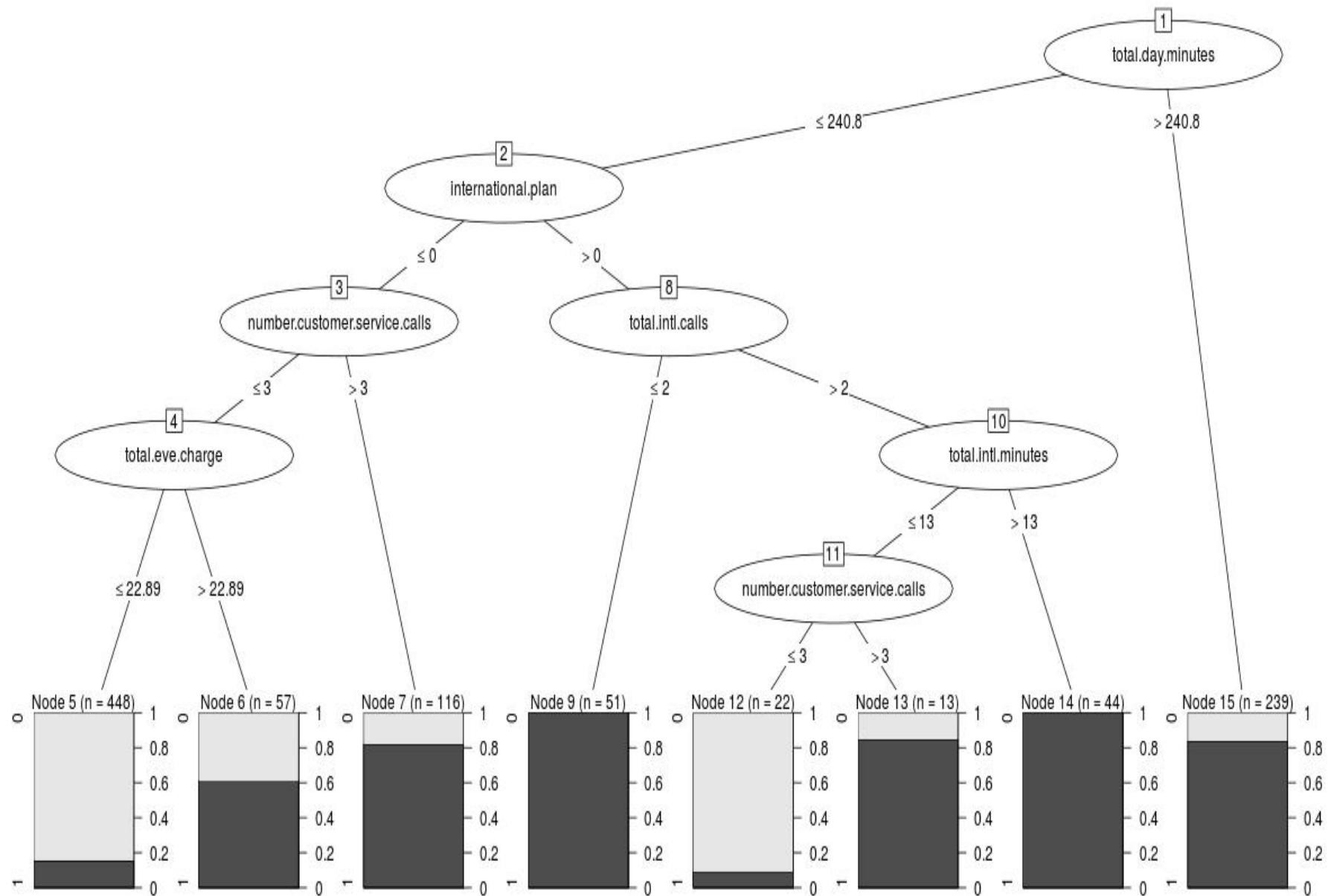
	Actual	
Prediction	False (0)	True (1)
False (0)	200	29
True (1)	23	172

Table 5.10 records the accuracy, sensitivity, and specificity.

Table 5.10: Accuracy, Sensitivity, Specificity - C5.0 Boosted tree.

Accuracy	:	0.8774
Sensitivity	:	0.8969
Specificity	:	0.8557
Pos Pred Value	:	0.8734
Neg Pred Value	:	0.8821

Figure 5.4: C 5.0 decision tree.



5.7 Naive bayes

In the second model training, a naive bayes classifier was trained and tested for performance. The following are parameters of the naive bayes developed.

In the following Table 5.11 are the statistics of confusion matrix.

Table 5.11: Confusion matrix naive bayes.

	Actual	
Prediction	False (0)	True (1)
False (0)	177	41
True (1)	46	160

Below Table 5.12 records the accuracy, sensitivity, and specificity.

Table 5.12: Accuracy, Sensitivity, Specificity - Naive Bayes.

Accuracy	:	0.7948
Sensitivity	:	0.7937
Specificity	:	0.7960
Pos Pred Value	:	0.8119
Neg Pred Value	:	0.7767

5.8 Support vector machine

In classification problems, support vector machines are very popular and hence the algorithm is quite suitable in this thesis to classify the churners among customers. Table 5.13 shows the statistics of confusion matrix.

Table 5.13: Confusion matrix support vector machines.

	Actual	
Prediction	False (0)	True (1)
False (0)	183	30
True (1)	40	171

Table 5.14 records the accuracy, sensitivity, and specificity.

Table 5.14: Accuracy, Sensitivity, Specificity - Naive Bayes.

Accuracy	:	0.8349
Sensitivity	:	0.8286
Specificity	:	0.8507
Pos Pred Value	:	0.8592
Neg Pred Value	:	0.8184

5.9 Neural network

Popularity of neural networks for solving classification problems grown by the day. The algorithm is implemented in this thesis to study the performance. A series of networks were designed to study the accuracy. The most relevant network was with two hidden layers with 4 nodes in first layer and 3 nodes in second layer, is shown in Figure 5.5. Table 5.15 shows the statistics of confusion matrix.

Table 5.15: Confusion matrix neural networks.

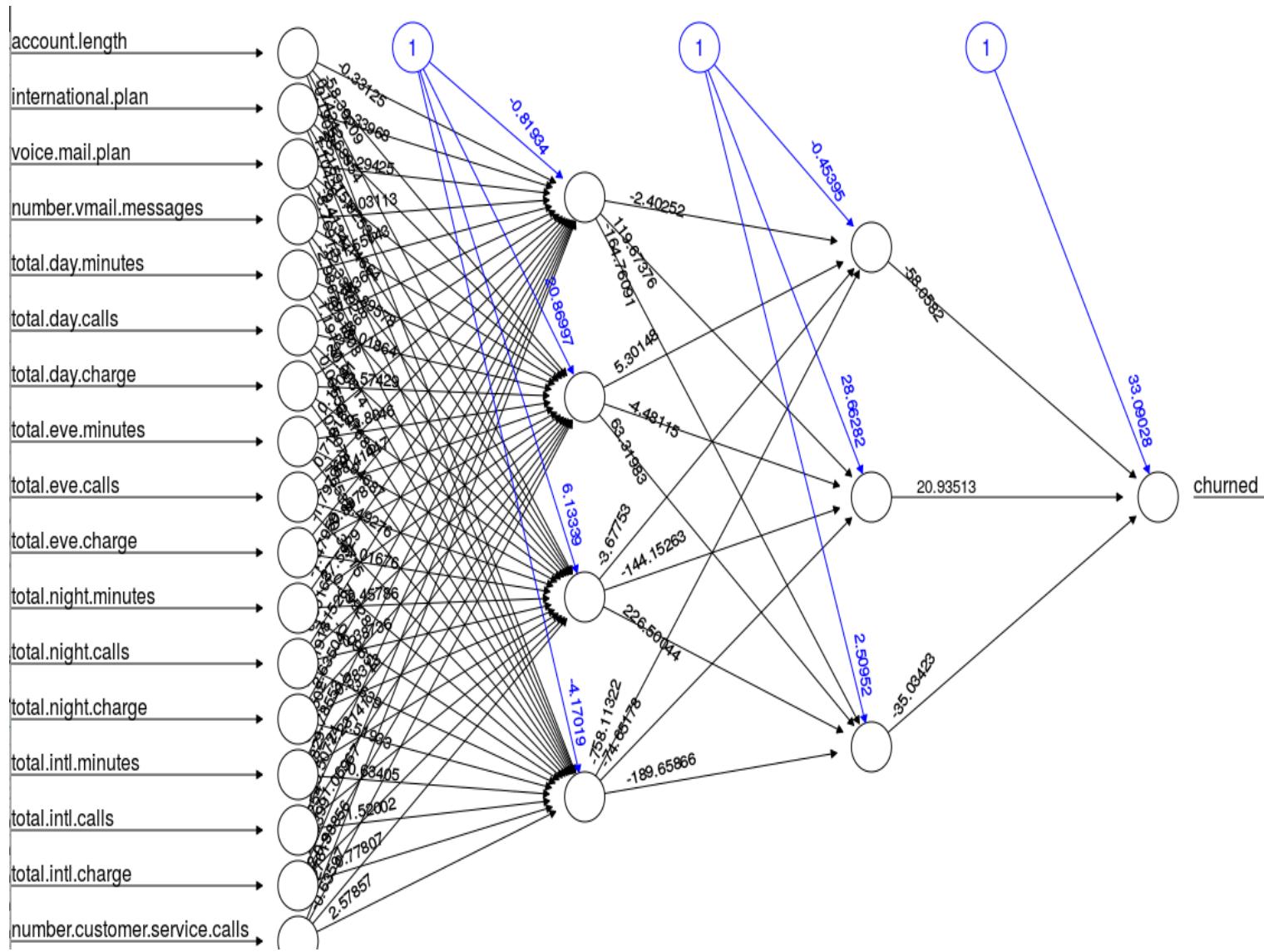
	Actual	
Prediction	False(0)	True (1)
False(0)	208	34
True (1)	15	167

Table 5.16 records the accuracy, sensitivity, and specificity.

Table 5.16: Accuracy, Sensitivity, Specificity - Naive Bayes.

Accuracy	:	0.8655
Sensitivity	:	0.7668
Specificity	:	0.8557
Pos Pred Value	:	0.8429
Neg Pred Value	:	0.8956

Figure 5.5: Neural network 2 hidden layers.



Chapter 6

Results and Discussion

This chapter covers the data analysis, and classification model comparisons and selection. Subsequently rules and decision table for customer retention are presented.

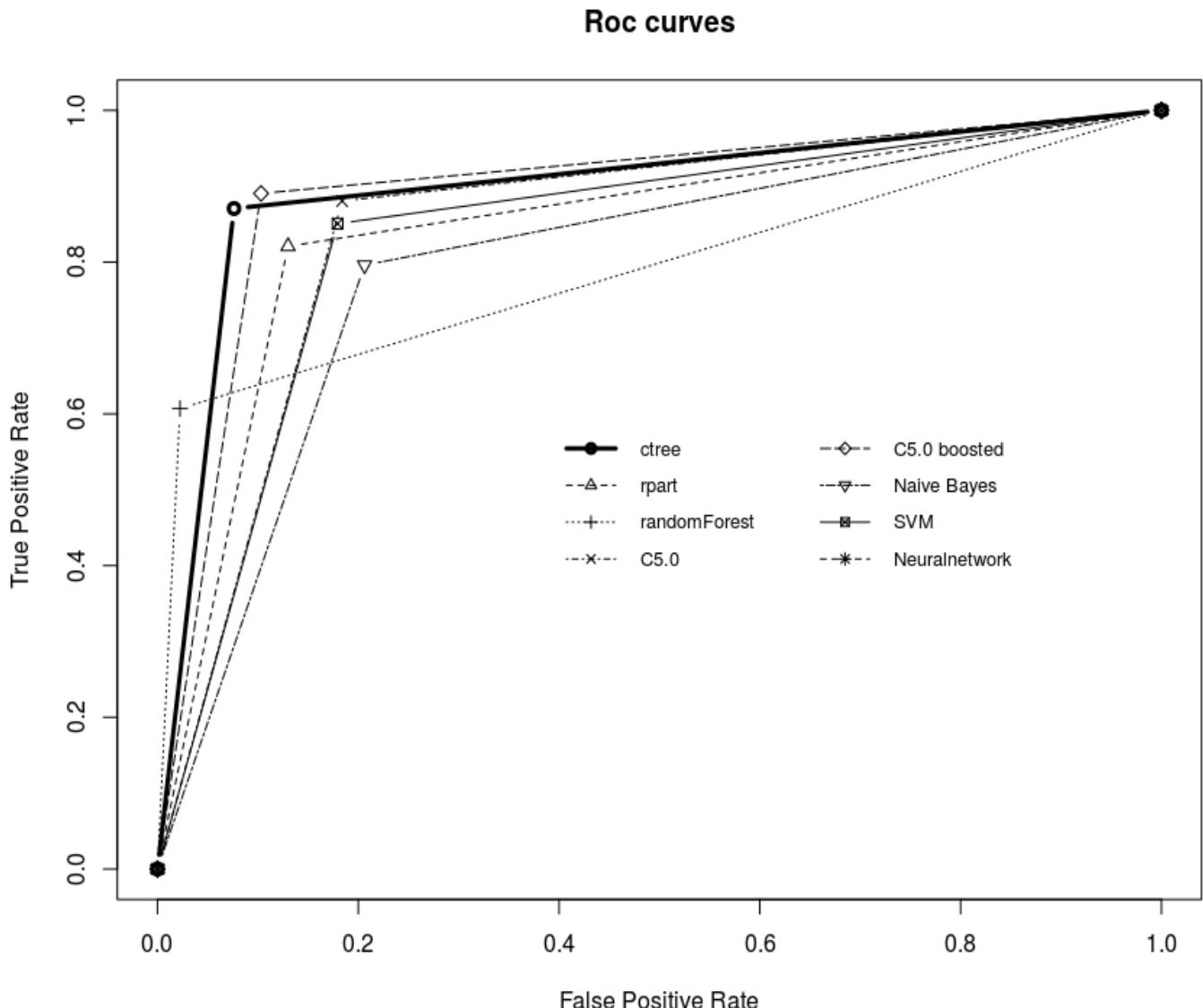
6.1 Model performance comparisons

In the previous chapter, eight models were developed trained and tested. The measurement statistics are shown in Table 6.1 below.

Table 6.1: Performance comparisons for prediction models.

SNo	Model	Accuracy	Sensitivity	Specificity	Pos Pred Rate	Neg Pred Rate
1	Tree - Rpart	0.8467	0.8700	0.8209	0.8435	0.8505
2	Tree - Ctree	0.8986	0.9238	0.8706	0.8879	0.9115
3	C5.0	0.8467	0.8161	0.8806	0.8835	0.8119
4	C5.0 Boosted	0.8938	0.8968	0.8905	0.90090	0.8861
5	Random Forest	0.8889	0.9686	0.6318	0.7448	0.9477
6	Naive Bayes	0.7948	0.7937	0.7960	0.8119	0.7767
7	Support Vector Machines	0.8089	0.8286	0.8557	0.8550	0.7678
8	Neural Network	0.8655	0.7668	0.8557	0.8429	0.8956

Figure 6.1: Model ROC curves.



From the metrics in Table 6.1 and from Figure 6.1, it is advisable to choose a model which has better sensitivity, specificity and a good accuracy rate. By comparing the values in Table 6.1, Conditional inference tree model from ctree library was chosen for prediction and the rules were generated from its decision tree for customer retention.

To test for customer retention a few 100 randomly selected records are picked from the data set, and predicted for churning. A decision table is generated from the rules of the Ctree decision tree. This table helps in identifying the features responsible for the particular classification label. Following are the rules which the system suggest to control churning. These rules are laid out below and tabular format in Table 6.2 for understanding the churn label classification.

```

[1] root
| [2] international.plan <= 0
| | [3] total.day.minutes <= 240.8
| | | [4] number.customer.service.calls <= 3
| | | | [5] total.eve.minutes <= 269.3: 0 (n = 448, err = 15.1786%)
| | | | [6] total.eve.minutes > 269.3
| | | | | [7] total.day.minutes <= 207.1: 0 (n = 25, err = 28.0000%)
| | | | | [8] total.day.minutes > 207.1: 1 (n = 32, err = 12.5000%)
| | | | [9] number.customer.service.calls > 3
| | | | | [10] total.day.charge <= 33.44
| | | | | | [11] total.eve.minutes <= 233.2: 1 (n = 90, err = 6.6667%)
| | | | | | [12] total.eve.minutes > 233.2: 1 (n = 12, err = 41.6667%)
| | | | | [13] total.day.charge > 33.44: 0 (n = 14, err = 28.5714%)
| | | [14] total.day.minutes > 240.8
| | | | [15] voice.mail.plan <= 0
| | | | | [16] total.eve.minutes <= 185.2
| | | | | | [17] total.day.charge <= 44.91: 0 (n = 15, err = 20.0000%)
| | | | | | [18] total.day.charge > 44.91: 1 (n = 36, err = 22.2222%)
| | | | | | [19] total.eve.minutes > 185.2: 1 (n = 142, err = 4.2254%)
| | | | | [20] voice.mail.plan > 0: 0 (n = 13, err = 15.3846%)
| [21] international.plan > 0: 1 (n = 163, err = 14.7239%)

```

Table 6.2: Decision table for classification.

No	Churning/Not Churning	Rule
1	Not Churning & Label 0	<ol style="list-style-type: none"> 1. <i>international.plan <= 0</i> 2. <i>total.day.minutes <= 240.8</i> 3. <i>number.customer.service.calls <= 3</i> 4. <i>total.eve.minutes <= 269.3</i>
2	Not Churning & Label 0	<ol style="list-style-type: none"> 1. <i>international.plan <= 0</i> 2. <i>total.day.minutes <= 240.8</i> 3. <i>number.customer.service.calls <= 3</i> 4. <i>total.eve.minutes > 269.3</i> 5. <i>total.day.minutes <= 207.1</i>

No	Churning/Not Churning	Rule
3	Churning & Label 1	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes \leq 240.8$ 3. $number.customer.service.calls \leq 3$ 4. $total.eve.minutes > 269.3$ 5. $total.day.minutes > 207.1$
4	Churning & Label 1	<ol style="list-style-type: none"> 1. $international.plan \leq 0 \ \&$ 2. $total.day.minutes \leq 240.8 \ \&$ 3. $number.customer.service.calls > 3 \ \&$ 4. $total.day.charge \leq 33.44 \ \&$ 5. $total.eve.minutes \leq 233.2$
5	Churning & Label 1	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes \leq 240.8$ 3. $number.customer.service.calls > 3$ 4. $total.day.charge \leq 33.44$ 5. $total.eve.minutes > 233.2$
6	Not Churning & Label 0	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes \leq 240.8$ 3. $number.customer.service.calls > 3$ 4. $total.day.charge > 33.44$

No	Churning/Not Churning	Rule
7	Not Churning & Label 0	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes > 240.8$ 3. $voice.mail.plan \leq 0$ 4. $total.eve.minutes \leq 185.2$ 5. $total.day.charge \leq 44.91$
8	Churning & Label 1	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes > 240.8$ 3. $voice.mail.plan \leq 0$ 4. $total.eve.minutes \leq 185.2$ 5. $total.day.charge > 44.91$
9	Churning & Label 1	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes > 240.8$ 3. $voice.mail.plan \leq 0$ 4. $total.eve.minutes > 185.2$
10	Not Churning & Label 0	<ol style="list-style-type: none"> 1. $international.plan \leq 0$ 2. $total.day.minutes > 240.8$ 3. $voice.mail.plan > 0$
11	Churning & Label 1	<ol style="list-style-type: none"> 1. $international.plan > 0$

6.2 Discussion

The decision table in Table 6.2 has 11 rules, with 6 rules deciding the churning classification. The decision table is an integral part of the ICPCR system, since it is used to explain the classification labels. Also it helps to identify the features which are responsible for customer churn. In order to retain the customer, upselling of benefits for the lowest consumed feature is suggested to the representative.

The ctree decision tree shows the most acceptable performance among the many models chosen to model prediction. The performance of ctree stands at 89%. The use of decision tree also helps in identifying the most important features and the features which have no significance on churn prediction. Like the feature "state", which feels absolutely perfect to manual inspection but has no significance on the prediction and hence is not considered in the decision tree.

Decision trees are also beneficial because it helps in determining the significance of features on the class variable. In this study the decision tree generated rules which were used to generate rule induction decision table to suggest retention strategies for the churning customers.

Chapter 7

System Development and Evaluation

This chapter presents the web service designed for ICPCR system. The objective of the study is to enable data visualizations and prediction of churning customers via a web service, hence this chapter describes the web site and its components along with screenshots.

7.1 ICPCR web service

The ICPCR web data dashboard with prediction and retention facility are designed and implemented with Shiny server's code base. The data crunching, visualizations and prediction modeling are developed in R language.

Tools used in development :

- R scripting
- RStudio
- Shiny server codebase
- R libraries for machine learning algorithms for preparing predictive model
- R Olap for presenting Pivot Table
- MySQL database management engine for data storage
- Plotly javascript library for presenting graphs
- Data table for creating rule engine to suggest retention techniques

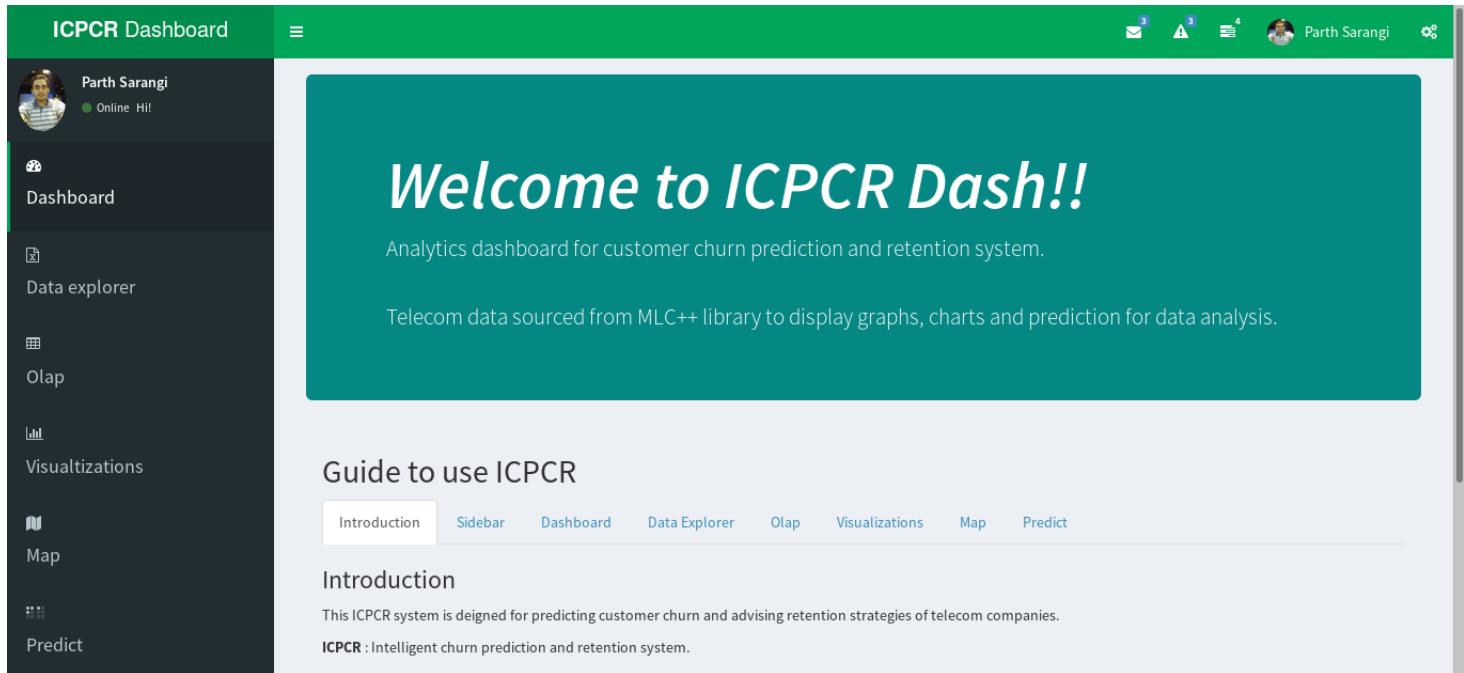
There are 6 tabs in the ICPCR web site :

- Dashboard
- Data explorer
- OLAP
- Graph visualizations
- Map visualizations
- Prediction functionality

7.2 Dashboard

The ICPCR dashboard tab is active by default when the application loads. On the dashboard the welcome banner is displayed. In the second part of the page a guide is available for guiding the user about the functionalities implemented. A screenshot of the dashboard is shown in Figure 7.1.

Figure 7.1: Dashboard - Data explorer.



7.3 Data explorer

In Figure 7.2 and figure 7.3 some analysis of the data fields are printed for the user to understand basic data structure. The print contains the mean value, the 1st quartile, median, 3rd quartile, minimum, maximum and class values.

Figure 7.2: Dashboard - Key performance index.

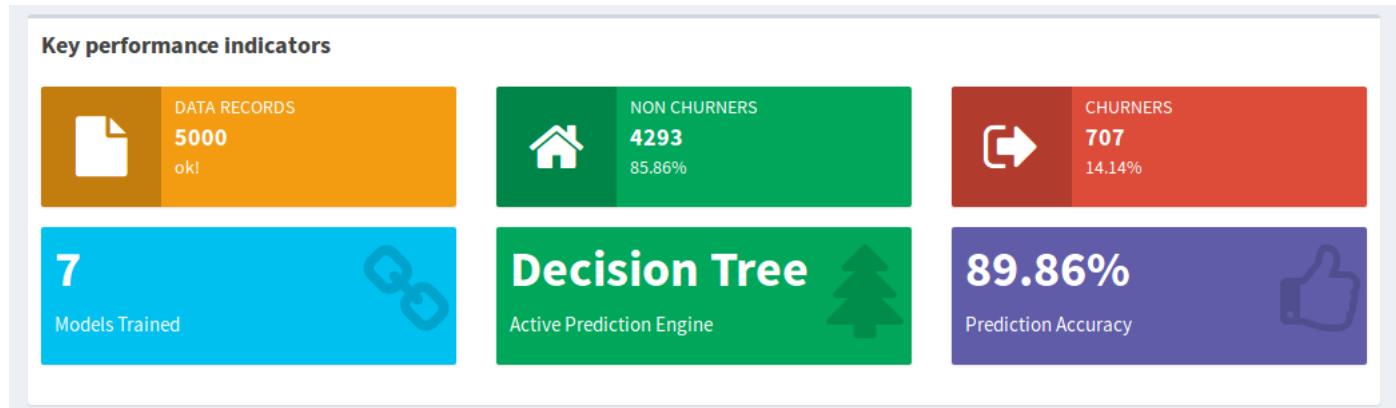
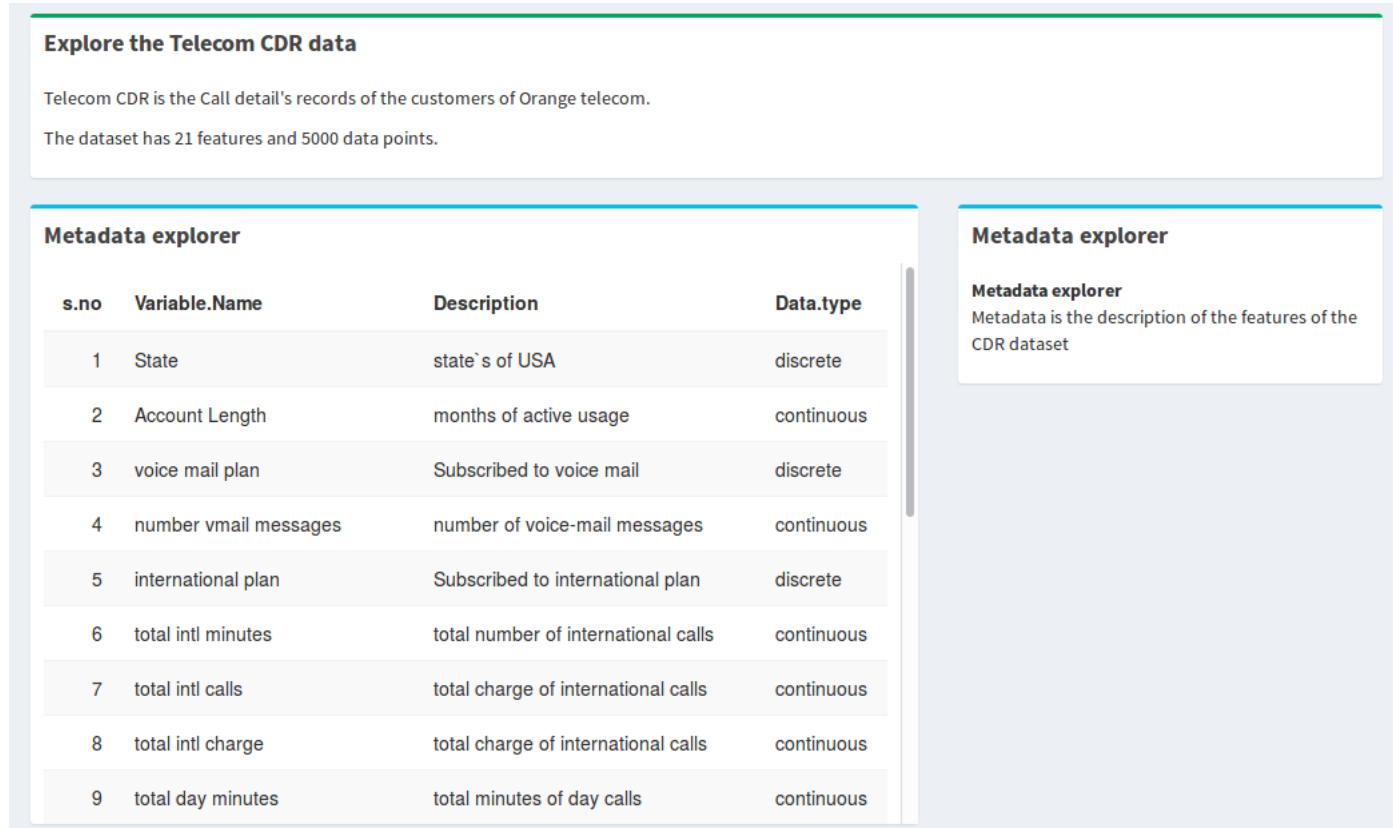


Figure 7.3: Dashboard - field analysis.



7.4 OLAP

In the ICPCR the OLAP functionality allows the customer relationship managers to perform data drill-down and roll-up activities. The OLAP functionality is made available via Pivot table in R. Pivot table has functionality to drag and drop features onto the axes. A screenshot of the ICPCR OLAP pivot table is shown in Figure 7.4.

Figure 7.4: OLAP - functionality helps to investigate data via drill-down and roll-up.

Investigate Data											
Row Heatmap	Count	churn_value	international_plan	voice_mail_plan							
account_length	state	churn_value	False		True					Totals	
area_code	international_plan	no		yes		no		yes			
phone_number	voice_mail_plan	no	yes	no	yes	no	yes	no	yes		
number_vmail_messages	state	no	yes	no	yes	no	yes	no	yes		
total_day_minutes	AK	44	19	2	2	3		1	1	72	
total_day_calls	AL	77	29	5		9		3	1	124	
total_day_charge	AR	43	29	3	1	9	1	5	1	92	
total_eve_minutes	AZ	48	28	2	2	7		2		89	
total_eve_calls	CA	28	10	2		9	2	3		52	
total_eve_charge	CO	54	25	5	1	7	1	3		96	
total_night_minutes	CT	52	26	4	1	10	1	4	1	99	
total_night_calls	DC	44	27	6	2	6	1	2		88	
total_night_charge	DE	52	21	4	2	9		4	2	94	
total_intl_minutes	FL	49	23	5	1	7	3	2		90	
total_intl_calls	GA	49	21	1	2	5	1	3	1	83	
total_intl_charge	HI	59	17	2	3	4		1		86	
	IA	43	17		1	6		1	1	69	
	ID	67	29	4	5	11	1	1	1	119	
	IL	53	18	8	2	1		4	2	88	
	IN	57	24	1	2	13		1		98	

7.5 Visualizations

In this tab, various graphs are plotted to help user visualize the data features. Figure 7.5 displays the bar plot region wise for the churners to non-churners.

Figure 7.5: Dashboard - Bar plot of Total, Churner and non-Churner customer counts region wise.

Bar plot for region

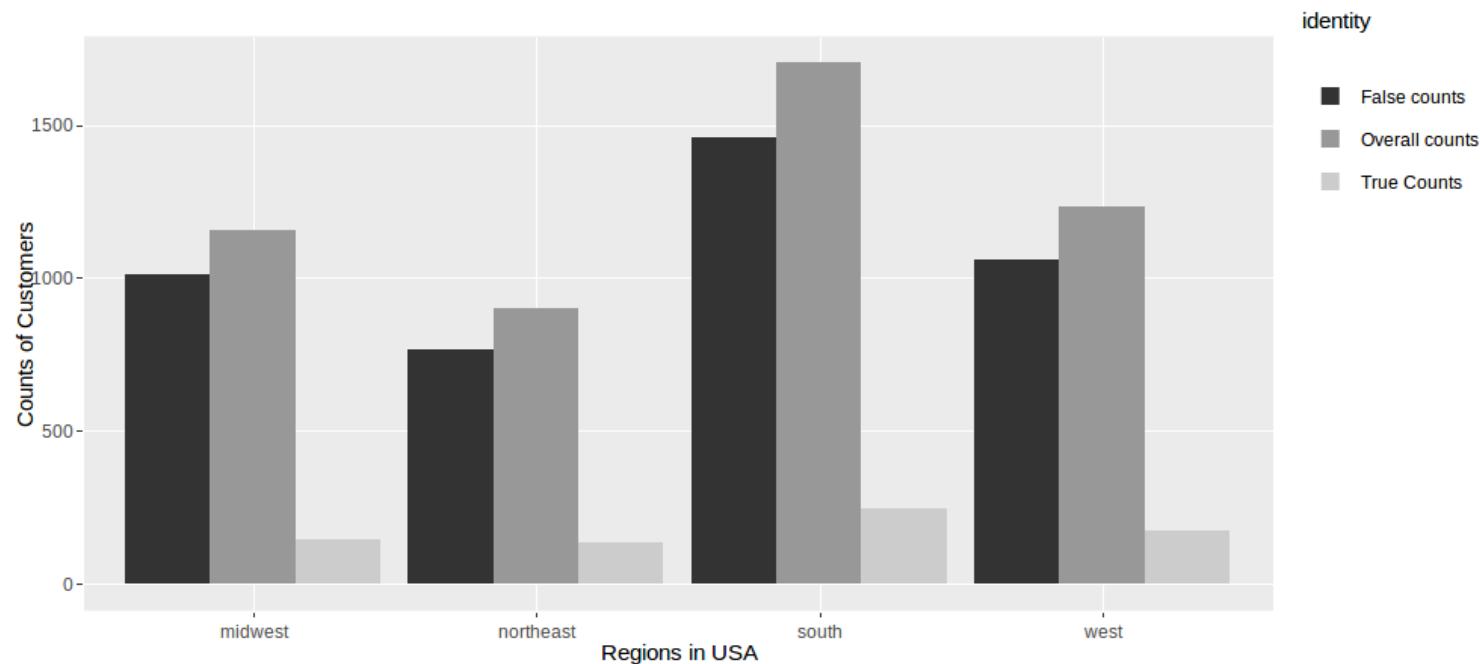


Figure 7.6: Dashboard - Bar plot of Churner customers to Non-churners State wise .

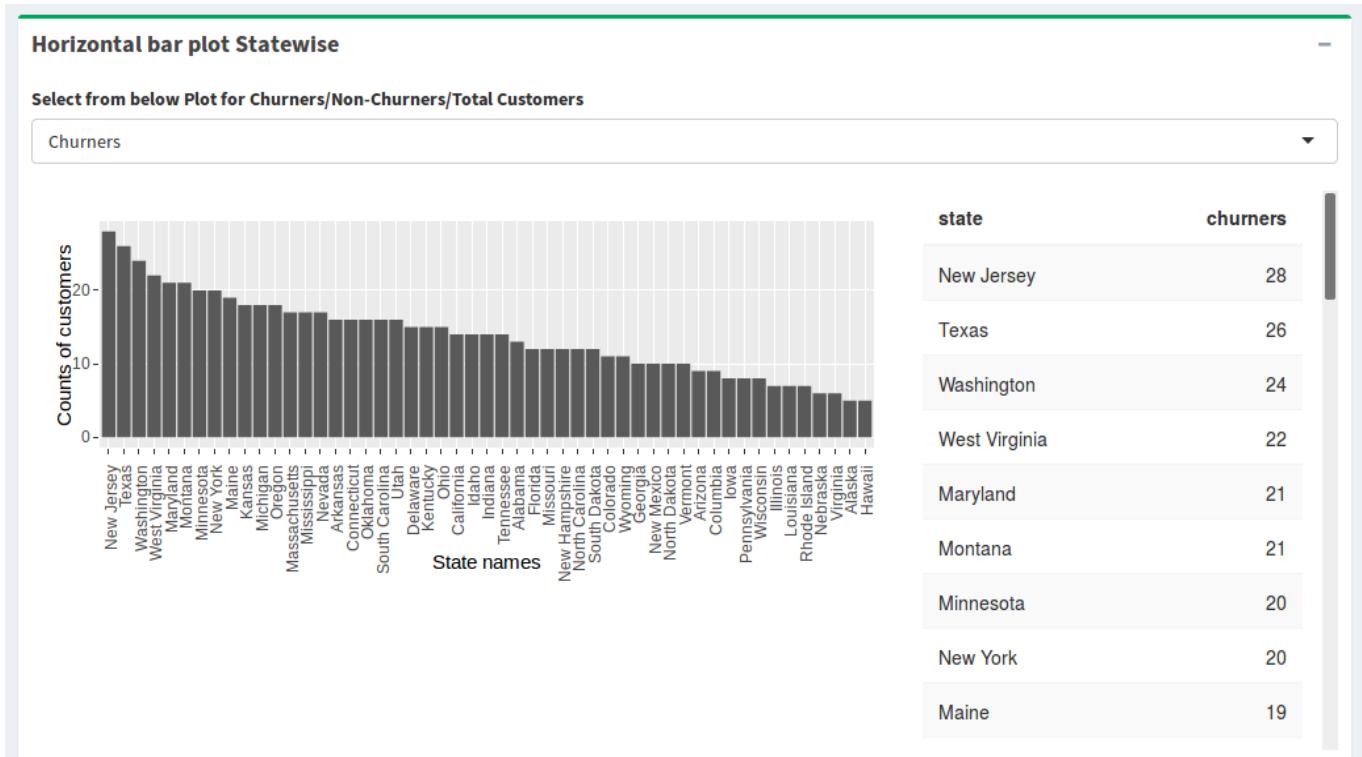


Figure 7.6 plots a bar chart of the state wise churners to non-churners. It can be understood that some states have higher percentage of churner customers.

7.6 Map

In Figures 7.7 and 7.8 of the dashboard map visualizations show state wise distribution of telecom population and user will be able to customize to select churning or retaining populations.

Figure 7.7: Dashboard - Map to visually show the density of customers state wise.

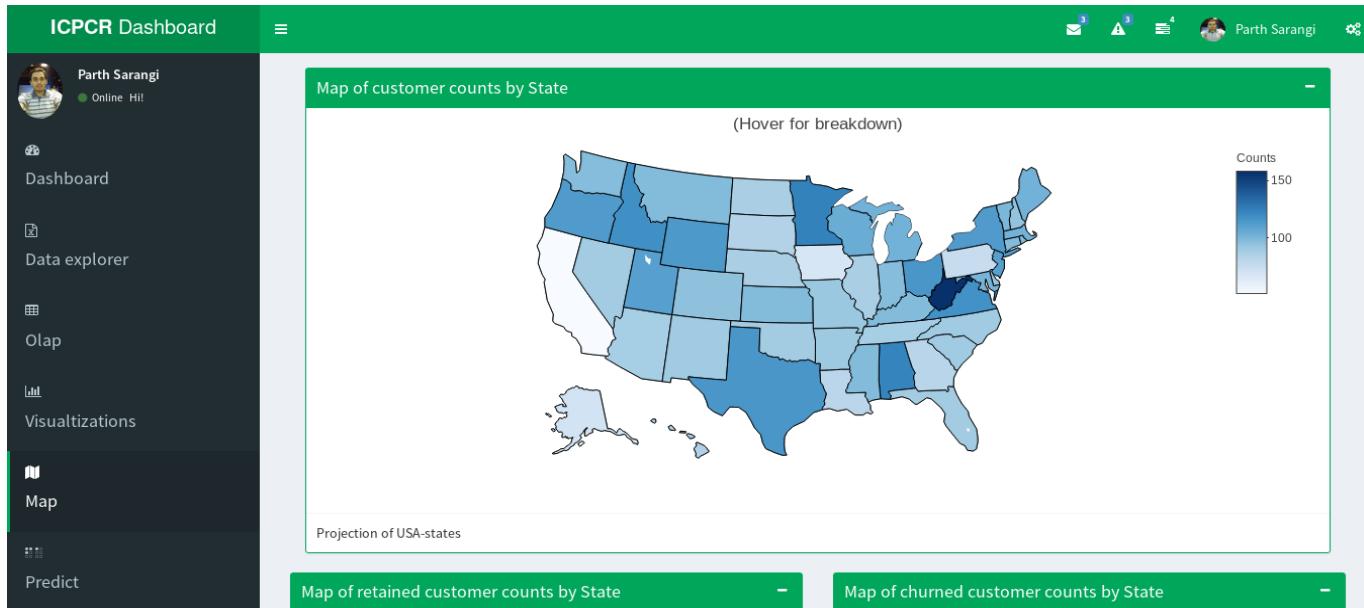
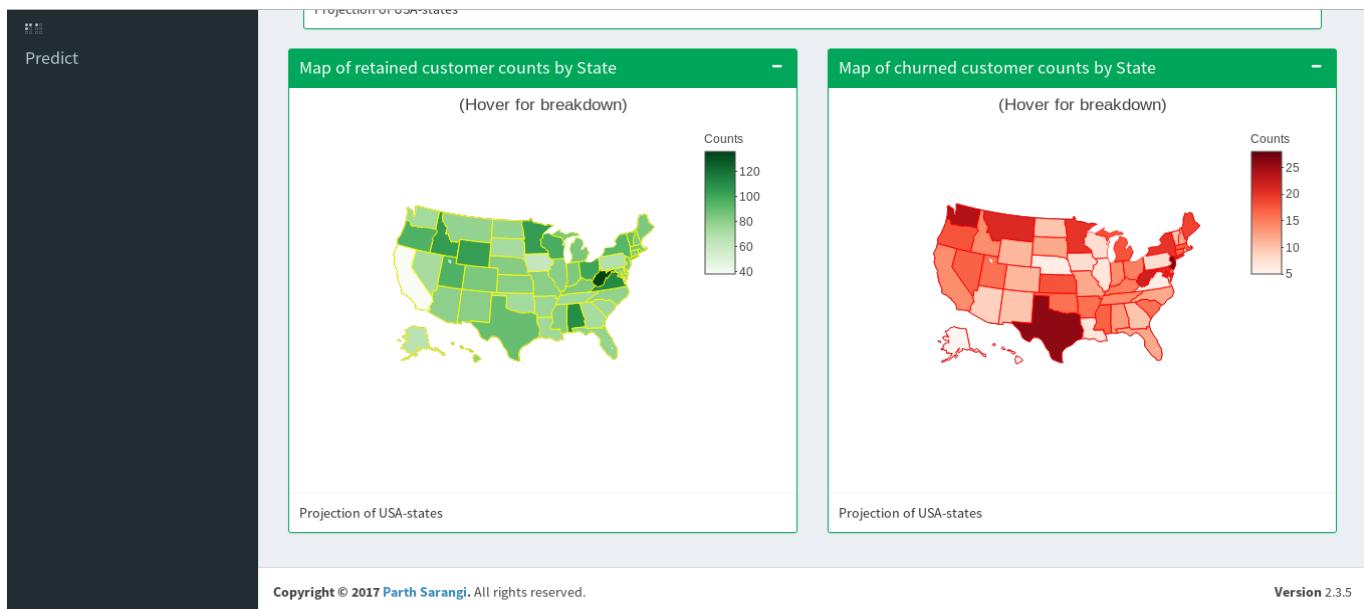


Figure 7.8: Map visualization of churners and non-churners state wise.



7.7 Prediction functionality

The screenshot in Figure 7.9 shows the prediction utility of the churn prediction system. A file in csv format is to be uploaded. A sample file format is also provided for reference. Prediction model will be executed when the file gets successfully uploaded.

Figure 7.9: Dashboard - Upload of file for prediction.

The screenshot displays a web-based dashboard titled "Prediction and Retention". At the top, there are four tabs: "File upload", "View input File", "Prediction", and "Retention". The "File upload" tab is active, indicated by a blue border. Below the tabs, there is a message: "Upload sample data in a given csv format and then check prediction. Sample file for download -> [Sample csv file](#)".

Choose CSV File

A file input field shows "test_prediction.csv" with a "Browse..." button. Below it, a blue bar indicates "Upload complete".

Clear selected file

A button labeled "Click Me" is available to clear the selected file.

Download the predicted results. Click the button below

A "Download" button with a downward arrow icon is present.

File has header?

A checkbox labeled "Header" is checked.

Separator

Radio buttons for "Separator": "Comma" (selected), "Semicolon", and "Tab".

Quote

Radio buttons for "Quote": "None", "Double Quote" (selected), and "Single Quote".

Display

Radio buttons for "Display": "Head" (selected) and "All".

Figure 7.10 shows the utility to view the file uploaded to the prediction functionality, and Figure 7.11 show the results of the prediction model.

Figure 7.10: Dashboard - View uploaded file contents.

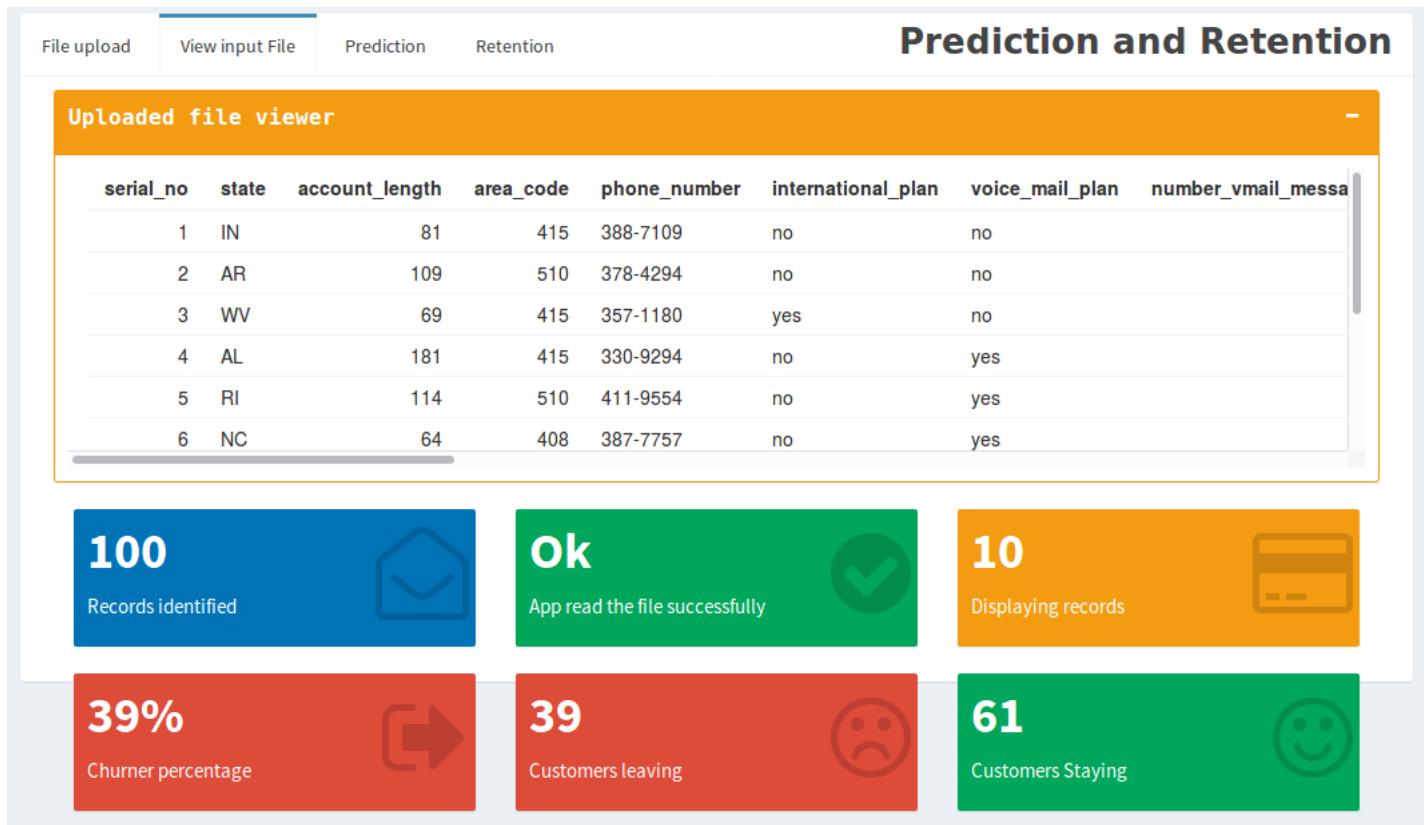


Figure 7.11: Dashboard - View prediction results.

The screenshot shows a web-based dashboard with a light gray header bar. In the top left, there are three buttons: 'File upload', 'View input File', and 'Prediction' (which is highlighted in blue). To the right of these is the title 'Prediction and Retention' in a bold, dark gray font. Below the header is a table with a yellow header row labeled 'Prediction result viewer'. The table has ten rows of data, each representing a customer record with columns for serial_no, churn_predicted, state, account_length, area_code, phone_number, international_plan, voice_mail_plan, and n.

serial_no	churn_predicted	state	account_length	area_code	phone_number	international_plan	voice_mail_plan	n
1	True	WV	69	415	357-1180	yes	no	
2	True	ME	100	510	351-2815	no	no	
3	True	UT	49	415	394-4520	no	no	
4	True	MI	108	408	341-9890	yes	no	
5	True	FL	113	415	395-3867	no	no	
6	True	KS	70	415	331-5650	no	no	
7	True	TX	52	415	364-9904	no	no	
8	True	WV	131	510	350-7785	no	no	
9	True	WA	171	408	419-1863	no	no	
10	True	VT	95	510	378-3508	yes	yes	

7.8 Retention functionality

Retention of churning customers is modeled by rules of the decision table. Figure 7.12 shows the retention strategies suggested by the system. Upselling of plans and benefits with freebies are some retention schemes suggested to retain the customer.

Figure 7.12: Dashboard - View retention strategies.

File upload	View input File	Prediction	Retention	Prediction and Retention							
Retention Strategies										-	
serial_no	Comments	churn_predicted	state	account_length	area_code	phone_number	international_plan	voice_mail_plan	customer_lifetime_value		
1	Upsell with more benefits of international plan usage		True	WV	69	415	357-1180	yes	no		
2	Upsell with more benefits of international plan usage		True	MI	108	408	341-9890	yes	no		
3	Upsell with more benefits of international plan usage		True	VT	95	510	378-3508	yes	yes		

Chapter 8

Conclusion and Recommendations

In this section the conclusion and recommendations for thesis is presented.

8.1 Conclusion

This study presents a web based system with dashboard and prediction ability of managing churn for a telecommunications company. Telecommunications companies experience high rate of customer churn in established markets. Increasing churn means the company is indifferent to the requests of customers, and it increases average revenue per user (ARPU). A higher ARPU means a higher dependency on the reduced customer base, hence it is most important to control churn. Acquiring a new customers is expensive and hence companies should be able to retain their existing customers. Retained customers help generate more revenue. We have presented a platform integrated with KPI's, charts, maps, OLAP and facility to predict churn and retain customers. The system is developed on R's shiny server framework and the machine learning models are developed with open source libraries for R.

In this study, churn prediction was performed with Orange telecom call details records (CDR) dataset, available in SGI website. The churn dataset has 5000 records and the positive class records count to 707. Hence we used a randomly selected 707 records from the negative class and prepared the sample dataset of 1414 records. The dataset has 21 features, of which the area-code and phone number have been singled out as churn determinants. Remaining 19 features are processed for machine learning.

We modeled eight data mining models viz.; rpart decision tree, ctree decision tree, random forest, C5.0, Boosted C5.0, naive bayes, support vector machine, and neural networks. In the training of SVM and neural networks, features with continuous data were normalized with the z-score, and "state". In addition, 3 nominal features were translated from yes/no or true/false to 1/0 format, for easier mapping by classification algorithms. In conclusion it was found that conditional inference decision tree performed very well with 89% accuracy and high positive class prediction accuracy.

In the final part of the study, a retention system was designed with rule induction via decision table derived from the decision tree used for classification. The decision table was constructed from the decision tree and rules were used to build the retention system. The retention system works only on the records labeled as churners. It identifies the features which are responsible for triggering the churn. Those features are mostly the least consumed by that particular customer.

Thus the system also suggests to control those specific features. Retention strategies are suggested based on the least used feature. It is the intention of this thesis to help customer relationship managers to predict and design strategies to retain customers, since the cost of retention is less than that of acquiring new customers.

8.2 Recommendations

Customer retention and churn prediction are quite interesting case studies. Machine learning algorithms are being used widely to access the churn. The thesis implemented the customer churn with the call detail records data for a telecommunications company. In this new age of internet and electronic economy, it is suggested to perform churn analysis with internet consumption data over a period of time and predict the churn in telecom domain. Customer churn prediction can be performed with data from other domains such as advertising, and website design structures. Whether a certain placement of ad banner on home page affects the affinity of the reader. How clicks does a certain page get.

In the field of analytics the field is wide open to various domains and in the current trend of applying analytic's with decentralized applications or with a combination of blockchain can be explored. In addition to application of standard models for machine learning, usage of deep neural networks can be explored.

Also methods of churn retention can be explored. It is advised to perform study on upselling and cross-selling techniques employed by companies.

References

- Berson, A., Smith, S., & Thearling, K. (1999). *Building data mining applications for CRM* (1st ed.). McGraw-Hill Professional.
- Bhattacharya, C. B. (1998). When customers are members: Customer retention in paid membership contexts. *Journal of the Academy of Marketing Science*, 26(1), 31-44.
- Coussement, K., & Poel, D. Van den. (2008). Churn prediction in subscription services: An application of support vector machines while comparing two parameter-selection techniques. *Expert systems with applications*, 34(1), 313–327.
- Cronin, F. J., Colleran, E. K., Herbert, P. L., & Lewitzky, S. (1993). Telecommunications and growth: The contribution of telecommunications infrastructure investment to aggregate and sectoral productivity. *Telecommunications policy*, 17(9), 677–690.
- Gerpott, T. J., Rams, W., & Schindler, A. (2001, may). Customer retention, loyalty, and satisfaction in the German mobile cellular telecommunications market. *Telecommunications Policy*, 25(4), 249–269.
- Gibert, K., Sànchez-Marrè, M., & Codina, V. (2010). Choosing the right data mining technique: classification of methods and intelligent recommendation.
- Han, J. (1997). Olap mining: An integration of olap with data mining. 2, 1–9.
- karpathy@cs.stanford.edu. (n.d.). *Cs231n convolutional neural networks for visual recognition*. <http://cs231n.github.io/neural-networks-1/>. (Accessed: 2017-03-29)
- Kylin, A. (n.d.). *Apache kylin — home*. <http://kylin.apache.org/>. (Accessed: 2017-03-27)
- Larivière, B., & Poel, D. Van den. (2005). Predicting customer retention and profitability by using random forests and regression forests techniques. *Expert Systems with Applications*, 29(2), 472–484.
- Liao, K.-H., & Chueh, H.-E. (2011). Applying fuzzy data mining to telecom churn management. 259–264.
- Mizerski, R. W. (1982). Journal of consumer research. *An Attribution Explanation of the Disproportionate Influence of Unfavourable Information*, 301-310.
- O, A. A., O, A. B., O, A. I., & R, A. E. (2015). Journal of Emerging Trends in Computing and Information Sciences Modeling & Simulation of a Predictive Customer Churn Model for Telecommunication Industry. 6(11).
- Olle, G. D. O., & Cai, S. (2014). A hybrid churn prediction model in mobile telecommunication industry. *International Journal of e-Education, e-Business, e-Management and e-Learning*, 4(1), 55.
- Poel, D. Van den, & Lariviere, B. (2004). Customer attrition analysis for financial services using proportional hazard models. *European journal of operational research*, 157(1), 196–217.
- PredictionIO. (n.d.). *Predictionio — a quick intro*. <http://predictionio.incubator.apache.org/start/>. (Accessed: 2017-03-27)
- Reinartz, W. J., & Kumar, V. (2003). The impact of customer relationship characteristics on profitable lifetime duration. *Journal of marketing*, 67(1), 77–99.

- Reuters. (2017). *India's reliance jio signs up 72 million paying customers* — reuters. <http://www.reuters.com/article/us-reliance-jio-prime-idUSKBN1722BJ>. (Accessed: 2017-04-01)
- Rstudio, S. by. (n.d.). *Shiny - scaling and performance tuning with shiny apps*. <https://rstudio.com/articles/scaling-and-tuning.html>. (Accessed: 2017-03-27)
- Scikit. (n.d.). *Scikit-learn choose the right estimator*. http://scikit-learn.org/stable/tutorial/machine_learning_map/. (Accessed: 2017-03-27)
- TRAI - Telecom Regulatory Authority of India. (n.d.). *Telecom Subscriptions Reports — Telecom Regulatory Authority of India*. <http://www.trai.gov.in/release-publication/reports/telecom-subscriptions-reports>. (Accessed: 2017-04-01)
- Tsai, C.-F., & Lu, Y.-H. (2009). Customer churn prediction by hybrid neural networks. *Expert Systems with Applications*, 36(10), 12547–12553.
- Tutorials Point. (n.d.). *Data warehousing schemas*. https://www.tutorialspoint.com/dwh/dwh_schemas.htm. (Accessed: 2017-03-31)
- Vafeiadis, T., Diamantaras, K. I., Sarigiannidis, G., & Chatzisavvas, K. C. (2015). A comparison of machine learning techniques for customer churn prediction. *Simulation Modelling Practice and Theory*, 55, 1–9.
- Wei, C.-P., & Chiu, I.-T. (2002). Turning telecommunications call details to churn prediction: a data mining approach. *Expert systems with applications*, 23(2), 103–112.
- Xie, Y., Li, X., Ngai, E., & Ying, W. (2009). Customer churn prediction using improved balanced random forests. *Expert Systems with Applications*, 36(3), 5445–5449.