

**COMPREHENSIVE STUDY OF METHODS TECHNIQUES AND TOOLS USED
FOR DATA ANALYTICS IN FIELDS OF FINANCE**

by

Parth Sarangi

A special study report 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 2017

Table of Contents

Chapter	Title	Page
	Title Page	i
	Table of Contents	ii
	List of Figures	iii
	List of Tables	iv
1	Introduction	1
	1.1 Overview	1
	1.2 Objectives	1
2	Literature Review	2
	2.1 Customer Churn & Retention	2
	2.2 OLAP & Datawarehouse	2
	2.3 Data Mining	3
	2.4 Model Evaluation Metrics	11
	2.5 Review of Selected Research Papers	13
	2.6 Summary of Selected Research Studies	14
3	Methodology	19
	3.1 Research Methodology	19
	3.2 Data Preprocessing and Datawarehouse Development	20
	3.3 Development and Evaluation of the Prediction Models	22
	3.4 System Development & Evaluation	23
	3.5 Timeline	25

List of Figures

Figure	Title	Page
2.1	OLAP Solution - Apache Kylin	3
2.2	Mammal classification problem	4
2.3	A sample neural network	5
2.4	Kohonen SOM	6
2.5	Select the Right Mining Technique	7
2.6	Another approach to select the Data Mining. Reprinted from Scikit	8
2.7	PredictionIO Engine interaction with Apps and Prediction Engine	9
2.8	R Shiny architecture	10
2.9	Confusion Matrix	12
3.1	Research Methodology	20
3.2	Data preprocessing	21
3.3	OLAP Star Schema	21
3.4	The Intelligent Churn Prediction Architecture	23
3.5	Gantt chart tasks	25

List of Tables

Table	Title	Page
2.1	Previous literature review.	14

Chapter 1

Introduction

1.1 Overview

Data analysis is an very robust topic in the field of data science and encompasses the various mathematical functions. The functions are statistical in nature and are performed on the data obtained. The goal of data analytics is to support (or reject) the hypothesis which the data scientist postulates. “ By processing a steady stream of real-time data, organizations can make time-sensitive decisions faster than ever before, monitor emerging trends, course-correct rapidly and jump on new business opportunities.” (?, ?).

This paper tries to enlist most of the up-to-date techniques used by researchers and mathematicians to make sense of the data. Also the paper presents them in three groups of analytics. But then there arises a question which is, why the need for data analytics ? Well, to answer that the paper proposes the literature from an article of Kdnuggets (?, ?).

Data analytics is already the next big disrupter in the financial sector (?, ?). New startups are build solutions to automate mundane manual tasks of reconciliations and consolidations. Artificial intelligence is helping companies to redesign the traditional process flows and restructure work-flows for optimization. With complex algorithms of machine learning, artificial intelligence, “Bot’s“ are designed which learn the product usage and popularity among customers. These “Bot’s“ can also interact autonomously to users and deduce patterns of interaction. Such an existing example is cited in Cnbc’s article (?, ?), about Bank of America deploying “Erica“ an digital assistant based on models of artificial intelligence.

1.2 Objectives

The overall objective of the study report is to build an understanding of the tools and techniques of analytics used by financial analysts and institutions. This study presents a brief of the current state of affairs analytics and subsequently a study of the techniques. Also there is a demonstration of some common methods.

Chapter 2

Literature Review

This thesis 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 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 (? , ?). 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 (? , ?).

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 a 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 inaccessible 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 (? , ?). Olap consists of four main processes viz., Drill-down, Roll-up, slice and dice. Multi-dimensional 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™ (? , ?). Shown in Figure 2.1 is the architecture of the product.

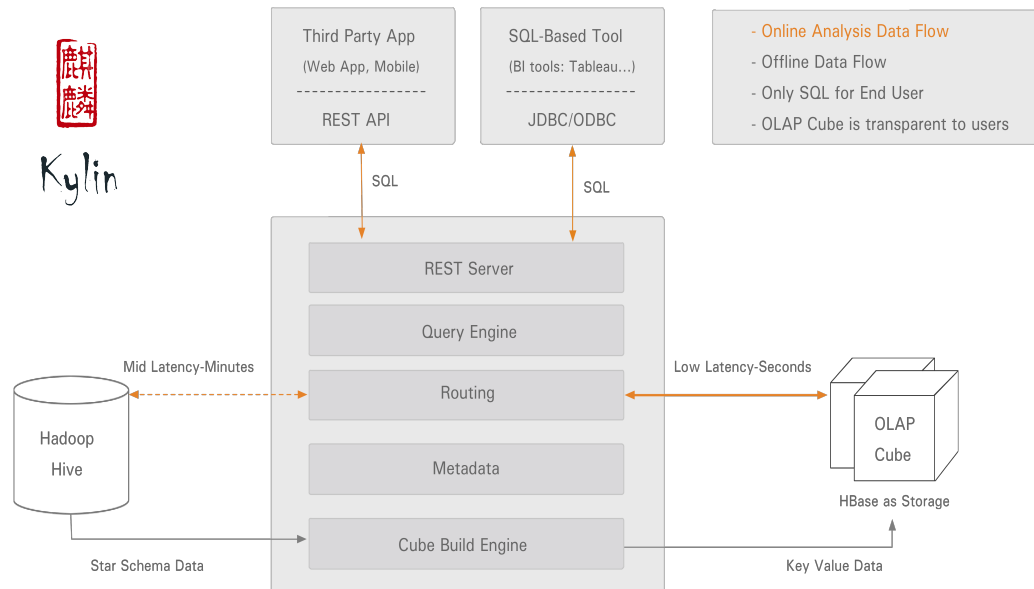


Figure 2.1: OLAP Solution - Apache Kylin

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, $y = \beta_0 + \beta_1x + \varepsilon$
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**.

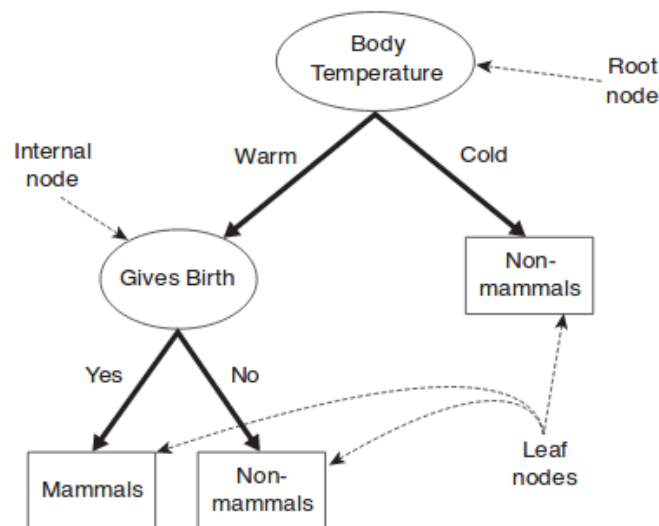


Figure 2.2: Mammal classification problem

- 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 (? , ?). 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.

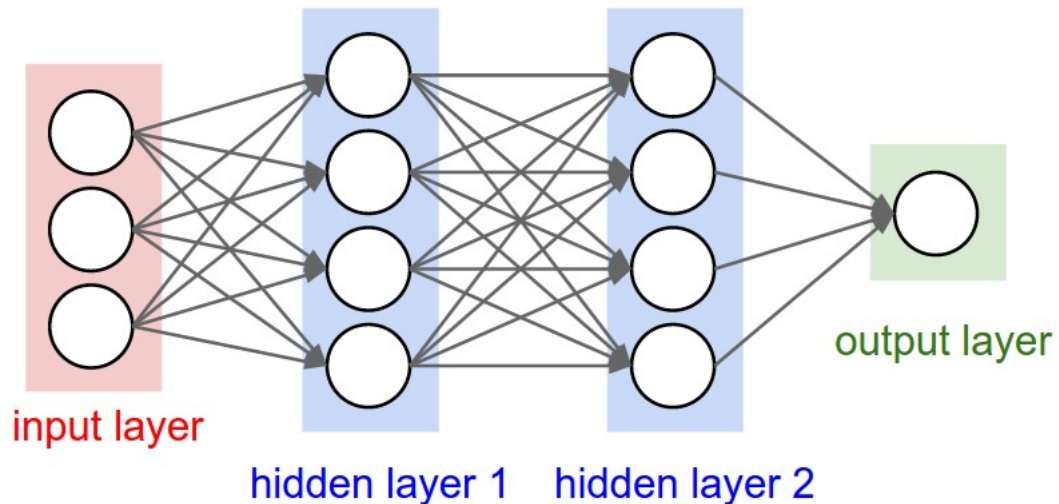


Figure 2.3: A sample neural network

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 see 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.

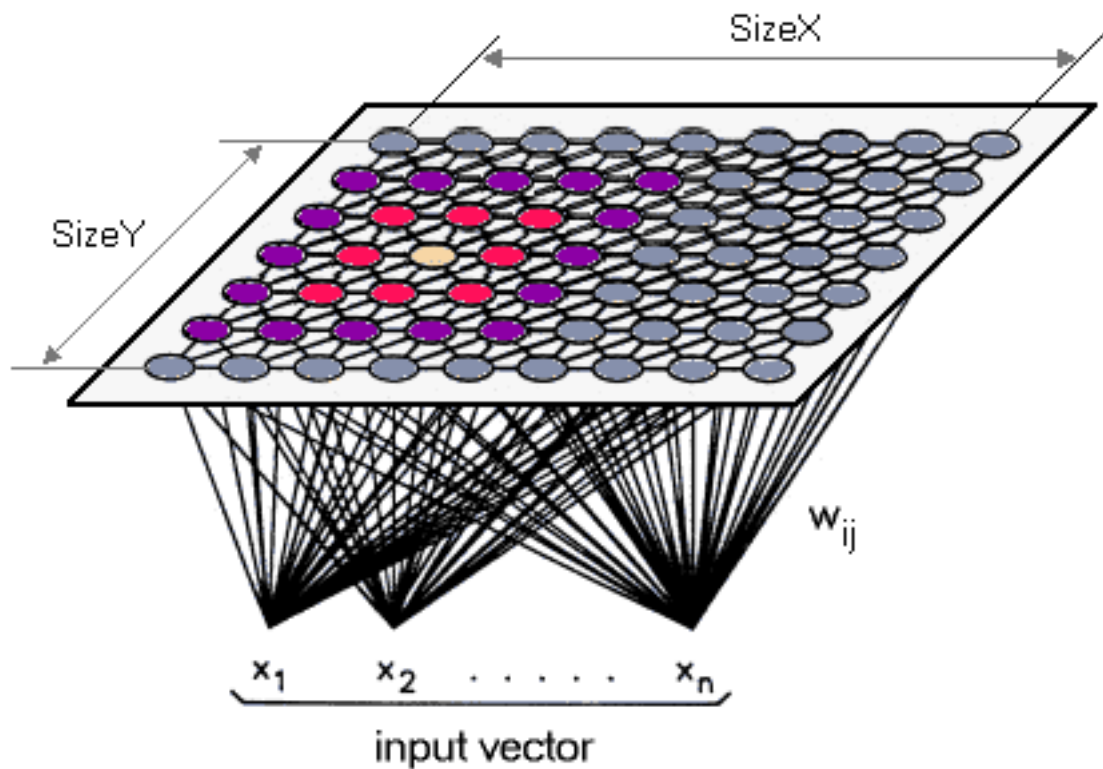


Figure 2.4: Kohonen SOM

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” (?, ?) shows the approach which could be taken to select between the various models available for data mining.

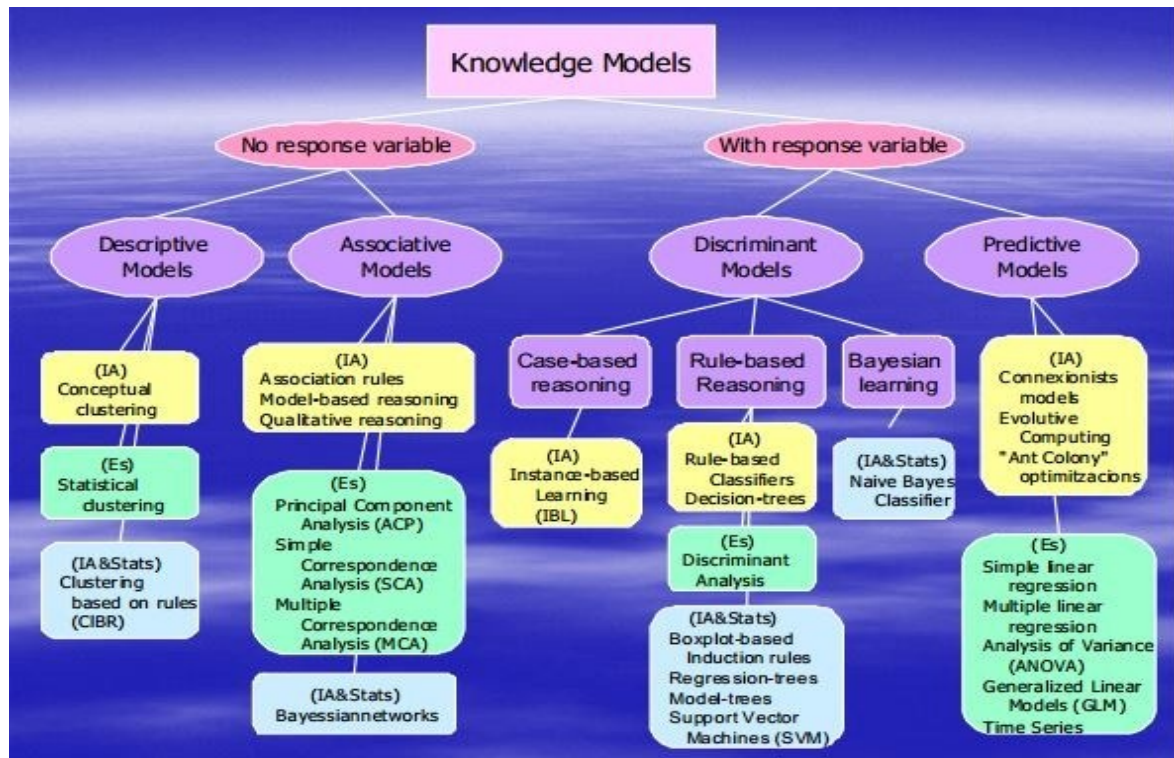


Figure 2.5: Select the Right Mining Technique

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 (?, ?).

- 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 provides various templates for varied mining algorithms. Classification templates like Decision trees, Logistic Reressiong, NLP are available for use.

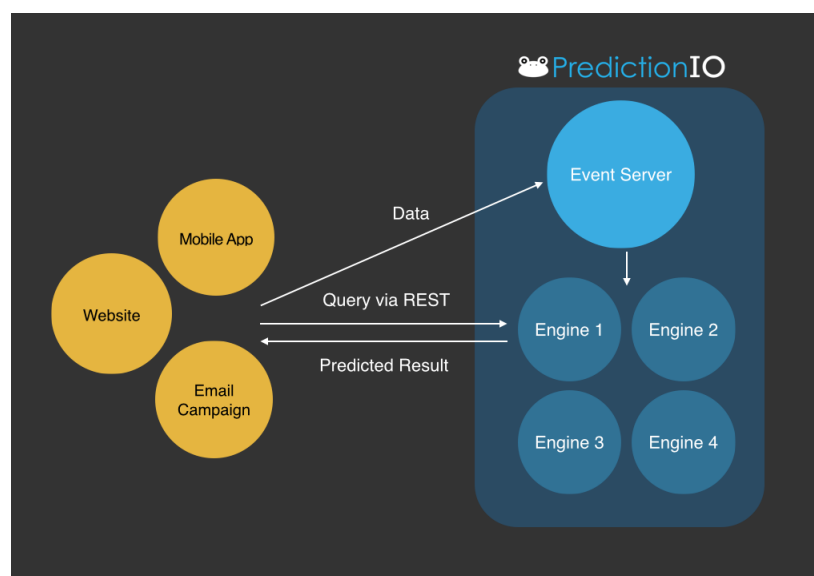


Figure 2.7: PredictionIO Engine interaction with Apps and Prediction Engine

- **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 (? , ?).

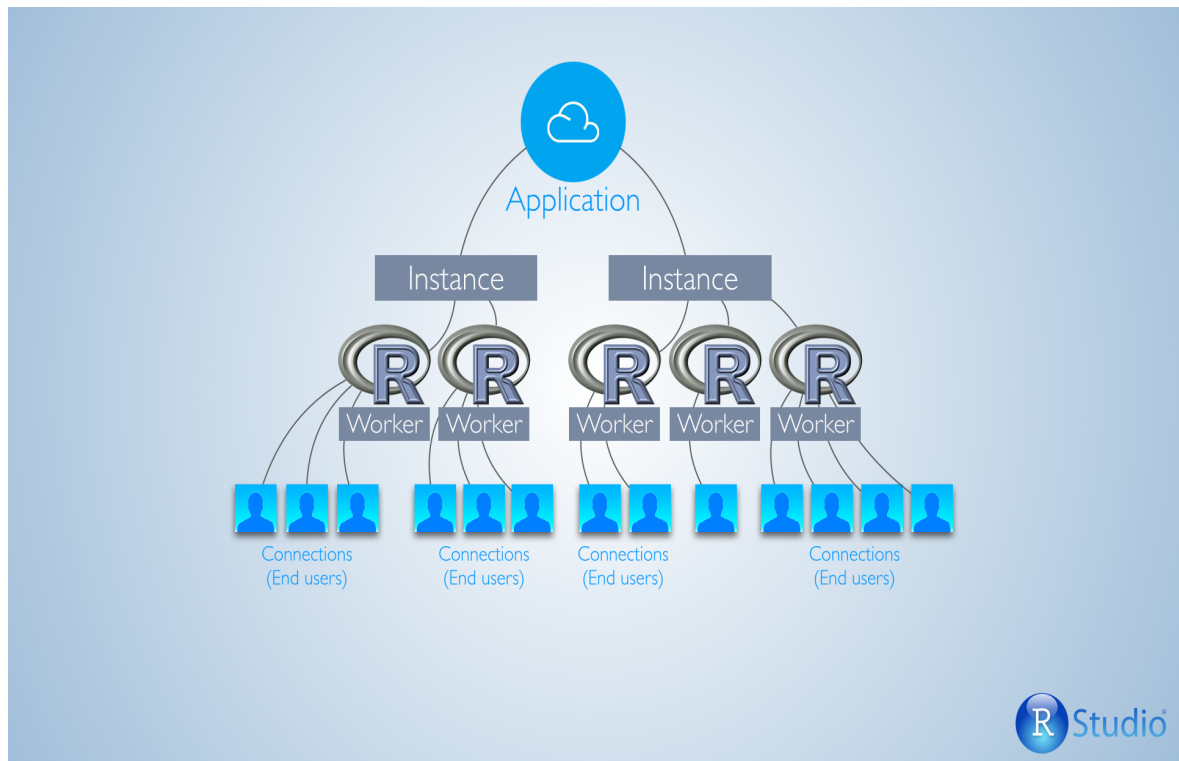


Figure 2.8: R Shiny architecture

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**.

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

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 (?, ?). 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 ” (?, ?) 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 (?, ?) 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 review.

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
1	Modeling & Simulation of a Predictive Customer Churn Model for Telecommunication Industry (?, ?)	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 Precision 80.86% recall 92.7% and predicted accuracy 95.8%.	None suggested
2	A Hybrid Churn Prediction Model in Mobile Telecommunication Industry (?, ?)	A model combined with VotedPerceptron and Logistic 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 separately.	None suggested

SNo	Title & Author	Objective	Data & Methodology	Outcome	Further Research
3	A comparison of machine learning techniques for customer churn prediction (?, ?)	The normal model functions were performance compared to their corresponding boosted models.	<p>Data : publicly hosted churn dataset at UCI machine learning repository.</p> <p>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.</p>	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 (?, ?)	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.	<p>Data : Telecom company of Taiwan. Contractual and call details of subscribers Oct 2000 – Jan 2001. 9100000 records.</p> <p>Method : Multi classifier class combiner, Decision tree C4.5</p>	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 (?, ?).	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 (?, ?).	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 (?, ?)	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 practionerx 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 (?, ?)	Very few studies for hybrid data mining ap- praoch 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 : $ANN + ANN > single ANN$ $3 * 3$ 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 (?, ?)	The paper discusses more than one variable of retention and profit outcome.	<p>Data : 100,000 Belgian finance company. Divided into 2 random parts, one for estimation other for evaluation.</p> <p>Method : Authors used random forest for regression to predict profitability, next purchase and defection decision. Benchmarked to linear regression model.</p>	Random forest are better than logit and linear regression.	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.

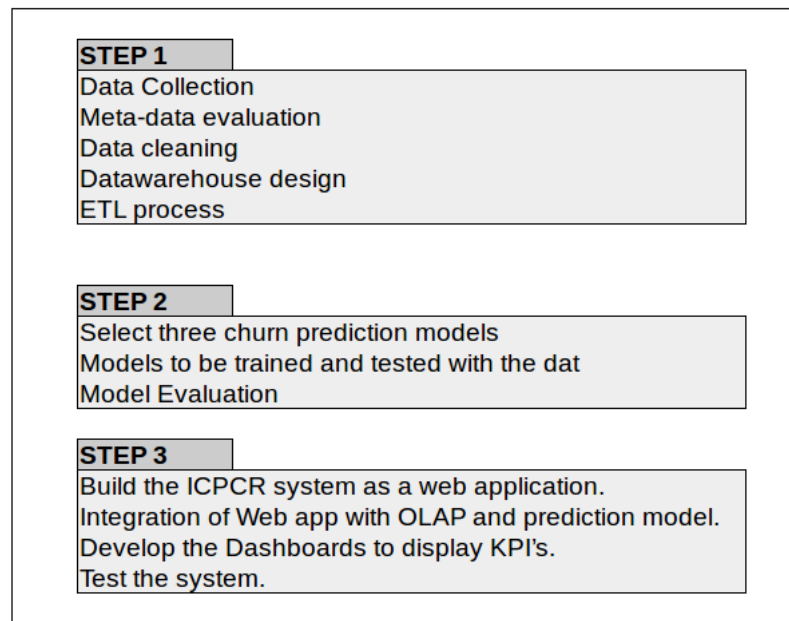


Figure 3.1: Research Methodology

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-usable data, either by replacing with suitable or by entirely removing it. Un-usable 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.

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

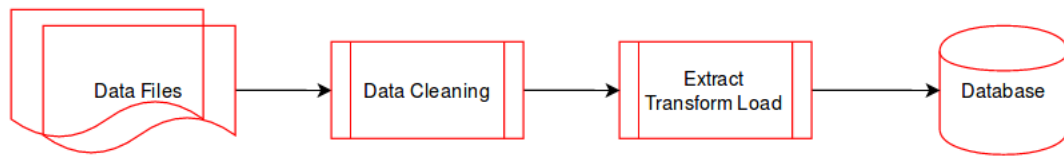


Figure 3.2: Data preprocessing

OLAP will be designed with the star schema. Figure 3.3 shows a typical implementation of the star schema (? , ?). 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.

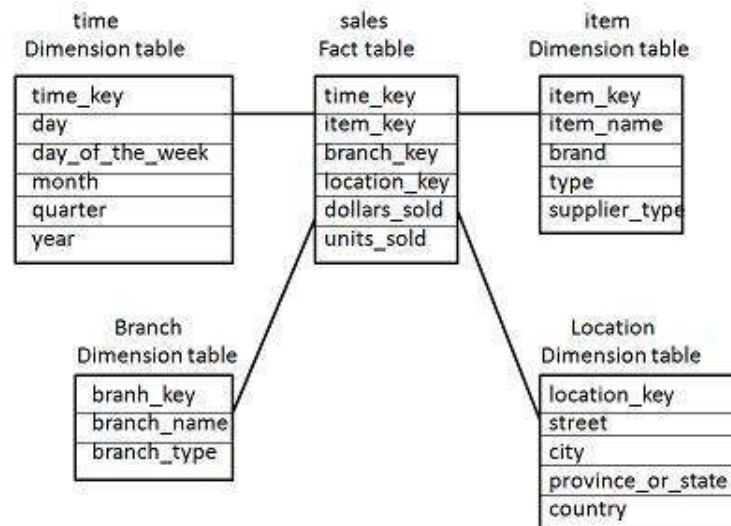


Figure 3.3: OLAP Star Schema

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 Vector 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 MLib 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 (S, S).

- K-Fold Cross Validation : It is proposed to perform this process to make the classification model more accurate. From previous literature it is learned that $k = 10$ 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 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.

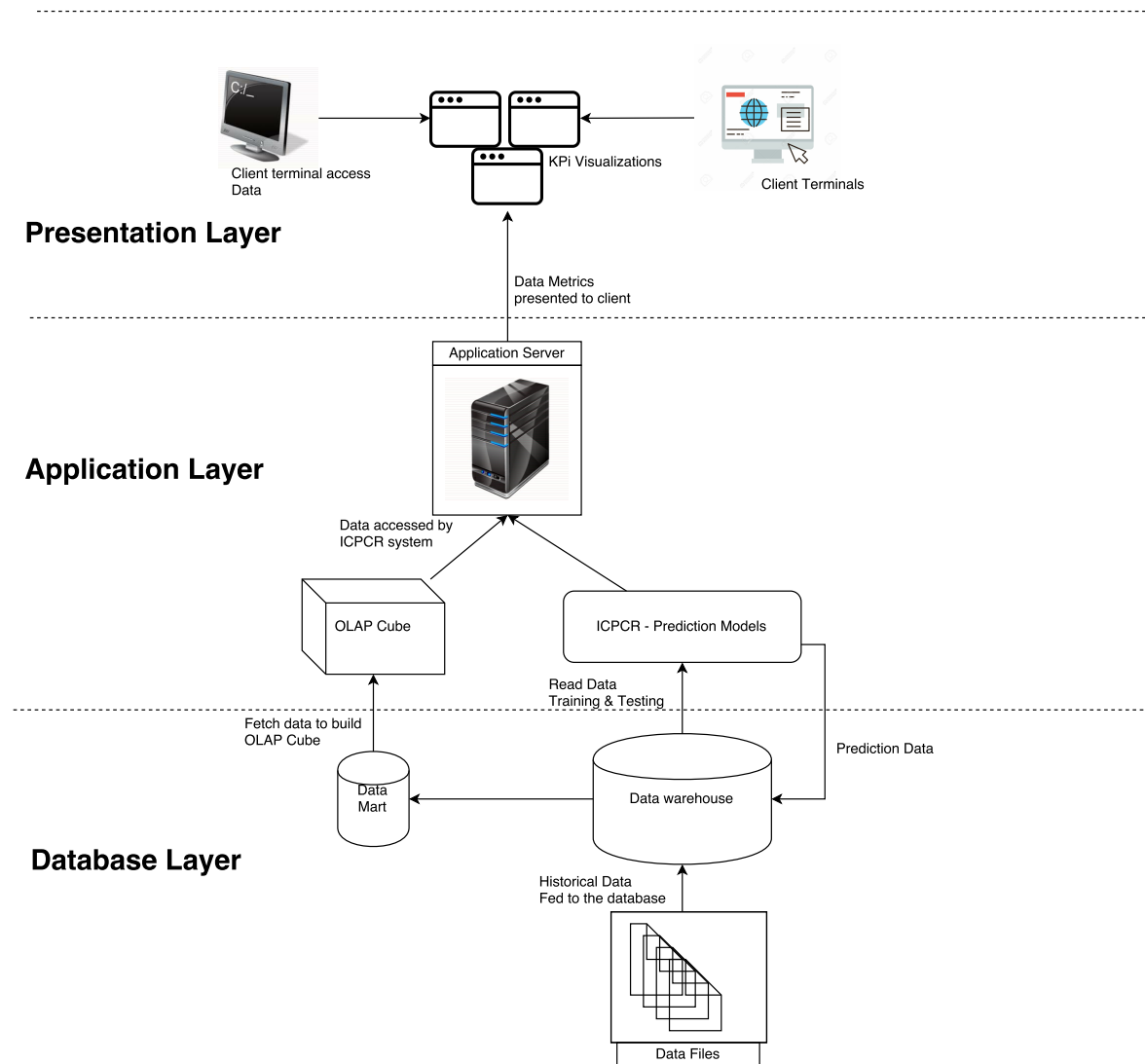


Figure 3.4: The Intelligent Churn Prediction Architecture

3.4.1 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.

3.4.2 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.

3.4.3 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.

3.4.4 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.

3.5 Timeline

The forecast of the tasks to be carried out in this thesis are shown below in a Gantt chart Figures 3.5.

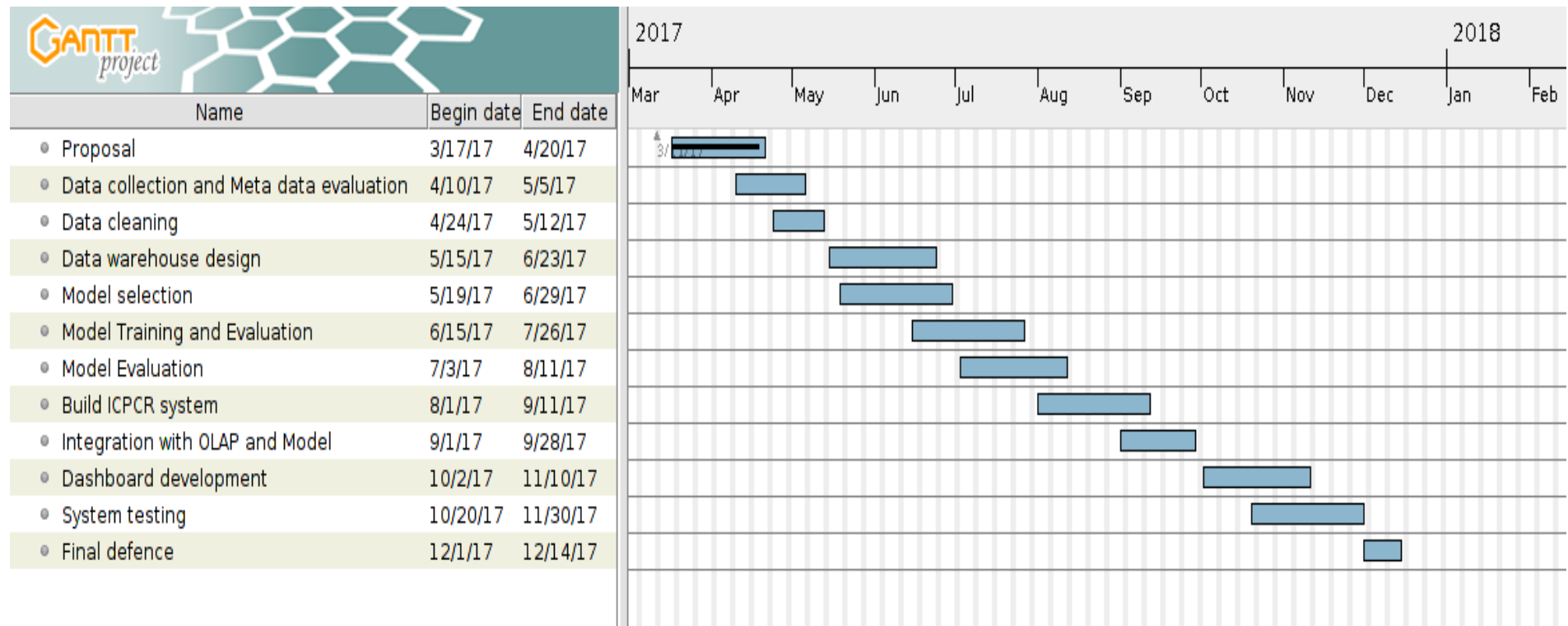


Figure 3.5: Gantt chart tasks

References

- Balm, P. (2015). *10 reasons why i love data and analytics* — *kdnuggets*. <http://www.kdnuggets.com/2015/06/10-reasons-love-data-analytics.html>. (Accessed: 2017-05-10)
- Culp, S. (2017, Feb). *Artificial intelligence is becoming a major disruptive force in banks' finance departments* — *forbes*. <https://www.forbes.com/sites/steveculp/2017/02/15/artificial-intelligence-is-becoming-a-major-disruptive-force-in-banks-finance-departments/39f25de14f62>. (Accessed: 2017-06-27)
- Das, T., & Kumar, P. M. (2013). Big data analytics: A framework for unstructured data analysis. *International Journal of Engineering Science & Technology*, 5(1), 153.
- Taylor, H. (2016, October). *Bank of america launches chat bot erica* — *cnbc*. <http://www.cnbc.com/2016/10/24/bank-of-america-launches-ai-chatbot-erica--heres-what-it-does.html>. (Accessed: 2017-06-27)