# A Comparative Study of Classification Techniques for Fire Data Set

Rachna Raghuwanshi
*M.Tech CSE*
*Gyan Ganga Institute of Technology & Science, Jabalpur*

**Abstract:**Classification of data has become an important research area. The process of classifying documents into predefined categories based on their content is Text classification. It is the automated assignment of natural language texts to predefined categories. The primary requirement of text retrieval systems is text classification, which retrieve texts in response to a user query, and text understanding systems, which transform text in some way such as answering questions, producing summaries or extracting data In this paper, we study classification techniques and use of them in FIRE data set is presented. The experimental results show that the proposed system works as a successful text classifier.

**Keywords: Classification, Rapid Miner tool, Decision tree, Naïve Bayes , K-NN, Fire data set, Cross Validation.**

## INTRODUCTION:

There are numerous text documents available in paper format such as books, news paper, magazines etc. More and more are becoming available every day. We also have FIRE data set which represents a massive amount of information that is easily accessible. In this huge collection seeking a value requires organization; data mining automated much of the work of organizing documents. The accuracy and our understanding of such systems greatly influence their usefulness. The task of data mining is to automatically classify documents into predefined classes based on their content. Many algorithms have been developed to deal with automatic text classification [10]. Various techniques of classification are ANN, Genetic Algorithms (GAs) / Evolutionary Programming (EP), Decision Tree, Naive Bayes, KNN, Clustering, Support Vector Machine, Rough set, Logistic Regression etc. But in this paper only three methods of classification are used to calculate various results with the help of Rapid Miner tool.

*Rapid Miner:* Rapid Miner is the most popular open source software in the world for data mining and strongly supports text mining and other data mining techniques that are applied in combination with text mining. The power and flexibility of Rapid Miner is due to the GUI-based IDE (integrated development environment) it provides for rapid prototyping development of data mining models, as well as its strong for scripting based on XML (extensible markup language). The visual modeling in the Rapid Miner IDE is based on the defining of the data mining process in terms of operators and the flow of process through these operators. User specify the expected inputs , the delivered outputs, the mandatory and optional parameters and the core functionalities of the operators and the complete process is automatically executed by Rapid miner.

- *Decision Tree:* The decision tree is a structure that includes various nodes such as root node, branch and leaf node. Each leaf node holds the class label, each branch denotes the outcome of test and internal node denotes a test on attribute. The root node is topmost node in the tree. There are many popular decision tree algorithms CART, ID3, C4.5, CHAID, and J48. The decision tree approach is more powerful for classification problems. There are two steps in this techniques building a tree & applying the tree to the dataset.

- *Naive Bayes:* Naive Bayes classifier is based on Bayes theorem. Naïve Bayes classifier algorithm uses conditional independence properties, means it assumes that an attribute value on a given class is independent of the values of other attributes. The Bayes theorem is as follows: Let $X=\{x_1, x_2... x_n\}$ be a set of n attributes. In Bayesian, X is considered as evidence and H is some hypothesis means, the data of X belongs to specific class C. We have to determine P (H|X), the probability that the hypothesis H holds given evidence i.e. data sample X. According to Bayes theorem the P (H|X) is expressed as[1]

    P (H|X) = P (X| H) P (H) / P (X).

- *K-Nearest Neighbor:* The *k*-nearest neighbor's algorithm (*K*-NN) is a method for classifying objects based on closest training data in the feature space. K-NN is a type of instance-based learning. The *k*-nearest neighbor algorithm is amongst the simplest of all machine learning algorithms. But the accuracy of the *k*-NN algorithm can be severely degraded by the presence of noisy or irrelevant features, or if the feature scales are not consistent with their importance.[1]

### LITERATURE REVIEW:

✓ Survey of Classification Techniques in Data mining By Thair Nu Phyu (2009):     In this paper author explains the basic classification techniques. Various kind of classification methods such as decision tree induction, Bayesian networks, K-nearest neighbor classifier, case based reasoning, genetic algorithm and fuzzy logic. The goal of this survey is to provide comprehensive review of techniques.

✓ A Decision Tree Algorithm for distributed Data Mining: Towards Network Intrusion Detection By

Baik, S. Bala(2004): This paper presents preliminary works on an agent-based approach for distributed learning of decision trees. The distributed decision trees approach is applied to intrusion detection domain, the interest of which is recently increasing. In this, approach a network profile is built by applying a distributed data analysis method for the collection of data from distributed hosts. Several experiments show that distributed decision tree performance is much better than the non-distributed decision tree.
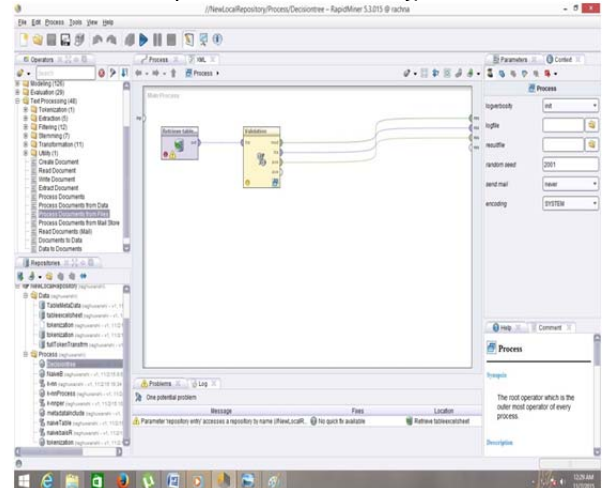
✓ Naïve Bayes Classifiers that perform well with continuous variables by Remco R. Bouckaert(2004): There are three main methods for handling continuous variables in naïve Bayes classifiers, namely the normal method(parametric approach), the kernel method (non-parametric approach) and discretisation. In this paper comparison is performed on these three methods, which shows large mutual differences of each of the methods and no single method is better. By using continuous variable the performance of Naïve Bayes can be increases.

✓ Kidney Disease Prediction using Data Mining Classification Technique by Suman Bala, Krishan Kumar(2014): In this paper data mining is used for predicting disease from datasets used Analytical methodology for detecting unknown and valuable information in health data. Various techniques are compared to find best method for prediction.

✓ Problem of Bayesian Network Classifier By Cheng(2002): He shows the various drawbacks of Bayesian Network in his paper and also give solutions. He concluded that Bayesian network is not suitable for dataset with many features. Reason is to construct large network is not feasible in terms of time and space.

✓ Automatic Discovery of Similar Words, in "Survey of Text Mining: Clustering, Berry Michael W (2004),: In this paper Berry explains about the text mining and how it is similar to data mining its various aspect and operations.

✓ A Survey of Text Mining Techniques and Applications By Vishal Gupta & Gurpreet S. Lehal(2009): Text Mining has become an important research area. Text Mining is the discovery by computer of new, previously unknown information, by automatically extracting information from different written resources. In this paper, a Survey of Text Mining techniques and applications have been s presented.

## METHODOLGY

**Classification:** - Classification is a very important data mining task, and the purpose of classification is to propose a classification function or classification model (called classifier).The classification model can map the data in the database to a specific class. Classification construction methods include: Decision Tree, Naive Bayes, KNN, ANN, Clustering, Support Vector Machine, Rough set, Genetic Algorithms (GAs) / Evolutionary Programming (EP), Logistic Regression etc. In this paper only three techniques
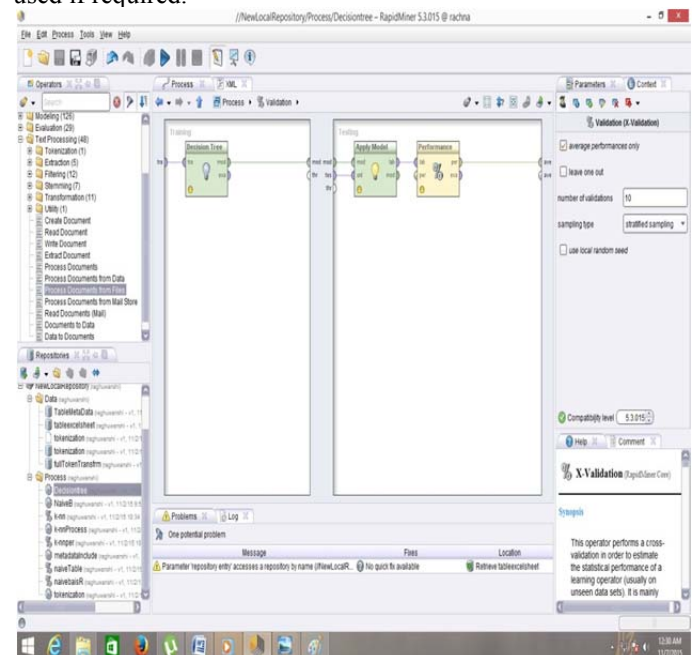
are used for text preprocessing and cross validation and its sub process testing and training are also applied.

*Validation-* First applies validation on table which is retrieve from data set. This operator performs cross validation in order to estimate the statical performance of a learning operator. The X-validation operator is a nested operator; it has two sub processes testing and training. The X-validation operator returns performance vector. This mainly used to estimate the accuracy of particular model. Cross Validation process is shown in figure 2.



**Fig 2-** Validation operator

*Testing & Training-* These are the sub processes of validation operator. The training sub process is used for training a model. The trained model is then applied in the testing sub process. During the testing phase performance of the model is also measured. The testing sub process must return a Performance Vector and other results also. This is usually generated during measuring performance of the model. Two such ports are provided but more can also be used if required.



**Fig 3-** Testing and Training

*Decision Tree:* In every recursion the DECISION TREE algorithm follows the following steps:

- An attribute A is selected to split on. Making a good choice of attributes to split on each stage is crucial to generation of a useful tree. The attribute is selected depending upon a selection criterion which can be selected by the criterion parameter.
- Examples in the FIRE Data Set are sorted into subsets, one for each value of the attribute A in case of a nominal attribute. In case of numerical attributes, subsets are formed for disjoint ranges of attribute values.
- A tree is returned with one edge or branch for each subset. Each branch has a descendant sub tree or a label value produced by applying the same algorithm recursively.

In general, the recursion stops when all the examples or instances have the same *label* value, i.e. the subset is pure. This is a generalization of the first approach; with some error threshold. However there are other halting conditions such as:

- There are less than a certain number of instances or examples in the current sub tree. This can be adjusted by using the minimal size for split parameter.
- No attribute reaches a certain threshold. This can be adjusted by using the minimum gain parameter.
- The maximal depth is reached. This can be adjusted by using the maximal depth parameter.

For splitting the attribute one of the attribute is selected for that there must be a criteria. It can have one of the following values:

- **information_gain**: All attributes entropy is calculated. The attribute with minimum entropy is selected for split. This method has a bias towards selecting attributes with a large number of values.
- **gain_ratio**: It is a variant of information gain. It adjusts the information gain for each attribute to allow the breadth and uniformity of the attribute values.
- **gini_index**: This is a measure of impurity of an Fire data set. Splitting on a chosen attribute gives a reduction in the average gain index of the resulting subsets.
- **Accuracy**: To maximizes the accuracy of the whole Tree an attribute is selected for split

*Naïve Bayes:* Studies comparing classification algorithms have fond a simple Baiyesian classifier known as the naïve Bayesian classifier to be comparable in performance with decision tree and selected neural network classifiers. Bayesian classifiers as comparison to other have also exhibited high accuracy and speed when applied to large databases.

Naïve Bayesian classifiers assume that the effect of an attribute value on a given class is independent of the values of the other attributes. This assumption is called class conditional independence. It is made to simplify the computations involved and, in this sense, is considered "naive". Bayesian belief networks are graphical models subsets of attributes. can also be used for classification.

*Naïve Bayes Operator:* - This operator generates a Naive Bayes classification model.

*K-Nearest Neighbor:* The k-nearest-neighbor method was first described in early 1950s, when increased computing power become available. Ti has since widely used in the area of pattern recognition.

It is based on learning by analogy, that is by comparing a given test tuple with training tuples that are similar to it. The Training tuples are described by n attributes each tuple represents a point in an n-dimensional space. In this way, all of the training tuples are stored in an n-dimensional space. When given an unknown tuple, a k-nearest-neighbor classifier searches the pattern space for the k training tuples that are closest to the unknown tuple. These k training tuples are the k "nearest-neighbor" of the unknown tuple.

"Closeness" is defined in terms of a distance metric, such as Euclidean distance. The Euclidean distance between two pints or tuples, say $X_1=(x_{11}, x_{12}, \ldots\ldots, x_{1n})$ and $X_2=(x_{21}, x_{22}, \ldots\ldots, x_{2n})$ , is

$$dist(X_1,X_2)=\sqrt{\sum(x_{1i}-x_{2i})^2} \qquad (1)$$

In other words, for each numeric attribute, we take the difference between the corresponding values of that attribute in tuple $X_1$ and in tuple $X_2$, square this difference, and accumulate it. The square root is taken of the total accumulated distance count. Typically, we normalize the values of each attribute before using Equation (1).

*K-NN Operator-* This operator generates a k-Nearest Neighbor model from the input ExampleSet. This model can be a classification or regression model depending on the input ExampleSet.

## RESULT:

On the newspaper data first process document from file then text processing operators are used such as tokenization, filter tokens, steaming, transformation and stop word. After that result is shown in data view format in rapid miner. As shown in fig 4. Then convert the data into excel table, in rapid miner import excel sheet and save it into the repository. After that various method of classification are applied. First apply validation on table which is retrieve from data set, cross validation is perform to estimate the statical performance of learning. The X-validation operator is a nested operator; it has two sub processes testing and training. These are the sub processes of validation operator. The training subprocess is used for training a model. The trained model is then applied in the testing subprocess. The performance of the model is also measured during the testing phase. The testing subprocess must return a Performance Vector.

**Decision Tree Result**

Now decision tree method is applied on the excel sheet and it also include validation and its sub process testing and training. Results of decision tree are shown in the below figures 5, 6, 7. It is the example set of decision tree. In rapid miner various types of view can see of data such as data view, Meta data view, plot view, advanced charts, annotations. In this figure data view of the data set is shown Example set consist (24 examples, special attribute and 5 regular attribute). First row of the table shows row no. second real, third name of attribute, fourth type of attribute etc. Figure 4 shows decision tree graph view in this graph

data set text type is shown there are two types of words in data set, Real and text. Figure 5 shows the performance vector of decision tree, accuracy rate of decision tree is 88.33%
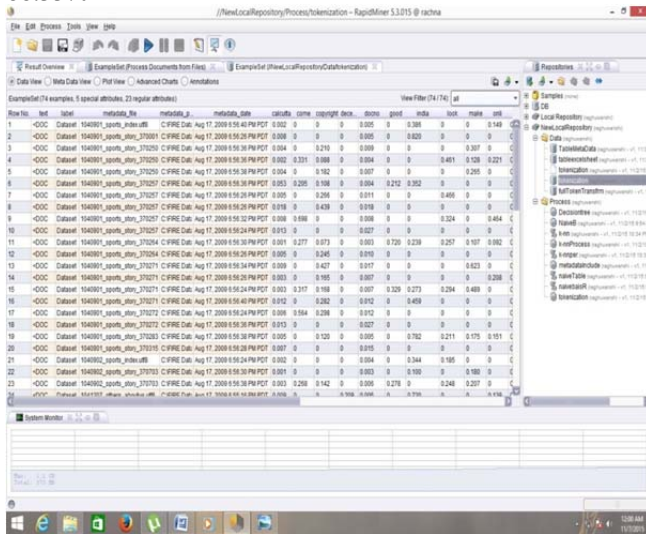


**Fig 4:** Result of Tokenization1

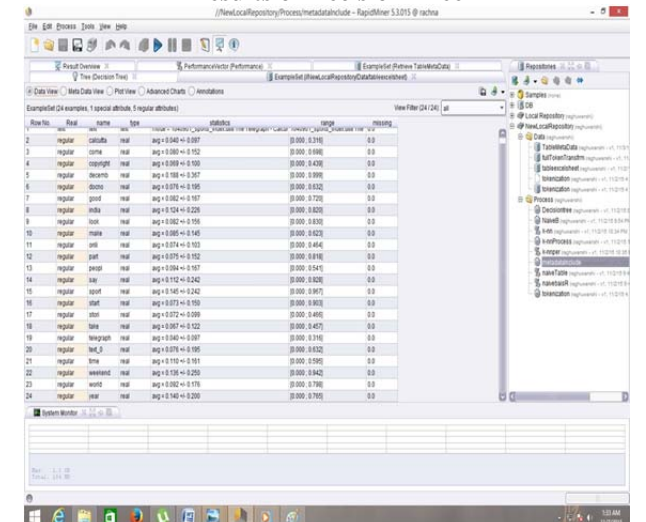### Results of Decision Tree
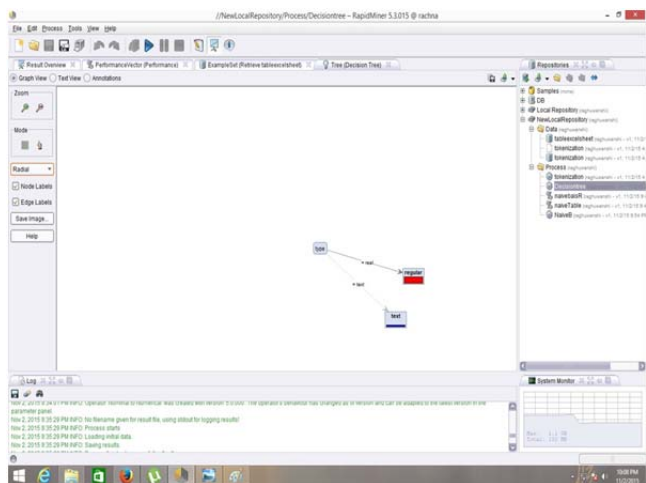


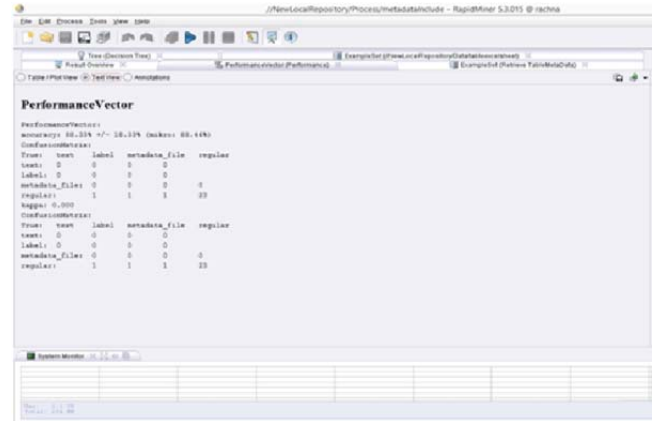**Fig 5 :-** Decision tree Example Set



**Fig 6 -**Decision Tree Graph view



**Fig 7-** Decision Tree Performance Vector Text View

**K-Nearest-Neighbor Classifier Results:**
 K-Nearest-Neighbor method is applied on the excel sheet and it also include validation and its sub process testing and training. Results of K-Nearest-Neighbor are shown in the below figures 8, 9.
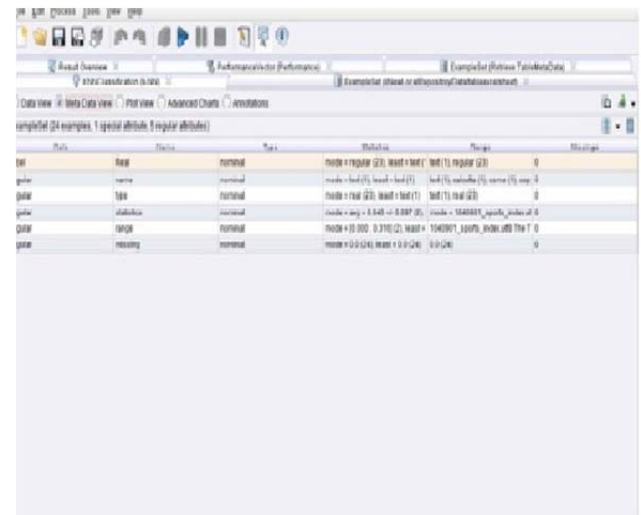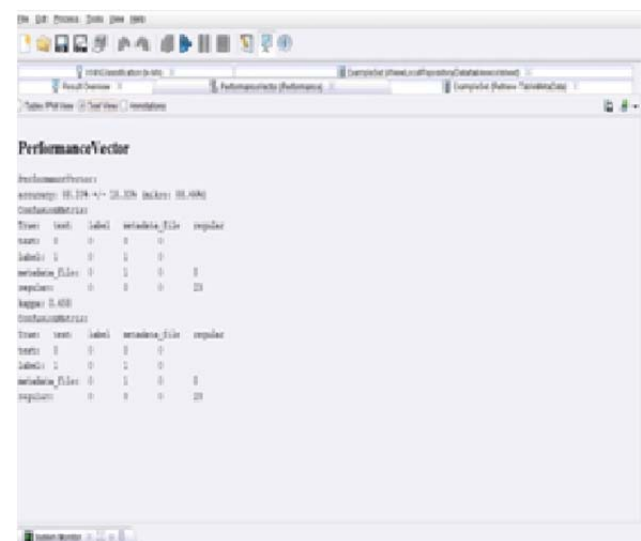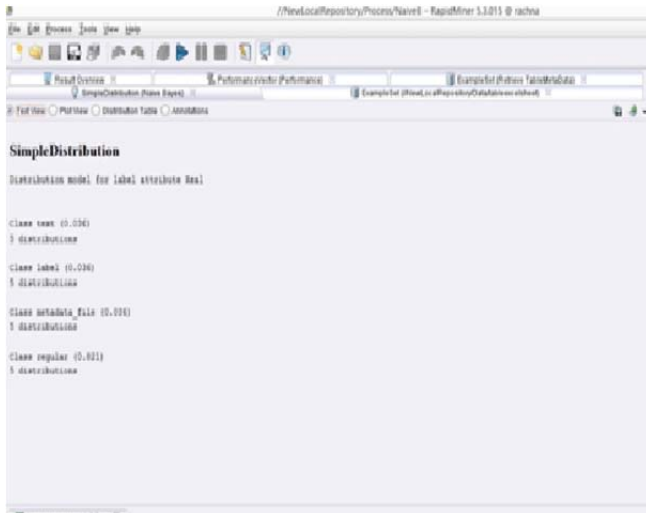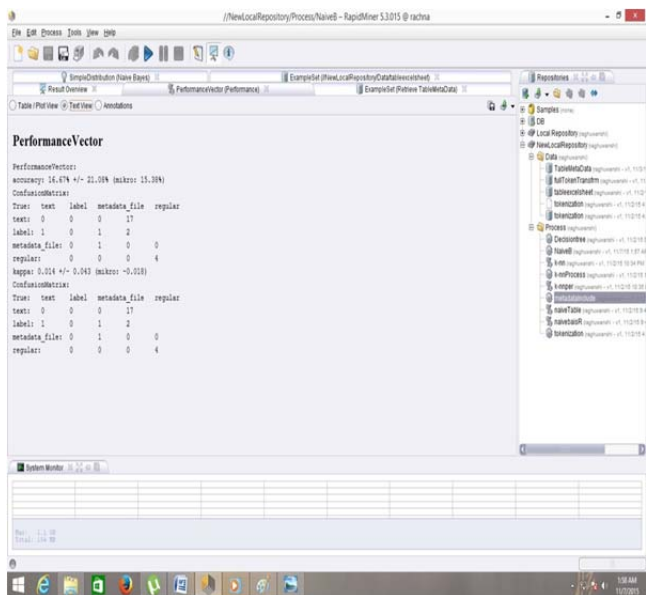


**Fig 8-**K-NN Meta Data View



**Fig 9-** K-NN Performance Vector (Text View)

**Results and output of Naïve Bayes Classification-**
Naïve Bayes method is applied on the excel sheet and it also include validation and its sub process testing and training. Results of Naïve Bayes are shown in the below figures. Figure 10 shows simple distribution of real data and figure Performance vector with accuracy rate 16.67%.



**Fig 10** - Naïve Bayes Simple Distribution



**Fig 11-** Naïve Bayes Performance Vector Text view

### CONCLUSION:

The overall objective of our is to study the various data mining classification techniques available to calculate the number of word most commonly used in the FIRE dataset and to compare them to find the best method which shows exact accuracy.

- We analyzed that the most commonly used DM technique such as Decision Trees, ANN and Naïve Bayes, resulting as well-performing on the FIRE data sets to find the number of word commonly used and also avoid small ones. Which single DM technique is best is very difficult to say. Depending on concrete situations, sometime some techniques perform better than others.

- We also analyzed that there is no single classifier which produce best result for FIRE dataset. The classifier performance is evaluated using testing data set. But there is also some problem with testing data set. Some time it is easy and some time it becomes complex to classify the testing data set. To avoid these problems we used cross validation method so that every record of FIRE data set is used for both training and testing.

### REFERENCES:

1) Suman Bala, Krishan Kumar A Literature Review on "Kidney Disease Prediction using Data Mining Classification Technique" IJCSMC, Vol. 3, Issue. 7, July 2014
2) Anurag Upadhayay, Suneet Shukla, Sudsanshu Kumar " Empirical Comparison by data mining Classification algorithms (C 4.5 & C 5.0) for thyroid cancer data set." International Journal of Computer Science & Communication Networks,Vol 3(1)
3) H. C. Koh and G. Tan, "*Data Mining Application in Healthcare*", Journal of Healthcare Information Management, vol. 19, no. 2, 2005.
4) Brijesh Kumar Bhardwaj, Saurabh Pal "Data Mining: A prediction for performance improvement using classification" (IJCSIS) International Journal of Computer Science and Information Security, Vol. 9, No. 4, April 2011
5) AI-Radaideh,Q. A., AI-Shawakfa, E.M., and AI-Najjar, M. I.,"Mining Student Data using Decision Trees", International Arab Conference on Information Technology(ACIT'2006),Yarmouk University, Jordan, 2006.
6) Alaa el-Halees, "Mining Students Data to Analyze e-Learning Behavior: A Case Study", 2009.
7) Vishal Gupta, Gurpreet S. Lehal Professor & Head, Department of Computer Science, Punjabi University Patiala, India " A Survey of Text Mining Techniques and Applications" JOURNAL OF EMERGING TECHNOLOGIES IN WEB INTELLIGENCE, VOL 1, AUGUST 2009
8) [Berry Michael W., (2004), "Automatic Discovery of Similar Words", in "Survey of Text Mining: Clustering, Classification and Retrieval", Springer Verlag, New York, LLC, 24-43.
9) Navathe, Shamkant B., and Elmasri Ramez, (2000), "Data Warehousing And Data Mining", in "Fundamentals of Database Systems", Pearson Education pvt Inc, Singapore, 841-872.
10) S. M. Kamruzzaman, Farhana Haider, Ahmed Ryadh Hasan "Text Classification Using Data Mining" ICTM 2005
11) K.Sudhakar & Dr. M. Manimekalai," *Study of Heart Disease Prediction using Data mining*", IJARCSSE, Volume 4, Issue 1, January 2014.
12) Frawley and Piatetsky-Shapiro, 1996. *Knowledge Discovery in Databases*: An Overview. The AAAI/MIT Press, Menlo Park, C.A.
13) DSVGK Kaladhar, Krishna Apparao Rayavarapu* and Varahalarao Vadlapudi,"*Statistical and Data Mining Aspects on Kidney Stones: A Systematic Review and Meta-analysis*", Open Access Scientific Reports, Volume 1 • Issue 12 • 2012
14) AI-Radaideh,Q. A., AI-Shawakfa, E.M., and AI-Najjar, M. I.,"Mining Student Data using Decision Trees", International
15) Arab Conference on Information Technology(ACIT'2006),Yarmouk University, Jordan, 2006.
16) Alaa el-Halees, "Mining Students Data to Analyze e-Learning Behavior: A Case Study", 2009.
17) Agarwal R., Mannila H., Srikant R., Toivonan H., Verkamo, "A Fast Discovery of Association Rules," *Advances in Knowledge Discovery and Data Mining,* 1996.