

# Survey of Classification Techniques in Data Mining

V. Krishnaiah<sup>1\*</sup>, Dr.G.Narsimha<sup>2</sup> and Dr.N.Subhash Chandra<sup>3</sup>

<sup>1</sup>Department of CSE, CVR College of Engineering, JNTUH, TS, India

<sup>2</sup>Department of CSE, JNTUH-Kondagattu, JNTUH, TS, India

<sup>3</sup>Department CSE, VBIT, JNTUH, TS, India

[www.ijcaonline.org](http://www.ijcaonline.org)

Received: Aug/24/2014

Revised: Sep/08/2014

Accepted: Sep/20/2014

Published: Sep/30/2014

**Abstract**— Classification in data mining is a technique based on machine learning algorithms which uses mathematics, statistics, probability distributions and artificial intelligence... To predict group membership for data items or to represent descriptive analysis of data items for effective decision making .Now a day's data mining is touching every aspect of individual life includes Data Mining for Financial Data Analysis, Data Mining for the Telecommunications Industry Data Analysis, Data Mining for the Retail Industry Data Analysis, Data Mining in Healthcare and Biomedical Research Data Analysis, and Data Mining in Science and Engineering Data Analysis, etc. The goal of this survey is to provide a comprehensive review of different classification techniques in data mining based on decision tree, rule based Algorithms, neural networks, support vector machines, Bayesian networks, and Genetic Algorithms and Fuzzy logic.

**Keywords**— Classifiers, Data Mining Techniques, Intelligent Data Analysis, Learning Algorithms, Artificial Intelligence, Decision Support System, Data Mining, KDD, Classification Algorithms, Fuzzy sets, Genetic Algorithm,etc.

## I. INTRODUCTION TO CLASSIFICATION

Classification is a data mining function that assigns items in a collection to target Categories or classes. The goal of classification is to accurately predict the target class for each case in the data. For example, a classification model could be used to identify loan applicants as low, medium, or high credit risks. A classification task begins with a data set in which the class assignments are known. For example, a classification model that predicts credit risk could be developed based on observed data for many loan applicants over a period of time. In addition to the historical credit rating, the data might track employment history, home ownership oriental, years of residence, number and type of investments, and so on. Credit rating would be the target, the other attributes would be the predictors, and the data for each Customer would constitute a case.

Classifications are discrete and do not imply order. Continuous, floating-point values would indicate a numerical, rather than a categorical, target. A predictive model with a numerical target uses a regression algorithm, not a classification algorithm. The simplest type of classification problem is binary classification. In binary classification, the target attribute has only two possible values: for example, high credit rating or low credit rating. Multiclass targets have more than two values: for example, low, medium, high, or unknown credit rating. In the model build (training) process, a classification algorithm finds relationships between the values of the predictors and the values of the target. Different Classification algorithms use different techniques for finding

relationships. These relationships are summarized in a model, which can then be applied to a different dataset in which the class assignments are unknown. Classification models are tested by comparing the predicted values to known target values in a set of test data. The historical data for a classification project is typically divided into two data sets: one for building the model; the other for testing the model.

### A. Building Classification Model

Building a classification model in data mining is divided into two categories, viz., supervised learning and unsupervised learning. where unsupervised is finding the class label for a record in training data set for which class label is not predefined.

Describing a set of predetermined classes each tuple is assumed to belong to a predefined class, as determined by the class label attribute (supervised learning).The set of tuples used for model construction: training set. The model is represented as classification rules, decision trees, or mathematical formulae.

### B. Testing the classification Model

For classifying previously unseen objects Estimate accuracy of the model using a test set of attributes The known label of test sample is compared with the classified result from the model. Accuracy rate is the percentage of test set samples that are correctly classified by the model. Test set is independent of training set, otherwise over-fitting will occur. or not for only few records where under fitting will occur.

A classification model is tested by applying it to test data with known target values and comparing the predicted values with the known values. The test data must be compatible with

Corresponding Author: V.Krishnaiah

the data used to build the model and must be prepared in the same way that the build data was prepared. Typically the build data and test data come from the same historical data set. A percentage of the records is used to build the model; the remaining records are used to test the model. Test metrics are used to assess how accurately the model predicts the known values. If the model performs well and meets the business requirements, it can then be applied to new data to predict the future.

There are at least three techniques which are used to calculate a classifier's accuracy. One technique is to split the training set by using two-thirds for training and the other third for estimating performance. In another technique, known as cross-validation, the training set is divided into mutually exclusive and equal-sized subsets and for each subset the classifier is trained on the union of all the other subsets. The average of the error rate of each subset is therefore an estimate of the error rate of the classifier. Leave-one-out validation is a special case of cross Validation. All test subsets consist of a single instance. This type of validation is, of course, more expensive computationally, but useful when the most accurate Estimate of a classifier's error rate is required.

## II. GENERAL ISSUES OF CLASSIFICATION LEARNING ALGORITHMS

Inductive machine learning is the process of learning a set of rules from instances (examples in a training set), or more generally speaking, creating a classifier that can be used to generalize from new instances. The process of applying supervised ML to a real-world problem is described in Fig-1.

The first step is collecting the dataset. If a requisite expert is available, then s/he could suggest which fields (attributes, features) are the most informative. If not, then the simplest method is that of "brute-force," which means measuring everything available in the hope that the right (informative, relevant) features can be isolated. However, a dataset collected by the "brute-force" method is not directly suitable for induction. It contains in most cases noise and missing feature values, and therefore requires significant pre-processing [1].

The second step is the data preparation and data pre-processing. Depending on the circumstances, researchers have a number of methods to choose from to handle missing data [7] have recently introduced a survey of contemporary techniques for outlier (noise) detection. These researchers have identified the techniques' advantages and disadvantages.

Instance selection is not only used to handle noise but to cope with the infeasibility of learning from very large datasets. Instance selection in these datasets is an optimization problem that attempts to maintain the mining quality while minimizing the sample size. It reduces data and enables a data mining algorithm to function and work effectively with very large datasets. There are varieties of procedures for sampling instances from a large dataset.

Feature subset selection is the process of identifying and removing as many irrelevant and redundant features as possible this reduces the dimensionality of the data and enables data mining algorithms to operate faster and more effectively. The fact that many features depend on one another often unduly influences the accuracy of supervised ML classification models. This problem can be addressed by constructing new features from the basic feature set; this technique is called feature construction/transformation. These newly generated features may lead to the creation of more concise and accurate classifiers. In addition, the discovery of meaningful features contributes to better comprehensibility of the produced classifier, and a better understanding of the learned concept.

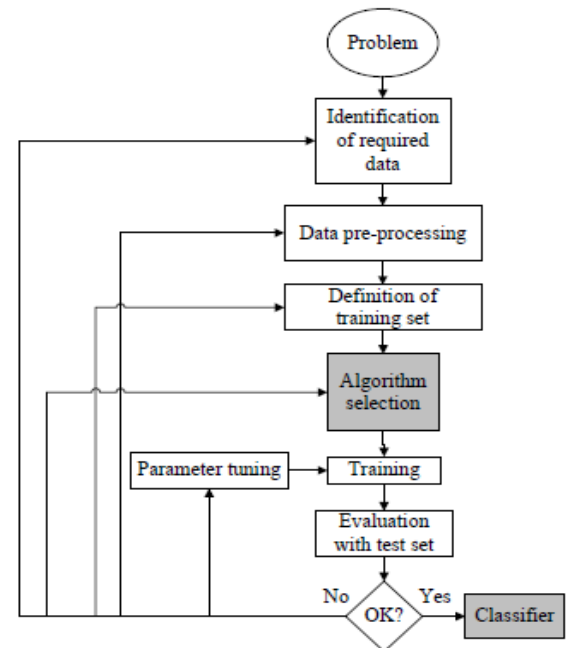


Figure 1. The process of supervised ML

The choice of which specific learning algorithm we should use is a critical step. Once preliminary testing is judged to be satisfactory, the classifier (mapping from unlabeled instances to classes) is available for routine use. The classifier's evaluation is most often based on prediction accuracy (the percentage of correct prediction divided by the total number of predictions). If the error rate evaluation is unsatisfactory, we must return to a previous stage of the supervised ML process (as detailed in Figure 1). A variety of factors must be examined: perhaps relevant features for the problem are not being used, a larger training set is needed, the dimensionality of the problem is too high, the selected algorithm is inappropriate or parameter tuning is needed. Another problem could be that the dataset is imbalanced.

A common method for comparing supervised ML algorithms is to perform statistical comparisons of the accuracies of trained classifiers on specific datasets. If we have sufficient supply of data, we can sample a number of training sets of size  $N$ , run the two learning algorithms on each of them, and estimate the difference in accuracy for each pair of classifiers on a large test set. The average of these

differences is an estimate of the expected difference in generalization error across all possible training sets of size  $N$ , and their variance is an estimate of the variance of the classifier in the total set. Our next step is to perform paired  $t$ -test to check the null hypothesis that the mean difference between the classifiers is zero. This test can produce two types of errors. Type I error is the probability that the test rejects the null hypothesis incorrectly (i.e. it finds a “significant” difference although there is none). Type II error is the probability that the null hypothesis is not rejected, when there actually is a difference. The test’s Type I error will be close to the chosen significance level.

In practice, however, we often have only one dataset of size  $N$  and all estimates must be obtained from this sole dataset. Different training sets are obtained by sub-sampling, and the instances not sampled for training are used for testing. Unfortunately this violates the independence assumption necessary for proper significance testing. The consequence of this is that Type I errors exceed the significance level. This is problematic because it is important for the researcher to be able to control Type I errors and know the probability of incorrectly rejecting the null hypothesis. Several heuristic versions of the  $t$ -test have been developed to alleviate this problem. here we explore some classification algorithms for supervised learning techniques.

### III. DECISION TREE INDUCTION BASED ALGORITHMS

Decision trees are trees that classify instances by sorting them based on feature values. Each node in a decision tree represents a feature in an instance to be classified, and each branch represents a value that the node can assume. Instances are classified starting at the root node and sorted based on their feature values. An example of a decision tree for the training set of Table I.

Table I. Training set

At1	At2	At3	At4	Class
a1	a2	a3	a4	Yes
a1	a2	a3	b4	Yes
a1	b2	a3	a4	Yes
a1	b2	b3	b4	No
a1	c2	a3	a4	Yes
a1	c2	a3	b4	No
b1	b2	b3	b4	No
C1	b2	b3	b4	No

Using the decision tree as an example, the instance  $At1 = a1, At2 = b2, At3 = a3, At4 = b4$  would sort to the nodes:  $At1, At2$ , and finally  $At3$ , which would classify the instance as being positive (represented by the values “Yes”). The problem of constructing optimal binary decision trees is an NP complete problem and thus theoreticians have searched for efficient *heuristics* for constructing near-optimal decision trees.

The feature that best divides the training data would be the root node of the tree. There are numerous methods for finding the feature that best divides the training data such as information gain, Relief algorithm estimates them in the context of other attributes. However, a majority of studies have concluded that there is no single best method. Comparison of individual methods may still be important when deciding which metric should be used in a particular dataset. The same procedure is then repeated on each partition of the divided data, creating sub-trees until the training data is divided into subsets of the same class.

The basic algorithm for decision tree induction is a greedy algorithm that constructs decision based on divide and conquers strategy. The algorithm, summarized as follows:

#### Pseudo Code: Decision Tree algorithm

1. create a node  $N$ ;
2. if samples are all of the same class,  $C$  then
3. return  $N$  as a leaf node labeled with the class  $C$ ;
4. if attribute-list is empty then
5. return  $N$  as a leaf node labeled with the most common class in samples;
6. select test-attribute, the attribute among attribute-list with the highest information gain;
7. label node  $N$  with test-attribute;
8. for each known value  $a_i$  of test-attribute
9. grow a branch from node  $N$  for the condition test-attribute =  $a_i$ ;
10. let  $s_i$  be the set of samples for which test-attribute =  $a_i$ ;
11. if  $s_i$  is empty then
12. attach a leaf labeled with the most common class in samples;
- 13.else attach the node returned by Generate\_decision\_tree(misattribute-list test-attribute).

**There are many specific decision-tree algorithms. Few of them are given below:**

- ID3 (Iterative Dichotomiser 3)
- C4.5 (successor of ID3)
- CART (Classification And Regression Tree)
- CHAID (Chi-squared Automatic Interaction Detector). Performs multi-level splits when computing classification trees
- MARS extends decision trees to handle numerical data better.

#### A. Iterative Dichotomiser 3(ID3) Algorithm

ID3 algorithm begins with the original set as the root node. On each iteration of the algorithm, it iterates through every unused attribute of the set and calculates the entropy (or information gain  $IG(A)$ ) of that attribute. Then selects the attribute which has the smallest entropy (or largest information gain) value. The set is  $S$  then split by the selected attribute (e.g.  $age < 50, 50 \leq age < 100, age \geq 100$ ) to produce subsets of the data. The algorithm continues to recurse on each subset, considering only attributes never selected before. Recursion on a subset may stop in one of these cases:

- ❖ Every element in the subset belongs to the same class (+ or -), then the node is turned into a leaf and labeled with the class of the examples

- ❖ There are no more attributes to be selected, but the examples still do not belong to the same class (some are + and some are -), then the node is turned into a leaf and labeled with the most common class of the examples in the subset
- ❖ There are no examples in the subset, this happens when no example in the parent set was found to be matching a specific value of the selected attribute, for example if there was no example with age  $\geq 100$ . Then a leaf is created, and labeled with the most common class of the examples in the parent set.

Throughout the algorithm, the decision tree is constructed with each non-terminal node representing the selected attribute on which the data was split, and terminal nodes representing the class label of the final subset of this branch.

**Pseudo code: ID3 (Examples, Target Attribute, and Attributes):**

1. Create a root node for the tree

If all examples are positive, Return the single-node tree Root, with label = +. If all examples are negative, Return the single-node tree Root, with label = -.

2. If number of predicting attributes is empty, then Return the single node tree Root, with label = most common value of the target attribute in the examples. Otherwise Begin

A ← The Attribute that best classifies examples.

Decision Tree attribute for Root = A. For each possible value,  $V_i$  of A,

3. Add a new tree branch below Root, corresponding to the test  $A = V_i$ .

Let Examples ( $V_i$ ) be the subset of examples that have the value  $V_i$  for A

If Examples ( $V_i$ ) is empty

4. Then below this new branch add a leaf node with label = most common target value in the examples Else below this new branch add the subtree ID3 (Examples ( $V_i$ ), Target Attribute, Attributes – {A}) end.

#### B. C4.5 Algorithm

C4.5 is an algorithm used to generate a decision tree developed by Ross Quinlan. C4.5 is an extension of Quinlan's earlier ID3 algorithm. The decision trees generated by C4.5 can be used for classification, and for this reason, C4.5 is often referred to as a statistical classifier. One limitation of ID3 is that it is overly sensitive to features with large numbers of values. This must be overcome if you are going to use ID3 as an Internet search agent. I address this difficulty by borrowing from the C4.5 algorithm, an ID3 extension. ID3's sensitivity to features with large numbers of values is illustrated by Social Security numbers. Since Social

Security numbers are unique for every individual, testing on its value will always yield low conditional entropy values. However, this is not a useful test. To overcome this problem, C4.5 uses a metric called "information gain," which is defined by subtracting conditional entropy from the base entropy; that is,  $\text{Gain}(PIX) = E(P) - E(PIX)$ . This computation does not, in itself, produce anything new. However, it allows you to measure a gain ratio. Gain ratio, defined as  $\text{Gain Ratio}(PIX) = \text{Gain}(PIX) / E(X)$ , where  $E(X)$  is the entropy of the examples relative only to the attribute. It has an enhanced method of tree pruning that reduces misclassification errors due noise or too-much details in the training data set. Like ID3 the data is sorted at every node of the tree in order to determine the best splitting attribute. It uses gain ratio impurity method to evaluate the splitting attribute. Decision trees are built in C4.5 by using a set of training data or data sets as in ID3. At each node of the tree, C4.5 chooses one attribute of the data that most effectively splits its set of samples into subsets enriched in one class or the other. Its criterion is the normalized information gain (difference in entropy) that results from choosing an attribute for splitting the data. The attribute with the highest normalized information gain is chosen to make the decision.

#### Pseudo Code:

1. Check for base cases.
2. For each attribute a calculate:
  - i. Normalized information gain from splitting on attribute
3. Select the best a, attribute that has highest information gain.
4. Create a decision node that splits on best of a, as root node.
5. Recurs on the sub lists obtained by splitting on best of a and add those nodes as children node.

#### C CART Algorithm

CART (Classification and Regression Trees) is very similar to C4.5, but it differs in that it supports numerical target variables (regression) and does not compute rule sets. CART constructs binary trees using the feature and threshold that yield the largest information gain at each node. It is a Binary decision tree algorithm Recursively partitions data into 2 subsets so that cases within each subset are more homogeneous Allows consideration of misclassification costs, prior distributions, cost-complexity pruning.

#### Pseudo Code:

1. The basic idea is to choose a split at each node so that the data in each
  - Subset (child node) is "purer" than the data in the parent node. CART
  - Measures the impurity of the data in the nodes of a split with an impurity measure  $i(t)$ .
2. If a split  $s$  at node  $t$  sends a proportion  $p_L$  of data to its left child node  $t_L$  and a corresponding proportion  $p_R$  of data

to its right child node  $tR$ , the decrease in impurity of split  $s$  at node  $t$  is defined as

$$\Delta i(s,t) = i(t) - pL_i(tL) - pR_i(tR)$$

= impurity in node  $t$  – weighted average of impurities in nodes  $tL$  and  $tR$

3. A CART tree is grown, starting from its root node (i.e., the entire training data set)  $t=1$ , by searching for a split  $s^*$  among the set of all possible candidates  $S$  which give the largest decrease in impurity.

4. The above split searching process is repeated for each child node.

5. The tree growing process is stopped when all the stopping criteria are met.

#### D. CHAID

CHAID is a type of decision tree technique, based upon adjusted significance testing (Bonferroni testing). The technique was developed in South Africa and was published in 1980 by Gordon V. Kass, who had completed a PhD thesis on this topic. CHAID can be used for prediction (in a similar fashion to regression analysis, this version of CHAID being originally known as XAID) as well as classification, and for detection of interaction between variables. CHAID stands for CHi-squared Automatic Interaction Detection, based upon a formal extension of the US AID (Automatic Interaction Detection) and THAID (THeta Automatic Interaction Detection) procedures of the 1960s and 70s, which in turn were extensions of earlier research, including that performed in the UK in the 1950s.

In practice, CHAID is often used in the context of direct marketing to select groups of consumers and predict how their responses to some variables affect other variables, although other early applications were in the field of medical and psychiatric research.

Like other decision trees, CHAID's advantages are that its output is highly visual and easy to interpret. Because it uses multiday splits by default, it needs rather large sample sizes to work effectively, since with small sample sizes the respondent groups can quickly become too small for reliable analysis. One important advantage of CHAID over alternatives such as multiple regressions is that it is non-parametric.

#### E. MARS (Multivariate adaptive regression splines)

Multivariate adaptive regression splines (MARS) is a form of regression analysis introduced by Jerome H. Friedman in 1991. It is a non-parametric regression technique and can be seen as an extension of linear cubic spline models that automatically models non-linearity's and interactions between variable[10].

#### Decision tree Advantages:

- Simple to understand and interpret. People are able to understand decision tree models after a brief explanation.
- Requires little data preparation. Other techniques often require data normalization, dummy variables need to be created and blank values to be removed.
- Able to handle both numerical and categorical data. Other techniques are usually specialized in analyzing datasets that have only one type of variable. (For example, relation rules can be used only with nominal variables while neural networks can be used only with numerical variables.)
- Uses a white box model. If a given situation is observable in a model the explanation for the condition is easily explained by Boolean logic. (An example of a black box model is an artificial neural network since the explanation for the results is difficult to understand.)
- Possible to validate a model using statistical tests. That makes it possible to account for the reliability of the model.
- Robust. Performs well even if its assumptions are somewhat violated by the true model from which the data were generated.
- Performs well with large datasets. Large amounts of data can be analyzed using standard computing resources in reasonable time.

#### Limitations:

1. The problem of learning an optimal decision tree is known to be NP-complete under several aspects of optimality and even for simple concepts. Consequently, practical decision-tree learning algorithms are based on heuristics such as the greedy algorithm where locally-optimal decisions are made at each node. Such algorithms cannot guarantee to

2. return the globally-optimal decision tree. To reduce the greedy effect of local-optimality some methods such as the dual information distance (DID) tree were proposed.

3. Decision-tree learners can create over-complex trees that do not generalize well from the training data. (This is known as over fitting. Mechanisms such as pruning are necessary to avoid this problem (with the exception of some algorithms such as the Conditional Inference approach that does not require pruning)

4. There are concepts that are hard to learn because decision trees do not express them easily, such as XOR, parity or multiplexer problems. In such cases, the decision tree becomes prohibitively large. Approaches to solve the problem involve either changing the representation of the problem domain (known as propositionalisation) or using learning algorithms based on more expressive representations (such as statistical relational learning or inductive logic programming).

5. For data including categorical variables with different numbers of levels, information gain in decision trees is biased in favor of those attributes with more levels. However, the issue of biased predictor selection is avoided by the Conditional Inference approach.

#### IV. NAIVE BAYES ALGORITHM

The Naive Bayes classifier is a simple probabilistic classifier based on applying Bayes' theorem with strong (naive) independence assumptions. A more descriptive term for the underlying probability model would be "independent feature model". Bayesian classification provides practical learning algorithms and prior knowledge and observed data can be combined. Bayesian Classification provides a useful perspective for understanding and evaluating many learning algorithms. It calculates explicit probabilities for hypothesis and it is robust to noise in input data. A naive Bayes classifier assumes that the presence or absence of a particular feature is unrelated to the presence or absence of any other feature, given the class variable. For some types of probability models, naive Bayes classifiers can be trained very efficiently in a supervised learning setting. It also called idiot's Bayes, simple Bayes, and independence Bayes. This method is important for several reasons. It is very easy to construct, not needing any complicated iterative parameter estimation schemes. This means it may be readily applied to huge data sets. It is easy to interpret, so users unskilled in classifier technology can understand why it is making the classification it makes. And finally, it often does surprisingly well: it may not Probabilistic approaches to classification typically involve modeling the conditional probability distribution  $P(C|D)$ , where  $C$  ranges over classes and  $D$  over descriptions, in some language, of objects to be classified. Given a description  $d$  of a particular object, we assign the class  $\operatorname{argmax}_c P(C = c|D = d)$ . A Bayesian approach splits this posterior distribution into a prior distribution  $P(C)$  and a likelihood  $P(D|C):P(D = d|C = c)P(C = c)$

$$\operatorname{argmax}_c P(C = c|D = d) = \operatorname{argmax}_c P(C = c) \prod_{i=1}^n P(D_i = d_i | C = c) \quad (1)$$

The denominator  $P(D = d)$  is a normalizing factor that can be ignored when determining the maximum a posteriori class, as it does not depend on the class. The key term in Equation (1) is  $P(D = d|C = c)$ , the likelihood of the given description given the class (often abbreviated to  $P(d|c)$ ). A Bayesian classifier estimates these likelihoods from training data, but this typically requires some additional simplifying assumptions. For instance, in an attribute-value representation (also called propositional or single-table representation), the individual is described by a vector of values  $a_1, \dots, a_n$  for a fixed set of attributes  $A_1, \dots, A_n$ . Determining  $P(D = d|C = c)$  here requires an estimate of the joint probability  $P(A_1 = a_1, \dots, A_n = a_n|C = c)$ , abbreviated to  $P(a_1, \dots, a_n|c)$ . This joint probability Distribution is problematic for two reasons:

- (1) Its size is exponential in the number of attributes  $n$ , and

- (2) It requires a complete training set, with several examples for each possible description. These problems vanish if we can assume that all attributes are independent

Given the class:

$$P(A_1 = a_1, \dots, A_n = a_n|C = c) = \prod_{i=1}^n P(A_i = a_i|C = c)$$

This assumption is usually called the Naive Bayes assumption, and a Bayesian classifier using this assumption is called the naive Bayesian classifier, often abbreviated to Naive Bayes". Effectively, it means that we are ignoring interactions between attributes within individuals of the same class.

#### V. SUPPORT VECTOR MACHINES (SVM)

SVMs introduced in COLT-92 by Boser, Guyon & Vapnik. Theoretically well motivated algorithm developed from Statistical Learning Theory (Vapnik & Chervonenkis) since the 60s. The support vector machine usually deals with pattern classification that means this algorithm is used mostly for classifying the different types of patterns. Now, there is different type of patterns i.e. Linear and non-linear. Linear patterns are patterns that are easily distinguishable or can be easily separated in low dimension whereas non-linear patterns are patterns that are not easily distinguishable or cannot be easily separated and hence these type of patterns need to be further manipulated so that they can be easily separated. Basically, the main idea behind SVM is the construction of an optimal hyper plane, which can be used for classification, for linearly separable patterns. The optimal hyper plane is a hyper plane selected from the set of hyper planes for classifying patterns that maximizes the margin of the hyper plane i.e. the distance from the hyper plane to the nearest point of each patterns. The main objective of SVM is to maximize the margin so that it can correctly classify the given patterns i.e. larger the margin size more correctly it classifies the patterns. The equation shown below is the hyper plane:

$$\text{Hyper plane, } aX + bY = C$$

The given pattern can be mapped into higher dimension space using kernel function,  $\Phi(x)$ . i.e.  $x \rightarrow \Phi(x)$  selecting different kernel function is an important aspect in the SVM-based classification, commonly used kernel functions include LINEAR, POLY, RBF, and SIGMOID. For e.g.: the equation for Poly Kernel function is given as:

$$K(x, y) = \langle x, y \rangle^p$$

The main principle of support vector machine is that given a set of independent and identically distributed training sample  $\{(x_i, y_i)\}_{i=1}^N$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$ , denote the input and output of the classification.

The goal is to find a hyper plane  $w^T x + b = 0$ , which separate the two different samples accurately. Therefore, the problem of solving optimal classification now translates into solving quadratic programming problems. It is to seek a partition hyper plane to make the bilateral blank area  $(2/\|w\|)$  maximum, which means we have to maximize the weight of the margin. It is expressed as:

$$\text{Min } \Phi(w) = \frac{1}{2} \|w\|^2 = \frac{1}{2} (w, w),$$

$$\text{Such that: } y_i (w \cdot x_i + b) \geq 1$$

SVM can be easily extended to perform numerical calculations. Here we discuss two such extensions. To extend SVM to perform regression analysis, where the goal is to produce a linear function that can approximate that target function.

## VI. K-NEARESTNEIGHBOUR (KNN) ALGORITHM

The nearest neighbor (NN) rule identifies the category of unknown data point on the basis of its nearest neighbor whose class is already known. This rule is widely used in pattern recognition text categorization ranking models object recognition and event recognition applications. M. Cover and P. E. Hart purpose k-nearest neighbour (KNN) in which nearest neighbor is calculated on the basis of value of k that specifies how many nearest neighbors are to be considered to define class of a sample data point. It makes use of the more than one nearest neighbour to determine the class in which the given data point belongs to and hence it is called as K-NN. These data samples are needed to be in the memory at the run time and hence they are referred to as memory-based technique. T. Bailey and A. K. Jain improve KNN which is based on weights The training points are assigned weights according to their distances from sample data point. But still, the computational complexity and memory requirements remain the main concern always. To overcome memory limitation, size of data set is reduced. For this, the repeated patterns, which do not add extra information, are eliminated from training samples. To further improve, the data points which do not affect the result are also eliminated from training data set. The NN training data set can be structured using various techniques to improve over memory limitation of KNN. The KNN implementation can be done using ball tree, k-d tree, nearest feature line (NFL), tunable metric, principal axis search tree and orthogonal search tree. The tree structured training data is divided into nodes, whereas techniques like NFL and tunable metric divide the training data set according to planes. These algorithms increase the speed of basic KNN algorithm. Suppose that an object is sampled with a set of different attributes, but the group to which the object belongs is unknown. Assuming its group can be determined from its attributes.

Different algorithms can be used to automate the classification process. With the k-nearest neighbor technique, this is done by evaluating the k number of closest neighbors In pseudo code, k-nearest neighbor classification algorithm can be expressed,

```

K ← number of nearest neighbors
For each object X in the test set
do

```

```

Calculate the distance D(X,Y) between X and every
object Y in the training set neighborhood ← the k
neighbors in the training set closest to X
X.class ← Select Class (neighborhood)

```

**End for**

The k-nearest neighbors" algorithm is the simplest of all machine learning algorithms. It has got a wide variety of applications in various fields such as Pattern recognition, Image databases, Internet marketing, Cluster analysis etc In binary (two class) classification problems, it is helpful to choose k to be an odd number as this avoids tied votes. Here a single number k" is given which is used to determine the total number of neighbors that determines the classification. If the value of k=1, then it is simply called as nearest neighbor.. K-NN requires an integer k, a training data set and a metric to measure closeness.

## VII. RULE BASED CLASSIFIER

Rule-based classifier makes use of set of IF-THEN rules for classification. We can express the rule in the following from:

*IF condition THEN conclusion*

*Let us consider a rule R1*

*R1: IF age=youth AND student=yes THEN buy computer=yes*

The IF part of the rule is called rule antecedent or pre-condition. The THEN part of the rule is called rule consequent. In the antecedent part the condition consists of one or more attribute tests and these tests are logically ANDed. The consequent part consist class prediction.

*Note: We can also write rule R1 as follows: R1: (age = youth) ^ (student = yes)(buys computer = yes)*

If the condition holds the true for a given tuple, then the antecedent is satisfied. Here we will learn how to build a rule based classifier by extracting IF-THEN rules from decision tree. Points to remember to extract rule from a decision tree, one rule is created for each path from the root to the leaf node. To from the rule antecedent each splitting criterion is logically ANDed. The leaf node holds the class prediction, forming the rule consequent.

### A. Rule Induction Using Sequential Covering Algorithm

Sequential Covering Algorithm can be used to extract IF-THEN rules form the training data. We do not require generating a decision tree first. In this algorithm each rule for a given class covers many of the tuples of that class. Some of the sequential Covering Algorithms are AQ, CN2, and RIPPER. As per the general strategy the rules are

learned one at a time. For each time rules are learned, a tuple covered by the rule is removed and the process continues for rest of the tuples. This is because the path to each leaf in a decision tree corresponds to a rule. The Decision tree induction can be considered as learning a set of rules simultaneously. The Following is the sequential learning Algorithm where rules are learned for one class at a time. When learning a rule from a class  $C_i$ , we want the rule to cover all the tuples from class  $C$  only and no tuple form any other class.

#### B. Pseudo code: Sequential Covering

Input:

D, a data set class-labeled tuples,

Att\_vals, the set of all attributes and their possible values.

Output: A Set of IF-THEN rules.

Method:

```
Rule_set={ }; // initial set of rules learned is empty
  for each class c do
repeat
Rule = Learn_One_Rule(D, Att_vals, c);
remove tuples covered by Rule form D;
until termination condition;
Rule_set=Rule_set+Rule; // add a new rule to rule-set
end for
return Rule_Set;
```

#### C. Rule Pruning:

The rule is pruned is due to the following reason:

The assessments of quality are made on the original set of training data. The rule may perform well on training data but less well on subsequent data. That's why the rule pruning is required. The rule is pruned by removing conjunct. The rule R is pruned, if pruned version of R has greater quality than what was assessed on an independent set of tuples. FOIL is one of the simple and effective methods for rule pruning. For a given rule R, FOIL\_Prune = pos-neg/ posing Where pos and neg is the number of positive tuples covered by R, respectively. This value will increase with the accuracy of R on pruning set. Hence, if the FOIL Prune value is higher for the pruned version of R, then we prune R.

## VIII. GENETIC ALGORITHMS

The idea of Genetic Algorithm is derived from natural evolution. In Genetic Algorithm first of all initial population is created. This initial population consists of randomly

generated rules. We can represent each rule by a string of bits.

For example, suppose that in a given training set the samples are described by two Boolean attributes such as A1 and A2. And this given training set contains two classes such as C1 and C2.

We can encode the rule *IF A1 AND NOT A2 THEN C2* into bit string **100**. In this bit representation the two leftmost bit represent the attribute A1 and A2, respectively. Likewise the rule *IF NOT A1 AND NOT A2 THEN C1* can be encoded as **001**. If the attribute has K values where  $K > 2$ , then we can use the K bits to encode the attribute values. The classes are also encoded in the same manner.

Based on the notion of survival of the fittest, a new population is formed to consist of the fittest rules in the current population and offspring values of these rules a well. The fitness of the rule is assessed by its classification accuracy on a set of training samples. The genetic operators such as crossover and mutation are applied to create offspring. In crossover the substring from pair of rules are swapped to form a new pair of rules. In mutation, randomly selected bits in a rule's string are inverted.

#### A. Rough Set Approach

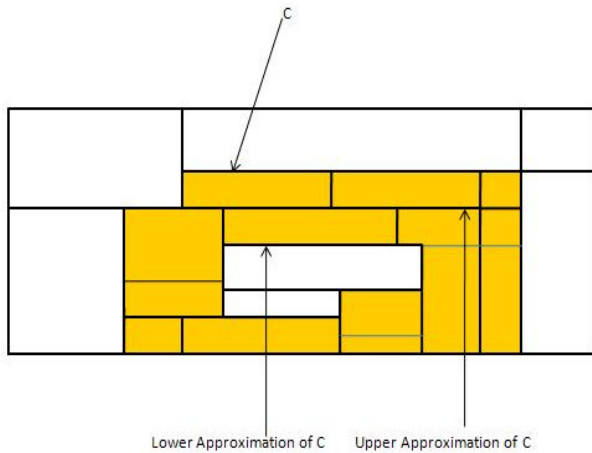
To discover structural relationship within imprecise and noisy data we can use the rough set. This approach can only be applied on discrete-valued attributes. Therefore, continuous-valued attributes must be discredited before its use. The Rough Set Theory is base on establishment of equivalence classes within the given training data. The tuples that forms the equivalence class are indiscernible. It means the samples are identical wrt to the attributes describing the data. There are some classes in given real world data, which cannot be distinguished in terms of available attributes. We can use the rough sets to **roughly** define such classes. For a given class, C, the rough set definition is approximated by two sets as follows:

*Lower Approximation of C:* The lower approximation of C consist of all the data tuples that Bases on knowledge of attribute. These attribute are certain to belong to class C.

*Upper Approximation of C:* The upper approximation of C consists of all the tuples that based on knowledge of attributes, cannot be described as not belonging to C.

The following diagram shows the Upper and Lower Approximation of class C:





**B. Fuzzy Set Approaches:**

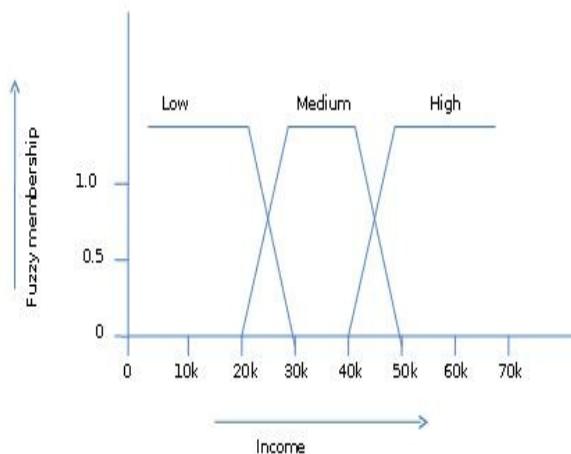
Fuzzy Set Theory is also called Possibility Theory. This theory was proposed by Lotfi Zadeh in 1965. This approach is an alternative **Two-value logic**. This theory allows us to work at high level of abstraction. This theory also provides us means for dealing with imprecise measurement of data.

The fuzzy set theory also allows dealing with vague or inexact facts. For example being a member of a set of high incomes is inexact (eg. if \$50,000 is high then what about \$49,000 and \$48,000). Unlike the traditional CRISP set where the element either belong to S or its complement but in fuzzy set theory the element can belong to more than one fuzzy set.

For example, the income value \$49,000 belong to both the medium and high fuzzy sets but to differing degrees. Fuzzy set notation for this income value is as follows:

$$m_{\text{medium\_income}}(\$49k)=0.15 \text{ and } h_{\text{high\_income}}(\$49k)=0.96$$

Where *m* is membership function that operates on fuzzy set of medium income and high-income respectively. This notation can be shown diagrammatically as follows:



Nowadays, let us deal by “fuzzy logic” in medicine in wide intellect. In the medicine, particularly, in oriental

medicine, nearly all medical concepts are fuzzy. The unfocused nature of medical concepts and their associations requires the use of “fuzzy logic”. It defines incorrect medical entities as fuzzy sets and provides a linguistic approach with an brilliant estimate to texts. “Fuzzy logic” presents analysis methods able of sketch estimated inferences.

*Table 2: Advantages and Disadvantages of Classification Algorithms*

S.NO	ALGORITHM MS	ADVANTAGES	DISADVANTAGES
1	C4.5 Algorithm	1. It produces the accurate result. 2. It takes the less memory to large program execution. 3. It takes less model build time. 4. It has short searching time.	1. Empty branches. 2. Insignificant branches. 3. Over fitting.
2	ID3	1. It produces the more accuracy result than the C4.5 algorithm. 2. ID3 algorithm generally uses nominal attributes for classification with no missing values. 3. It produces false alarm rate and omission rate decreased, increasing the detection rate and reducing the space Consumption	1. It has long searching time. 2. It takes the more memory than the C4.5 to large program execution.
3	Naive bayes Algorithm	1. To improves the classification performance by removing the irrelevant features. 2. Good performance. 3. It is short computational time	1. The naive Bayes classifier requires a very large number of records to obtain good results. 2. it is instance-based or lazy in that they store all of the training samples
4	Support vector machine Algorithm	1. Produce very accurate classifiers. 2. Less over fitting, robust to noise. 3. Especially popular in text classification problems where very high-dimensional spaces are the norm. 4. Memory-intensive.	1. SVM is a binary classifier. To do a multi-class classification, pair-wise classifications can be used (one class against all others, for all classes). 2. Computationally expensive, thus runs slow.
5	k-nearest neighbour Algorithm	1. It is an easy to understand and easy to implement classification technique. 2. Training is very fast. 3. Robust to noisy training data. 4. It is particularly well suited for multi-modal classes	1. It is sensitive to the local structure of the data. 2. Memory limitation. 3. Being a supervised learning lazy Algorithm i.e., runs slowly.

**IX. CONCLUSION**

This paper deals with various classification techniques used in data mining and a study on each of them. Data mining is a wide area that integrates techniques from various fields including machine learning, artificial intelligence, statistics and pattern recognition, for the analysis of large volumes of data. These classification algorithms can be implemented on

different types of data sets like data of patients, customers in banking, customers in telecom industry, customers in life, health insurance industry, customers in sales industry students in education, online social networks, web related files, audio files, voice files ,image files and video files and text files etc. Hence these classification techniques show how a data can be determined and grouped when a new set of data is available. Each technique has got its own pros and cons as given in the paper. Based on the needed Conditions each one as needed can be selected On the basis of the performance of these algorithms, these algorithms can also be used to detect and predict the natural disasters like cloud bursting, earth quake, tsunami, etc.

#### REFERENCES

- [1] B. Kotsiantis · I. D. Zaharakis · P. E. Pintelas, “Machine learning: a review of classification and combining techniques”, Springer Science 10 November 2007.
- [2] Raj Kumar, Dr. Rajesh Verma,” Classification Algorithms for Data Mining P: A Survey” IJJET Vol. 1 Issue August 2012, ISSN: 2319 – 1058.
- [3] Ms. Aparna Raj, Mrs. Bincy, Mrs. T.Mathu “Survey on Common Data Mining Classification Techniques”, International Journal of Wisdom Based Computing, Vol. 2(1), April 2012.
- [4] [http://www.tutorialspoint.com/data\\_mining/dm\\_rbc.htm](http://www.tutorialspoint.com/data_mining/dm_rbc.htm)
- [5] M.Archana Survey of Classification Techniques in Data Mining International Journal of Computer Science and Mobile Applications Vol.2 Issue. 2, February- 2014 **ISSN: 2321-8363.**
- [6] Thair Nu Phyu Survey of Classification Techniques in Data Mining Proceedings of the International Multi Conference of Engineers and Computer Scientists 2009, Vol I IMECS 2009, March 18 - 20, 2009, Hong Kong.
- [7] Jensen, F. (1996). An Introduction to Bayesian Networks. Springer.
- [8] Cover, T., Hart, P. (1967), Nearest neighbor pattern classification.
- [9] [http://docs.oracle.com/cd/B28359\\_01/datamine.111/b28129.pdf](http://docs.oracle.com/cd/B28359_01/datamine.111/b28129.pdf)
- [10] Text Book: Introduction to data mining by michael steinback, Pang-Ning Tan, Vipin Kumar.
- [11] [http://en.wikipedia.org/wiki/Decision\\_tree\\_learning.](http://en.wikipedia.org/wiki/Decision_tree_learning)