

## **Implementation and evaluation Of EAODT (Enhanced Algorithm of Decision Tree) on large data set**

**Rajdeep Kaur**

*M.tech Scholar CSE Dept, SBBSIET, Padhiana, India.*

*Jalandhar, India*

*E-mail: rajdeepatwal@yahoo.com*

**Harpreet Kaur**

*Asst Professor (CSE), SBBSIET, Punjab, India.*

*E-mail: er.harpreetarora@gmail.com*

### **Abstract**

Data mining is a process of identification of useful information from large amount of random data. It is used to discover meaningful pattern and rules from data. Classification, clustering, association rules are data mining techniques. Classification is a process of assigning entities to already defined class by examining the features. Decision tree is a classification technique in which a model is created that anticipates the value of target variable depends on input values. CART and C4.5 are commonly used decision tree algorithms. These algorithms are based on Hunt's algorithm. Goal of this study is to provide review of these decision tree algorithms. At first we present concept of Data Mining, Classification and Decision Tree. Then we present proposed algorithm EAODT, CART and C4.5 algorithms and we will make comparison of these two algorithms.

## **1. INTRODUCTION**

Data mining is a process of extraction useful information from large amount of data. It is used to discover meaningful pattern and rules from data. Data mining is a part of wider process called knowledge discovery [4]. The steps of knowledge discovery are

- Selection
- Processing
- Transformation
- Data mining
- Interpretation/Evaluation

Data mining uses two types of approaches i.e supervised learning or unsupervised learning.

## **CLASSIFICATION**

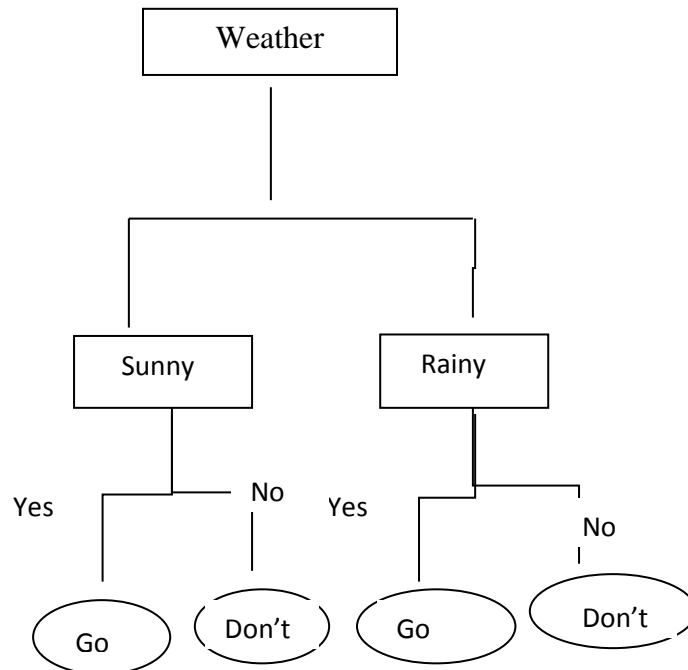
Classification is the process of assigning newly presented entities to already defined class by examining the features of entities. Classification is to make decision from unseen cases by building examples of past decisions [2]. There are two steps in classification process.

- In first step, model is built from training data in which value of class label is known. Classification algorithms are used to create model from training data sets.
- In second step, accuracy of model is checked by test data and if correctness of model is satisfactory then the model is used to classify data with unknown class label.

Among classification algorithm, decision tree algorithms is usually used because it is easy to follow and economical to implement.

## **DECISION TREES:**

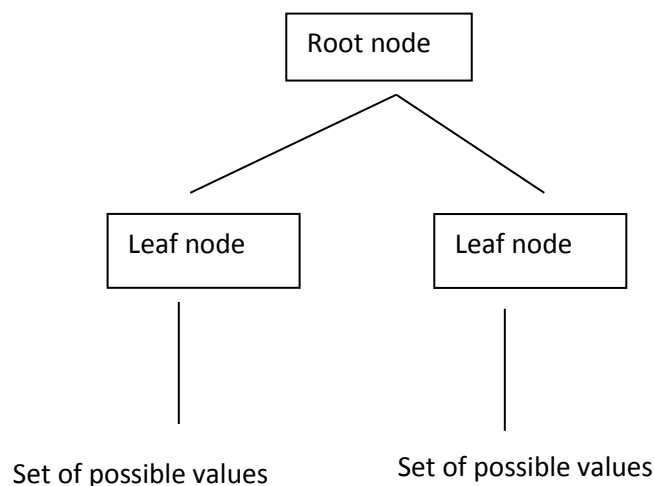
Decision tree is a classification technique. It is a tree like structure where internal node contains splits and splitting attributes. It represents test on an attribute. Arcs between internal node and its child contain consequences of test. Each leaf node is associated with a class label. Decision tree is constructed from training set. Then this decision tree is used to classify the tuples with unknown class label [2].



**Figure 1.** Decision Tree showing whether to go for trip or not depending on weather

**2. DECISION TREE ALGORITHMS**

Decision tree learning methods are most commonly used in data mining. The goal is create a model to predict value of target variable based on input values. Training dataset is used to create tree and test dataset is used to test accuracy of the decision tree. Each leaf node represents the target attribute’s value depend on input variables represented by path by path from root to leaf node. First, an attribute that splits data efficiently is selected as root node in order to create small tree. The attribute with higher information is selected as splitting attribute[4].



**Figure 2.** Decision tree induction

Decision tree algorithm involves three steps:

1. For a given dataset S, select an attribute as target class to split tuples in partitions.
2. Determine a splitting criterion to generate a partition in which all tuples belong to a single class. Choose best split to create a node.
3. Iteratively repeat above steps until complete tree is grown or any stopping criterion is fulfilled.
4. **CART:** CART algorithm is presented by J.R. Quinlan, 1986. CART uses Information gain as splitting criterion. Topmost decision node is the best predictor, it is called root node. The attribute with highest Information Gain is selected as split attribute. Information gain is used to create tree from training instances. This tree is used to classify test data. When information gain approaches to zero or all instances belong to single target then growing of tree stops. [1].

It grows tree classifiers in three steps:

1. Selection of target attribute and calculation of entropy of attributes.
2. Select attribute with highest information gain measure
3. Create node containing that attribute. Iteratively apply these steps to new tree branches and stop growing tree after checking of stop criterion.

The CART decision makes use of two concepts when creating a tree from top-down [1]:

1. Entropy
2. Information Gain (as referred to as just gain) Using these two concepts, the nodes to be created and the attributes to split on can be determined.

### Entropy

Entropy is degree of randomness of data. It is used to calculate homogeneity of data attribute. If entropy is zero then sample is totally homogeneous and if is one then sample is completely uncertain.

### Information Gain

Information gain is decrease in entropy. Attribute with highest information gain is selected as best splitting criterion attribute

$$ET(X, S) = \sum_{j=1}^k \frac{|S_j|}{|S|} \cdot ET(S_j)$$

$$IG(X, S) = E(S) - E(X, S)$$

### C4.5

C4.5 algorithm is enhancement to ID3. C4.5 can handle continuous input attribute..It follows three steps during tree growth [3]:

1. Splitting of categorical attribute is same to ID3 algorithm. Continuous attributes always generate binary splits.
2. Attribute with highest gain ratio is selected.
3. Iteratively apply these steps to new tree branches and stop growing tree after checking of stop criterion. Information gain bias the attribute with more number of values. C4.5 used a new selection criterion which is Gain ratio which is less biased.

The Gain ratio measure is a selection criterion which is used less biased towards selecting attributes with more number of values [3].

$$GR(X, S) = \frac{IG(X, S)}{SI(X, S)}$$

$$SI(X, S) = - \sum_{j=1}^k \frac{|S_j|}{|S|} \log \frac{|S_j|}{|S|}$$

### Advantages: C4.5 made improvements to ID3 [10]:

1. It can handle both discrete and numerical attributes.
2. It can handle missing value attribute.
3. It can avoid over fitting of decision tree by providing the facility of pre and post pruning.

### IMPLEMENTATION OF EAODT:

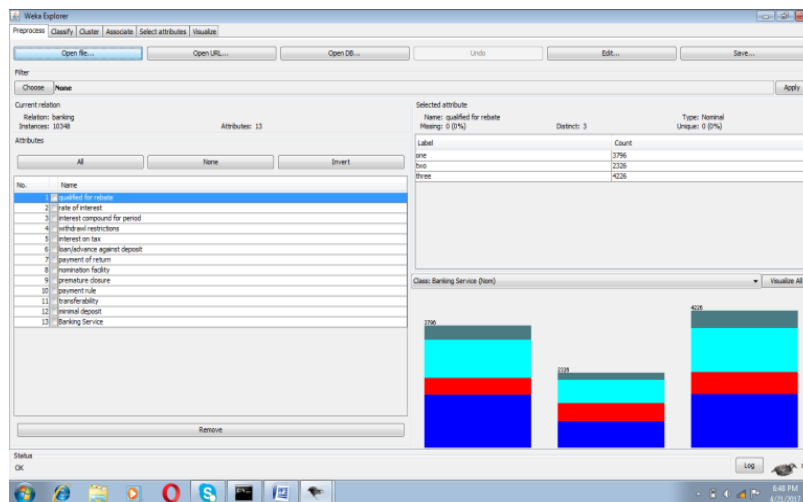


Figure 2: Description of First attribute of Dataset

Figure 2 gives the description of the first attribute (Qualified for rebate) in the dataset. It specifies that the attribute is numeric along with its minimum value, maximum value, Mean of the attribute and the standard deviation. The bar graph shows the distribution of four banking services in the first attribute

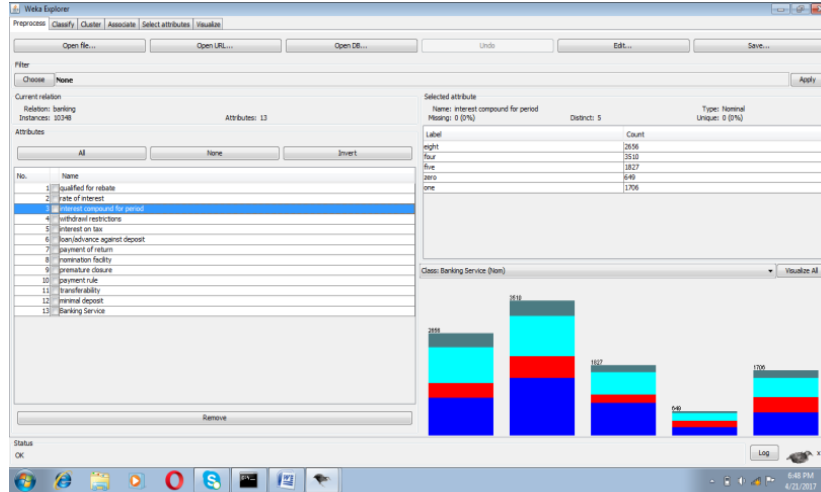


Figure 3: Description of 3<sup>rd</sup> attribute of dataset

Figure 3 shows that the attribute is numeric. So its min, max ,mean and standard deviation values are depicted in the figure. Graph shows the distribution of banking services in the attribute.

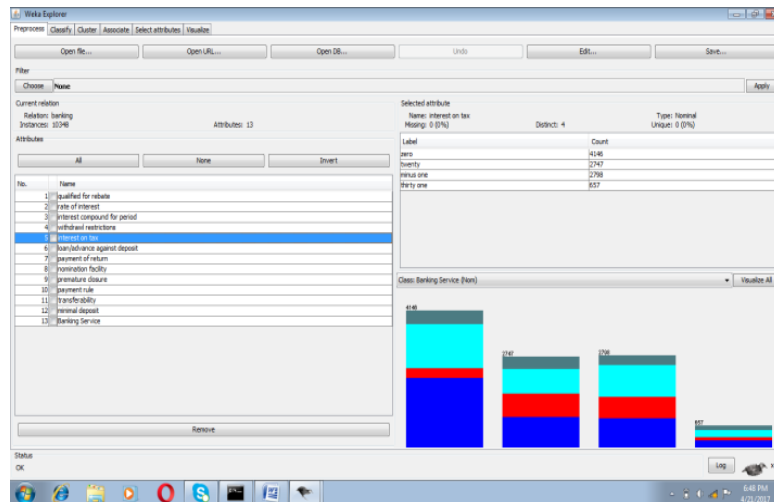
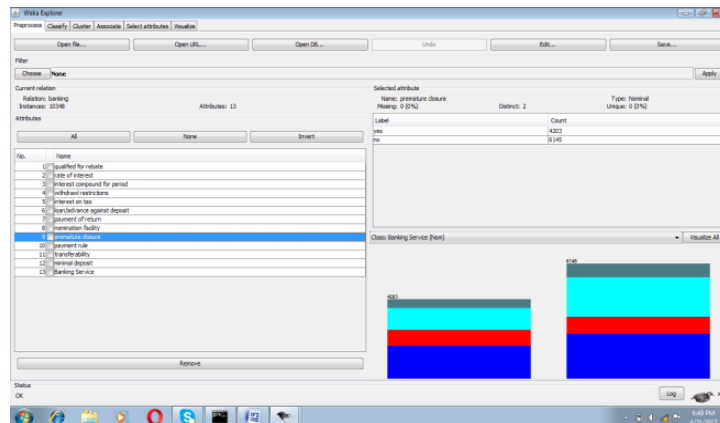


Figure 4 shows that the attribute is numeric. So its min, max ,mean and standard deviation values are depicted in the figure. Graph shows the distribution of banking services in the attribute.

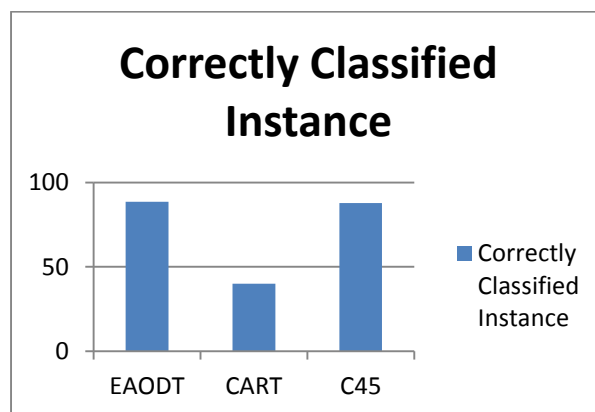


**Figure 5** Description of 10<sup>th</sup> attribute of dataset

Figure 5 describes the premature closure attribute which is a nominal attribute. In case of nominal attribute WEKA specifies the various values (labels) under that attribute along with number of instances under that label of the attribute. Graph shows the distribution of banking services in premature closure attribute.

**Table 1:** Analysis between correctly classified instance

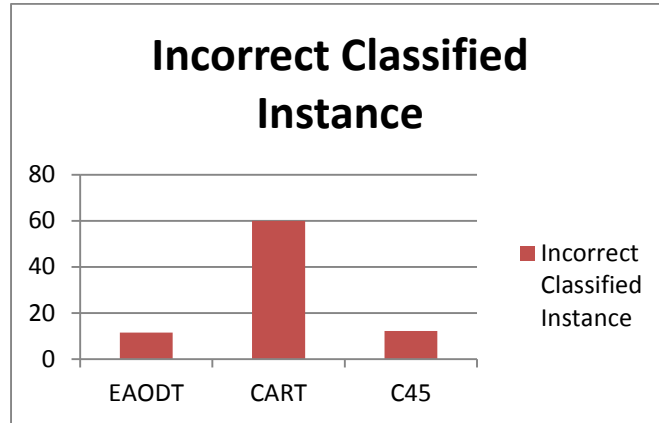
	EAODT	CART	C45
Correctly Classified Instance	88.51	39.96	87.78



**Figure 4.18** Shows Correctly Classified Instance Graphical representations of EAODT, CART, C45 results

**Table 2:** Analysis between Incorrect classified instance

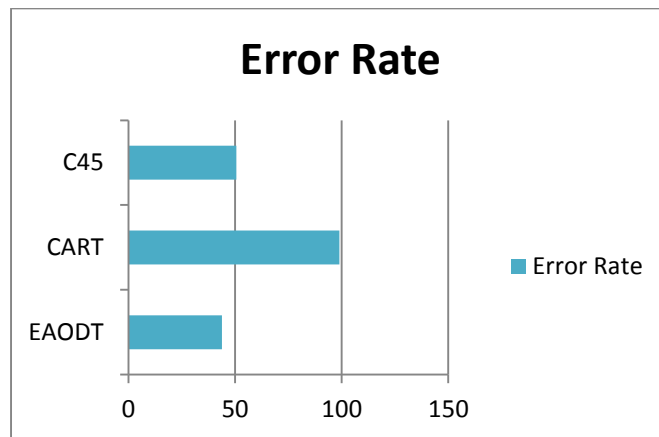
	EAODT	CART	C45
Incorrect Classified Instance	11.48	60.03	12.21



**Figure 4.19:** Shows Incorrect Classified Instance Graphical representations of EAODT, CART, C4.5 results

**Table 3:** Analysis between Incorrect classified instance

	EAODT	CART	C45
Error Rate	43.82	98.95	50.6



**Figure 4.20** Shows Error Rate Graphical representations of EAODT, CART, C45 results

**CONCLUSION:**

In this Research paper, we presented classification technique decision tree. We presented decision tree algorithm CART and C4.5. We focused on key elements of construction of decision tree. We did comparison of EAODT, CART AND C4.5 algorithms. It is concluded that EAODT is more accurate and consume less execution time to mine data with minimum error rate is a best algorithm for mining a data set.



**REFERENCES:**

- [1]. Fong, P.K. and Weber-Jhanke, J.H (2012), "Privacy Preserving Decision Tree Learning using Unrealized Data Sets", IEEE Transactions on knowledge and Data Engineering, Vol.24,No.2, February 2012, pp. 353-364
- [2]. Kabra, R.R. and Bichkar, R.S. (2011),"Performance Prediction of Engineering Students using Decision Tree", International Journal of Computer Applications, Vol.36, No.11, December 2011, pp. 8-12.
- [3]. Karaolis, M.A. &Moutiris, J.A (2010), "Assessment of the Risk Factors of Coronary Heart Events Based on Data Mining with Decision Trees", IEEE Transactions on Information Technology in Biomedicine, Vol.14, No.3, May 2010, pp. 559-566.
- [4]. Kesavraj, G. and Sukumaran, S. (2013), "A Study on Classification Technique in Data Mining", 4<sup>th</sup> ICCNT-2013.
- [5]. Sautikar, A.V., Bhujada, V., Bhagat, P.&Khparde, A.(2014)," A Review paper on Various Data Mining Techniques", International Journal of Advanced Research in Computer Science and Software Engineering, Vol.4,Issue 4, April 2014, pp. 98-101.
- [6]. Li, L. & Zhang, X. (2010), "Study of Data Mining Algorithm based on Decision Tree", 2010 International Conference on Computer Design and Applications (ICCD 2010), Vol.1, pp. 155-158.
- [7]. Yi-Yang, G. and Man-ping, R. (2009), "Data Mining and Analysis of Our Agriculture based on the Decision Tree", ISECS International Colloquium on Computing, Communication, Control and management, 2009, pp. 134-138.
- [8]. Zhang, X.F. and Fan, L.(2013)," A Decision Tree Approach for Traffic accident Analysis of Saskatchewan Highways", 26<sup>th</sup> IEEE Canadian Conference of Electrical and Computer Engineering(CCECE) 2013.
- [9]. Zhang, T., Fulk, G.D. & Tang, W.(2013),"Using Decision Tree to Measure Activities in People with stroke", 35<sup>th</sup> Annual International Conference of the IEEE EMBS, July 13, pp.6337-6340.
- [10]. Suknovic, .M, Delibasic, B., Jovanovic, M., Vukecevic, M., Obradovic, Z.(2011),"Reusable components in decision tree induction algorithm",Comp Stat Februaury 2011.
- [11]. Er. Harpreet Kaur" Classification of data using New Enhanced Decision Tree Algorithm (NEDTA)" in IJETCAS(International Association of Scientific Innovation & Research Issue 8, Volume 2, pp. 147-152, March-May, 2014
- [12]. Er. Harpreet Kaur "Optimizing Fuzzy Clustering using Swarm Intelligence in Data Mining" in international Journal of Advanced. Research in Computer Science Volume 2 No. 4.

- [13]. Er. Harpreet Kaur."Proposed Work for Classification and Selection of Best Saving Service for Banking Using Decision tree Algorithms" in International Journal of Advanced Research in Computer Science and Software Engineering Volume 3, Issue 9, September 2013 .
- [14]. Er. Harpreet Kaur."K-Mean Clustering and PSO: A Review" in International Journal of Engineering and Advanced Technology Volume/ Issue:3-5 Publication: June 30, 2014.
- [15]. Er.Harpreet Kaur.."Proposed Work for Classification and Selection of Best Saving Service for Banking Using Decision tree Algorithms" in International Journal of Advanced Research in Computer Science and Software Engineering Volume 3, Issue 9, September 2013 .