

International Journal of Advanced Research in Computer Science

**REVIEW ARTICLE** 

Available Online at www.ijarcs.info

# **Review of Decision Tree Data mining Algorithms: CART and C4.5**

Satbir Kaur and Harjit Kaur Department of Computer Science & Engineering Lovely Professional University, Phagwara, Punjab, India

## 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 an 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 CART and C4.5 algorithms and we will make comparison of these two algorithms.

# 1. INTRODUCTION

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. 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

### 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

## 1. 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].



Set of possible

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.
- **CART** : CART algorithm is presented by J.R. 4. 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.
- Select attribute with highest information gain measure 2.
- Create node containing that attribute. Iteratively apply 3. 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)---ET (Sj)

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

#### C4.5

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

- Splitting of categorical attribute is same to CART 1. algorithm. Continuous attributes always generate binary splits.
- Attribute with highest gain ratio is selected. 2.
- 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].

### Advantages: C4.5 made improvements to CART[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 CART AND C4.5:**







Figure 4: Define Class Variables



Figure 5: Generating Tree based on class attribute services

Error rate						0.4911				
Valu	es pred	iction	Confusion matrix							
Value	Recall	1-Precision		airport	store	movie	restaurant	station	Sur	
airport	0.2896	0.5640	airport	75	23	69	59	33		
store	0.3115	0.6049	store	0	81	91	36	52		
movie	0.5413	0.5054	movie	54	24	321	92	102		
restaurant	0.5802	0.4931	restaurant	22	53	114	293	23		
station	0.6130	0.4023	station	21	24	54	98	312		
			Sum	172	205	649	578	522		

© 2015-19, IJARCS All Rights Reserved

#### Figure 6: Calculate Error Rate

1 - M M									
D	efa.ik tile	HTML	Chart						
🗄 🎹 Detaset (mobile error	ment.bxt)		360	serviced Learning 1 (CS-CRT)					
🖃 🏠 Define status 2	1		Parameters						
Supervised Le	arning 1 (CS-CR1)	Classification to	ree (C-RT) parameters						
		Size before split	10						
		Pruning set size	(5) 33						
		30.5E rule	1.00						
		Random general	ior 1						
		Show all tree set	Show all tree seq (even if + 15) 0						
		100000000000000000000000000000000000000							
		Misclassific	ation Cost Matrix		,				
		Misclassific	Ation Cost Matrix		>				
Data visualization	Statistics	Misclassific	Attion Cost Matrix Meanuartier	Feature construction	Festure sélection				
Data visualization Regression	Statistics Factorial analysis	Nonparametric statistics PLS	Components Instance selection Clustering	Feature construction Spy learning	Festure selection Meta-spv learning				
Data visualization Regression Spv learning assessment	Statistics Factorial analysis Scoring	Nonparametric statistics FLS Association	ation Cost Matrix Meanwart or Components Instance selection Clustering	Feature construction	Festure selection Meta-spr learning				
Data visualization Regression Spo learning essessment de Binary logistic regression	Statistics Factorial analysis Scoring R.C.RT	Nonparametric statistics PES Association	ation Cost Matrix Measured or Components Instance selection Clustering	Feature construction Sov learning	Feature selection Meta-sprieaming				
Data visualization Regression Spriearing assessment delinary legistic regression Co.C.4.3	Statistics Factorial analysis Scoring A, C.RT A, CS-CRT	Misclassific	ation Cost Matrix Hearnart in Components Instance selection Clustering List	Feature construction Spr learning X-NN Unear discriminant analysis	Feature selection Acta-spr learning >-Multiliver perceptron 				
Data visuelization Regression Spo learning essessment C-R-S C-C-S C-C-S	Statistics Factorial analysis Scoring A_C.CRT A_CC-CRT A_CC-CRT A_CS-MC4	Nonparametric statistics FLS Association Egy Cover Ta Decision	ation Cost Matrix Wearupt or Components Instance selection Clustering	Feature construction Soviearning X-NN Unear discriminant analysis Ling-Reg TRIRLS	Feature telection Metacov learning ⇒ Multinonial Logistic Regression ⊒in Multinonial Logistic Regression ⊒in Multi bryos				

Figure 7: Parameters of C4.5 Algorithm

	1000		11774.01									
E Detaset (mobile ervor E Define status 2		Classifier performances										
Supervised Learning 1 (CS-CRT)			Errorrate			0.0461						
			Values prediction			Confusion matrix						
			Value	Recall	I-Precision		airport	store	movie	restaurant	station	S
			airport	0.9537	0.0500	airport	247	10	0	1	1	
			store	0.9538	0.1237	store	1	248	0	1	10	
			movie	0.9582	0.0190	reovie	1	12	567	11	z	
			restauran	e 0,9485	0.0264	restaurant	0	13	0	479	13	
			station	0.9568	0.0507	station	11	0		0	487	
				Com	iponents							
Data visualization	ata visualization Statistics Nonpara		Nonparametric statistics Instance		Instance se	election Feature construction			n	Feature selection		
Regression ipv learning assessment	Factorial analysis Scoring	Ad	PLS sociation		Cluster	ing	Spv	learning		Meta spv lear	ning	
C4.5 ALCS-CRT		Logic-Svc Mu Decision List		SK-NN Curver discriminant analysis chalor-Reg TBIRLS		A de	Multileyer perceptron		essio			

Figure 8: Calculate error rate

D	efa.ik tile	HTML	Chart						
<ul> <li>Detaset (mobile emo</li></ul>	rming 1 (CS-CBT)	Data parti Scowngust Prumg set Trees sequ H" # Leeves 30 1 8 11 1 113 Tree des	Data partition           Data partition           Instrument         Inst.           Inner int         Title           Press sequence (# 30)         Inst.           Inst.         Carlot         Carlot           Inst.         Carlot         Carlot         Carlot           Inst.         Carlot         Carlot         Carlot         Scott           Inst.         Carlot         Carlot         Scott         Scott         Scott           Inst.         Carlot         Carlot         Scott         Scott         Scott         Scott						
		<u>&lt;</u>	Comosoeria					0	
Data visualization	Statistics	Nonparametric statistics	Instance selection	Featur	e constructio	an i	Feature selection		
Regression Factorial analysis Spr learning assessment Scoring A		PLS Association	Clustering	S	pv learning		Meta-spv learning		
a Binary logistic regression C4.5 CPLS	ALC-RT ALCS-ORT ALCS-MC4	[25 ⊂ SVC ∰u Decision Agi 103	1 Det	EK-NN Culture discr CarLog-Reg TRI	iminant analy: RLS	A -3 01	Multilleyer perceptron Multinomial Logistic Regre Naive bayes	ustor	

Figure 9: Calculate no of nodes

Result:

	CART	C4.5
Error rate	0.40	0.04
Nodes	61	217
Leaves	31	109
Execution time	94ms	125ms

#### 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 CART AND C4.5 algorithms. It is concluded that C4.5 is more accurate and consume less execution time to mine data with minimum error rate 0.04. C4.5 is a best algorithm for

CONFERENCE PAPERS National Conference on Emerging Trends on Engineering & Technology (ETET-2017) On 21<sup>±</sup> April 2017 University Inst. of Engg. & Tech. & University Inst. of Computer, SBBS University, Punjab (India) mining a data set.

#### REFERENCES

- 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.& Khaparde, 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 (ICCDA 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 February 2011.