



Implementation and Analysis the Performance of (IDTA) Improved Decision Tree Algorithm on Mobile Services Dataset

Dr. Gurpreet Singh¹, Er. ManjotKaur²

Professor & Head¹, M.Tech Scholar²

Department of Computer Science & Engineering

St. Soldier Institute of Engineering & Technology, Jalandhar, Punjab, India

Abstract:

Data mining is a process of extraction of useful information and patterns from huge data. It is also called as knowledge discovery process, knowledge mining from data, knowledge extraction or data /pattern analysis. We present an improved approach to support nearest neighbor queries from mobile hosts by leveraging the sharing capabilities of wireless ad-hoc networks. We illustrate how previous query results cached in the local storage of neighboring mobile peers can be leveraged to either fully or partially compute and verify spatial queries at a local host. The feasibility and appeal of our technique is illustrated through extensive simulation results that indicate a considerable reduction of the query load on the remote database.

Keywords: mobile Services, CART, C4.5, IDTA

I. INTRODUCTION:

Various algorithms and techniques like Classification, Clustering, Regression, Artificial Intelligence, Neural Networks, Association Rules, Decision Trees, Genetic Algorithm, Nearest Neighbor method etc., are used for knowledge discovery from databases. But here we are going to discuss Association rules mining. So, having information about our data business and data mining techniques we can decide what we will use. Or we can try them all (if we have enough time, money and data) and find out which one is the best in our case. Decision tree is one of the important analysis methods in classification. It builds its optimal tree model by selecting important association features. While selection of test attribute and partition of sample sets are two crucial parts in building trees. Different decision tree methods will adopt different technologies to settle these problems. Traditional algorithms include C4.5, ID3, CART, SPRINT, SLIQ etc. ID3 is the representation of decision tree method. It is easy to understand and has fast classified speed which is applicable to large datasets. Many decision tree algorithms are improved based on it, like CART, SLIQ. But these algorithms more or less have some problems in selection of test features, type of samples, memory utilization of data and the pruning of trees etc. Presently, researchers have present many improvements. The dataset used in this research is based on mobile environment obtained from Wireless on software. This report is based on the findings maximum used of mobile service. The results in this report are based on data from mobile service related. As we look at Data Mining tools, we see that there are different algorithms used for creating a decision making (or predictive analysis) system. There are algorithms for creating decision trees such as C4.5 and CART along with algorithms for determining known nearest neighbor (KNN) or clustering when working on classification. The goal of this research is to look at one particular decision tree algorithm called enhanced algorithm and how it can be used with data mining for mobile service. The

purpose is to manipulate vast amounts of data and transform it into information that can be used to make a decision. In this work, I propose a technology based on data mining algorithms for the induction of decision trees. It is well suited in our context for various reasons.

1. To enhanced decision tree algorithm which will work on large scale high dimensional dataset- there is a problem of data mining in the classification of large datasets. There is no such algorithm stated that performs well in this problem. An algorithm can be made with certain split selection methods involved from the literature which includes algorithms like C4.5 and CART.
2. To enhance the efficiency with a new classifier that combines the k-Nearest Neighbor (CART) distance based algorithm with the classification tree paradigm based on the C4.5 algorithm.
3. To reducing presentsum of square error- the proposed algorithm gives reduced sum of square error as compare to the CART and C4.5 classification algorithm which means that the new algorithm gives more accuracy.
4. To enhancement in the efficiency of decision tree construction- various pruning techniques are proposed which can help in the improvement of decision tree construction.

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].

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

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

CART

The CART distance based algorithm with the classification tree paradigm based on the ID3 algorithm. The CART algorithm is used as a preprocessing algorithm in order to obtain a modified training database for the posterior learning of the classification tree structure. Then the incorrectly classified instances are duplicated with the previous data set and finally ID3 is applied to complete the classification procedure of biomedical data. In this approach a boosting technique is incorporated in such way that the incorrectly classified instances in the training set are identified using the k –NN algorithm. The performance of the proposed method is compared with the related algorithms. Experimental results show that the newly proposed approach performs better than the other existing techniques.

IDTA- Proposed Algorithm

```

create a node N;
if samples are all of the same class, C then
return N as a leaf node labeled with the class C;
if attribute-list is empty then
return N as a leaf node labeled with the most common class in
samples;
select test-attribute, the attribute among attribute-list with the
highest information gain;
label node N with test-attribute;
for each known value ai of test-attribute;
grow a branch from node N for the condition test-attribute = ai;
letsi be the set of samples in samples for which test-attribute =
ai; // a partition
ifsi is empty then
attach a leaf labeled with the most common class in samples;
else attach the node returned by Generate_decision_tree (si,
attribute-list- test-attribute);
    
```

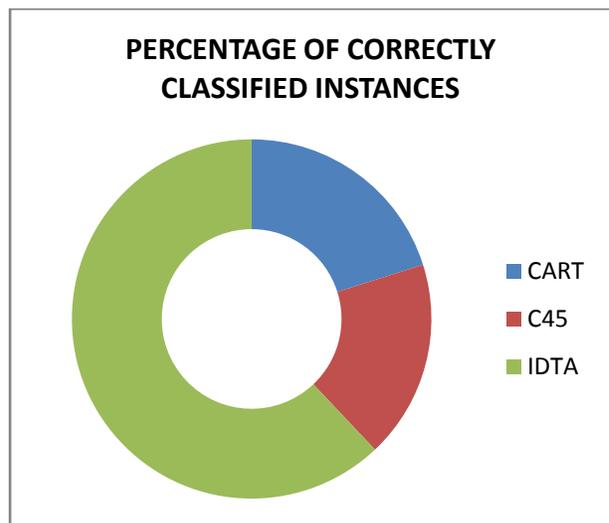
The basic strategy is as follows:

The tree starts as a single node representing the training samples (step 1). If the samples are all of the same class, then the node becomes a leaf and is labeled with that class (steps 2 and 3). Otherwise, the algorithm uses an entropy-based measure known as information gain as a heuristic for selecting the attribute that will best separate the samples into individual classes (step 6). This attribute becomes the “test” or “decision” attribute at the node (step 7). (All of the attributes are categorical or discrete value. Continues-valued attribute must be discretized.) A branch is created for each known value of the test attribute, and the samples are partitioned accordingly (steps 8-10). The algorithm uses the same process recursively to form a decision tree for the samples at each partition. Once an attribute has occurred at a node, it need not be considered in any of the node’s descendents (step 13). The recursive partitioning stops only when any one of the following conditions is true: All the samples for a given node belong to the same class (steps 2 and 3), or There are no

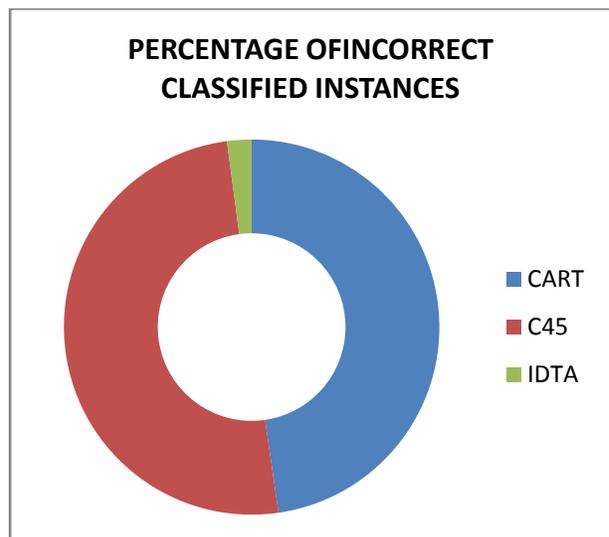
remaining attributes on which the samples may be further partitioned (step 4). In this case, majority voting is employed (step 5). This involves converting the given node into a leaf and labeling it with the class in majority among samples. Alternatively, the class distribution of the node samples may be stored. There are no samples for the branch test-attribute = a_i (step 11). In this case, a leaf is created with the majority class in samples (step 12).

Implementation and Analysis

	CART	C45	IDTA
PERCENTAGE OF CORRECTLY CLASSIFIED INSTANCES	31.45	27.89	96.97

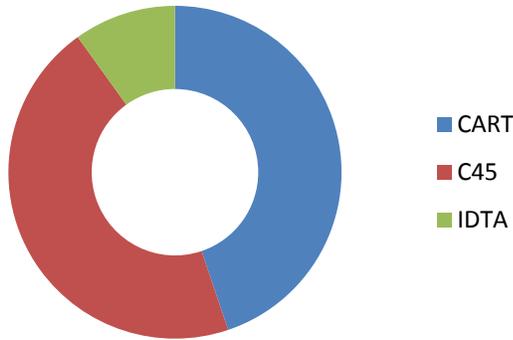


	CART	C45	IDTA
PERCENTAGE OF INCORRECT CLASSIFIED INSTANCES	68.54	72.1	3.02

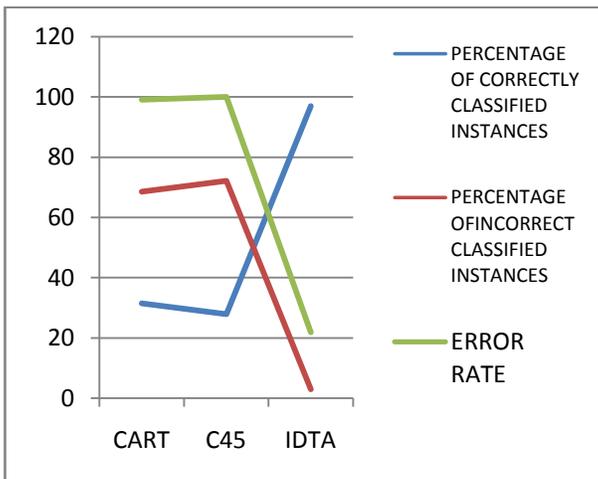


	CART	C45	IDTA
ERROR RATE	99.02	100	21.91

ERROR RATE



	CART	C45	IDTA
PERCENTAGE OF CORRECTLY CLASSIFIED INSTANCES	31.45	27.89	96.97
PERCENTAGE OF INCORRECT CLASSIFIED INSTANCES	68.54	72.1	3.02
ERROR RATE	99.02	100	21.91



II. CONCLUSION:

In this Research, I wanted to highlight the approaches for creating a decision tree. They are mainly available into academic tools from the machine learning community. I note that they are an alternative quite credible to decision trees and predictive association rules, both in terms of accuracy than in terms of error rate. After analysis Order C45, CART and Improved algorithm is more suitable to find accurate with minimum error rate. So enhanced algorithm is a best algorithm for mining a data on mobile services data set.

III. REFERENCES

[1]. Xiang Lian, Student Member, IEEE, and Lei Chen, Member, IEEE, "Ranked Query Processing in Uncertain Databases",

IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 22, NO.3, MARCH 2010.

[2].Stavroula G. Mougiakakou, Member, IEEE, "SMARTDIAB: A Communication and Information Technology Approach for the Intelligent Monitoring, Management and follow-up of Type 1 Diabetes Patients", IEEE TRANSACTIONS ON INFORMATION TECHNOLOGY IN BIOMEDICINE, VOL. 14, NO. 3, MAY 2010.

[3]. Eric Hsueh-Chan Lu, Vincent S. Tseng, Member, IEEE, "Mining Cluster-Based Temporal Mobile Sequential Patterns in Location-Based Service Environments", IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 23, NO. 6, JUNE 2011.

[4]. Mark N. Gasson, EleniKosta, Denis Royer, Martin Meints, and Kevin Warwick, "Normality Mining: Privacy Implications of Behavioral Profiles Drawn From GPS Enabled Mobile Phones", IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS—PART C: APPLICATIONS AND REVIEWS, VOL. 41, NO. 2, MARCH 2011.

[5]. Tzung-Shi Chen, Member, IEEE, Yen-Ssu Chou, and Tzung-Cheng Chen, "Mining User Movement Behavior Patterns in a Mobile Service Environment", IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS—PART A: SYSTEMS AND HUMANS, VOL. 42, NO. 1, JANUARY 2012.

[6]. R. Agrawal, T. Imielinski, and A. Swami, "Mining Associations between Sets of Items in Massive Databases," Proc. ACM SIGMOD, pp. 207-216, May 1993.

[7]. R. Agrawal and J. Shafer, "Parallel Mining of Association Rules," IEEE Trans. Knowledge and Data Eng, vol. 8, no. 6, pp. 866-883, Dec.1996.

[8]. R. Agrawal and R. Srikant, "Mining Sequential Patterns," Proc.11th Int'l Conf. Data Eng., pp. 3-14, Mar. 1995.