

Decision Tree Construction: A Continues Label Support Degree Based Approach

N.Madhuri¹, T.Nagalakshmi² and sujatha dandu³

¹ Auroras Technological and Research Institute (JNTU)

Received: 25 October 2011 Accepted: 22 November 2011 Published: 2 December 2011

Abstract

Data mining and classification systems utilize decision tree algorithms since they proffer rapid speediness, advanced exactness and also simple organization of those algorithms. An ideal decision can be built only when the appropriate attributes are chosen. This paper focuses on throwing light on choosing characteristics based on the theory of attribute support degree on account of which a unique decision tree construction algorithm is proposed on the basis of rough set and granular computing theory. It is henceforth proved that the decision tree proposed by the new approach yields far more better results in terms of precision and consistency as compared to the decision trees yielded by ID3, C4.5 and DTBAS.

Index terms— rough set: decision tree: granular computing: attribute support degree: attribute selection.

1 INTRODUCTION

Decision sets can be denoted using tree structures with the help of decision tree which is a unique, spontaneous, data illustration scheme and also a competent classifier. Quinlan et al [1] proposed ID3, decision tree algorithm and hence has been persistently augmented which have been advanced to C4.5 [2]. The preeminent attribute is chosen as the existing attribute which is then recursively inflates the decision tree branches unless and until the conditional statement is achieved, which ultimately makes use of top-down greedy algorithms. There are different classification schemes that can be achieved concerning different solutions which poses two issues [3] in decision tree construction. Choosing characteristics for crafting new branches in a tree is one issue while the other one is pruning which is all about omitting and decreasing the tree. DTBAS [10] considered the Assortment of attribute as main concern, which was refined and improved by considering assortment of continuous labels also that is discussed in this paper.

Z. Pawlak et al [4] recommended rough set theory which is an expansion of set theory for studying intelligent systems which is followed up by inadequate and partial data information. There is a thriving submission of the rough set theory in the disciplines of data mining, pattern recognition, machine learning, decision analysis etc in recent times. Models are categorized into various resembling classes that houses imperceptible objects in terms of few attributes. Issues pertaining to feature selection, data reduction and pattern extraction can be amicably taken care of such that it can liberate the system of redundant data in systems containing null values or missing data.

Lin et al [5] proposed the expression of Granular Computing which spans itself covering all aspects of concerning theories, tactics, practices and means essential in solving a problem that makes use of granules. Granular Computing has witnessed vast inputs from different practices such as fuzzy sets, rough sets, shadowed sets, probabilistic sets etc.

A crucial step that needs to be taken care of while building a decision tree is choosing characteristics of nodes of a tree that houses minimum number of branches. Decision tree based on continuous label support degree (DTBLS) algorithm is introduced which is considered as a splitting criterion on account of rough set theory and granular computing. Trial results have approved the usage of DTBLS algorithm that assures and provides uncomplicated structures and superior categorization accuracy.

45 The rest of the paper is organized as follows: Section 2 discusses concepts relevant to rough set theory and
 46 granular computing. Section 3 gives a basic introduction to our new method and presents a simple example.
 47 Experimental comparison of the proposed method with 103 and C4.S is given in section 4. The final section
 48 concludes the research work of this paper.

49 **2 II. BASIC CONCEPT**

50 Few fundamental concepts of rough set theory [6,7] and granular computing [8] are first initiated for ease of
 51 demonstration. Definition 1 (Information System) : An information system can be labeled as wherein U is a
 52 finite set of object known as the universe; A , which is a non-vacant finite group of attributes; C and D depict set
 53 of condition and decision attributes respectively as also v , which says that v is a value set of the attribute a and
 54 Where $|I|$ is the cardinal number of I . (I) 2011 December (, ,) $S U A V f = A C D = U$, $a v v a A = ? ? U a v$
 55 When in an equivalence relation R , the granular degree of R reaches the minimum size ; When R is a domain
 56 relation, the granular degree of R attains the maximum size Definition 5 : Assume R is knowledge of repository
 57 , the granular degree of basic knowledge is defined asIII.

58 **3 PROPOSED ALGORITHM**

59 This section aims at familiarizing the algorithm of building a decision tree on the basis of attribute support
 60 degree.

61 **4 a) The Principle of Label and Attribute Selection**

62 The label that represents least average uncertainty is supposed to be chosen as the test label and then choose
 63 attribute with less uncertainty as test attribute from the class represented by the selected label, which because
 64 it makes apt decisions when compared to existing test attribute selection in different decision tree algorithms. (
 65 Attribute Support Degree) Let S be an information system is a label contains subset of attributes represented as
 66 . Attribute support degree can be denoted as follows based on the definitions mentioned above.

67 **5 Where**

68 denotes the cardinal number of S .

69 **6 Global**

70 I) 2011 December (f U A V = \times ? (,) a f u a V ? a A ? , u U ? (()) X BND X (, , ,) S U A V f = R A ? X
 71 U ? () R X () R X () BND X () R X Ji(x), () R X () BND X () { | | } ... (1) R R X x U x X = ? ? () { | |
 72 } ... (2) () () () ... (3) R R X x U x X BND X R X R X ? = ? ? ? = ? (, , ,) S U A V f = P A ? () { (,) | , (
 73 ,) (.) ... (4) a IND P x y U U P f x a f y a = ? \times ? ? = () U IND P (,) K U R = R U U ? \times () G D R 2 | | |
 74 | () ... (5) | | | R R G D R U U U = \times | | R R U U ? \times 2 | | / | | 1 / | | U U U = 2 2 | | / | | 1 U U = (,) K
 75 U R = , 1 2 / { , , , , } n U R X X X = 2 1 2 | | () ... (6) | | n i i i x G D X U = = ? 1 2 { , , , , } l s a a a | | 1 (
 76) () ... (7) | | l i i u c u c a a v g l l = = ? () u c a v g l l l () i u c a i a l (, , ,) S U A V f = , A C D = U 1 Q C ? (
 77) | () | (,) ... (8) () () G D Q D IND QUD S Q D G D Q IND Q = = U | () | IND Q () IND Q U U ? \times

78 Decision Tree Construction: A Continues Label Support Degree Based Approach
 79 . Definition 5 states that whenever we get the relations among them, namely, when $GD(R)$ is smaller, the
 80 distinguishable degree is stronger and $S(Q,D)$ is greater, thereby Q is better sets of test attribute of D . On the
 81 contrary, the smaller $S(Q,D)$ is, the worse we get Q as sets of test attribute of D .

82 **7 b) The Description of DTBLS D**

83 The basic notion of DTBLS D expresses the point that whenever label support degree with association of label
 84 level attribute support degree is made use of as a customary for choosing a test attribute concerning every node
 85 in the decision tree. The attribute reduction set assists in selecting a condition attribute that possesses the
 86 highest degree of label level attribute which can be put to use at the root of the decision tree. There will be a
 87 testing of the remaining condition attributes on each and every branch of the root node $S(Q,D)$ D using Q can
 88 be estimated with the help of a measure addition of new sub-trees to every division until the leaf is reached.
 89 According to the above idea, using the $S(Q,D)$ as the splitting criterion, we propose our algorithm DTBLS D.
 90 Current sample set is depicted by T , set of labels depicted by L , condition attribute set of a label is depicted
 91 by the S . depicts the number of attributes in the condition attribute set of label S . All attributes of the condition
 92 attribute set are discrete and continuous values are discretized by continuous labeling. Following are the specific
 93 steps of the algorithm.

94 Algorithm : A decision tree is created by DTBLS D (T , attribute list) that using the given training data. The
 95 top-down recursive divide and conquer approach for construction of a decision tree wherein the recursion related
 96 division takes place only when any one criterion mentioned below is gratified. A common class contains :

97 **8 Input**

98 1. All specimens for a specific branch which restores a leaf that is termed with the concerned class.

99 Here, a large case of voting is provisioned to change the present working node into a leaf that is termed with the
100 concerned class that in in demand from amongst various specimens. 2. In addition, there are no more specimen
101 test attributes and the class division specimens can be placed wherein a leaf is generated and termed with the
102 most featured class in specimens.

103 IV.

104 9 EXPERIMENTS a) Example Analysis

105 Table ?? showcases a data tuple training group originated from All Electronics customer records that are
106 implemented using polic mentioned in reference (6). The first step is to estimate the degree of attribute support
107 for each situational attribute or characteristic. ()

108 10 U IND al

109 = { {1,2,8,9, 11 },{3,7,12,13},{

110 11 U IND az

111 = { {1, 2, 3, 13},{4, 8, 10, 11, 12, 14}, {
112 (U IND a = {{1, 2, 3, 4, 8, 12, 14}, {53)

113 12 U IND a

114 = { {1,3,4,5,8,9,10,13},{2,6,7,11,12,14}};

115 () Table ?? : Training data tuple from the AIIElectronics customer database Notations used in example
116 descriptions: Age???, Income???, Student???, Credit???, Buy?? is selected as the min root of the decision tree
117 and is tagged with age since is the maximum extent of a degree from amongst all the condition attributes as
118 also various number of divisions which are branched in reference to a range of different attributes. In case where
119 age=1, all the specimens that are grouped into this should belong to the same class and hence a leaf should
120 be generated at the end of every division and should be tagged with d=yes. The figure above depicts the final
121 decision tree that is built by DTBLSD.U IND d = { {1 ,2,6,8,14 },{3,4,5,7,9, 10, 11,12, 13}}; (,)

122 13 b) Experimentl Comparison

123 Experimental comparison of DTBLSD with respect to ID3 [1], C4.5 [2] and DTBAS [10] is discussed in this section.
124 The real datasets that are used in this are approved from University of California, Irvine (UCI), and is known
125 as the machine learning database repository where C++ design language is implemented to form the requisite
126 algorithm. WEKA 3.7 is used for successful accomplishment of ID3 [1] and C4.5 [2] which is a compilation of
127 machine learning algorithms used for data mining generated and procured by ??rank

128 14 CONCLUSION

129 The paper first focuses on explaining the basic notion of label support degree and attribute support degree
130 [10] and selecting it as a basic decisive factor on the basis of degree of involvement between condition attribute
131 and decision attribute accordingly where a unique decision algorithm tree based on continuous label support
132 degree and label level attribute support degree (DTBLSD) is recommended. Accordingly a suitable methodology
133 is devised which is flexible enough to accommodate and provides lower complexity and high level of accuracy
134 as compared to other algorithm generating methods. A disadvantage identified in [10] is issues pertaining to
135 adjustment with adaptability of samples, which has been overcome successfully in our model. ^{1 2}

¹© 2011 Global Journals Inc. (US)

²Global I) 2011 December (



6

Figure 1: Definition 6 :

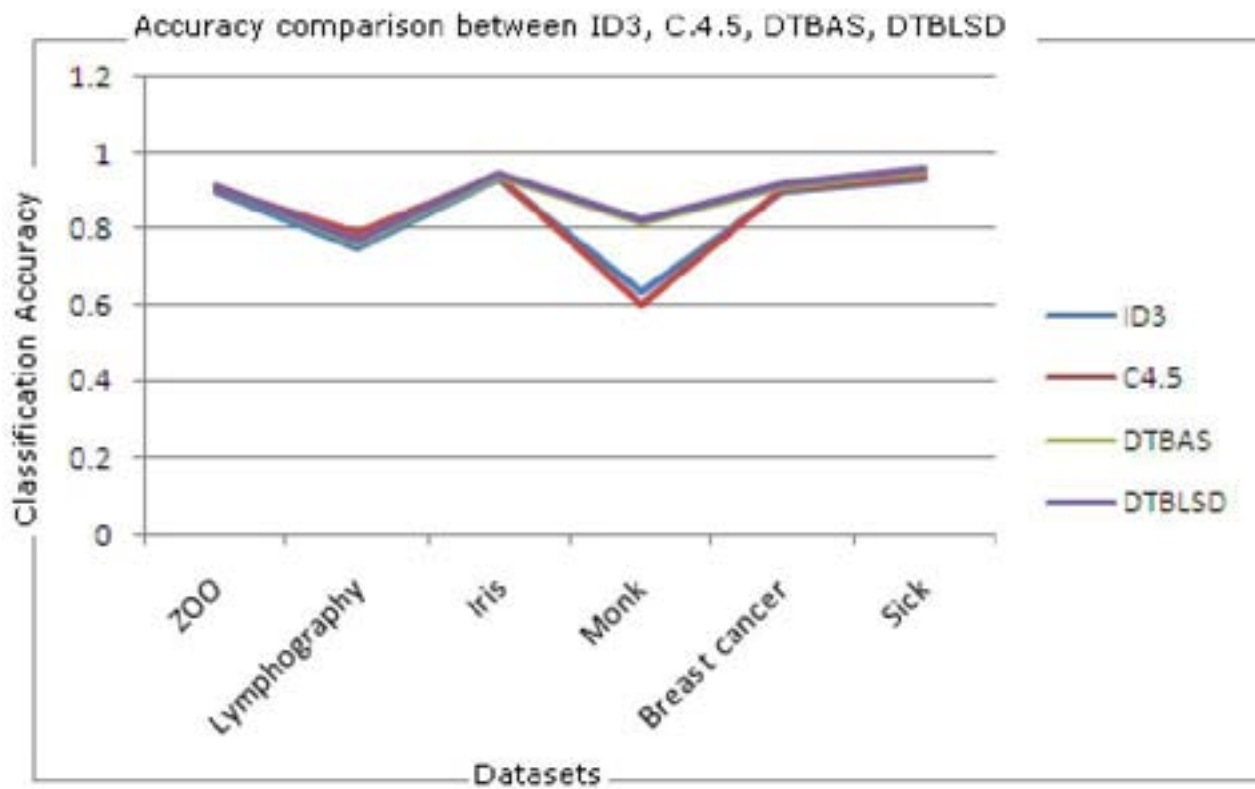


Figure 2:

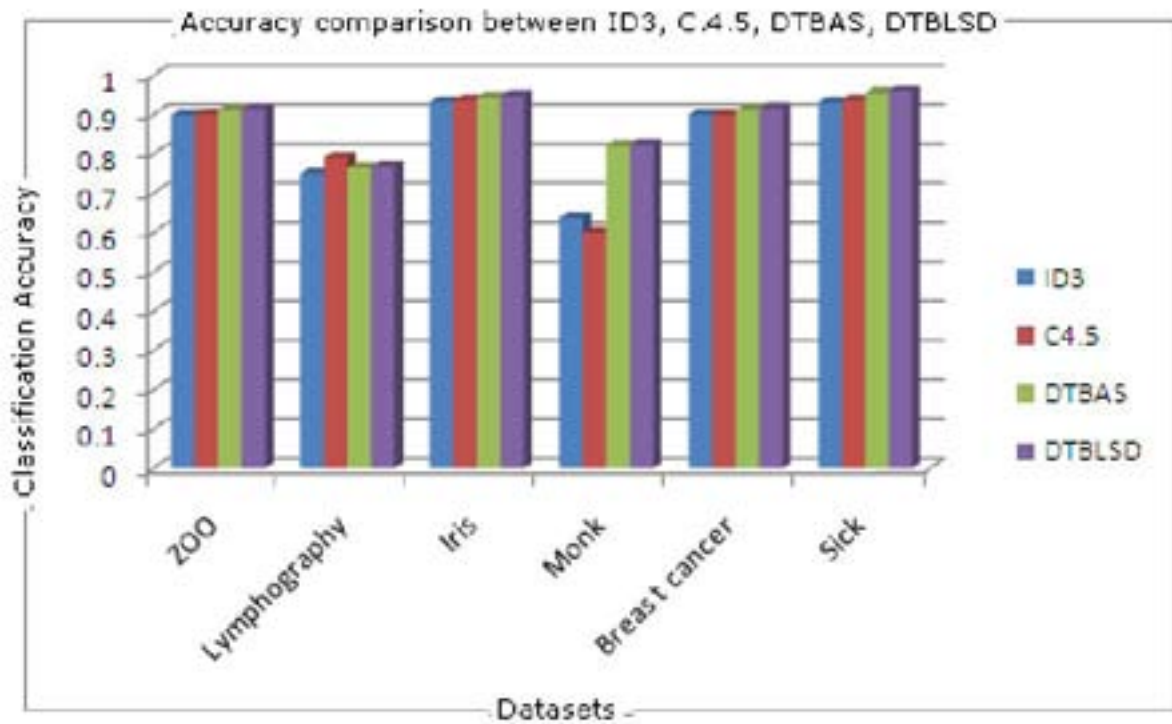


Figure 3:

U	age	income	student	credit
1	<=30	high	no	fair
2	<=30	high	no	excellent No
3	31<=age<=40	high	no	fair
4	>40	medium	no	fair
5	>40	low	yes	fair
2011 6 7 8 9 10	>40	low low	yes yes	excellent No excellent
De-	31<=age<=40	medium	no yes	
cem-	<=30 <=30	low	yes	
ber	>40	medium		
11	<=30	medium	yes	excellent Yes
12	31<=age<=40	medium	no	excellent Yes
13	31<=age<=40	high	yes	fair
14	>40	medium	no	excellent No
1 a				

1 (,) S a d

I)

(

Global that involved the 10 fold cross estimation to calculate classification authenticity. All experiments were

Figure 4:

II

Dataset	ID3	? ? C4.5	DTBA S	DTBLS D
ZOO	0.899	0.901	0.911	0.914
Lymphograph y	0.75	0.791	0.765	0.767
Iris	0.932	0.936	0.943	0.947
Monk	0.637	0.6	0.821	0.824
Breas t cancer	0.9	0.9	0.912	0.916
Sick	0.931	0.937	0.955	0.959

Figure 5: Table II :

-
- 136 [Ding] , B S Ding , YQ .
137 [Frank et al.] , E Frank , M Hall , Weka . <http://www.cs.waikato.ac.nz/ml/weka>
138 [Zheng and Zang ()] ‘A New Decision Tree Algorithm Based on Rough Set Theory’. S Y Zheng , Zang . *Proc.*
139 *Asia-Pacific Conference on Information Processing*, (Asia-Pacific Conference on Information essing) 2009. p.
140 .
141 [Yong and Zhou ()] ‘An algorithm of decision tree construction based on attribute support degree’. Jianping Yong
142 , ; Jun Zhou . *Educational and Information Technology (ICEIT)*, 2010.
143 [Chen et al. ()] ‘Constructing a Multi-Valued and Multi-Labeled Decision Tree’. Y L Chen , C L Hsu , S C Chou
144 . *Expert Systems with Applications* 2003. 25 p. .
145 [December Decision Tree Construction: A Continues Label Support Degree Based Approach] *December*
146 *Decision Tree Construction: A Continues Label Support Degree Based Approach*,
147 [Lin ()] ‘Granular Computing II Infrastructure for AI-Engineering Examples, Intuitions and Modeling’. T Y Lin
148 . *Proc. 2006 IEEE International Conference on Granular Computing*, (2006 IEEE International Conference
149 on Granular Computing) GrC 2006. 2006. p. 2011. (I)
150 [Han and Kamber ()] J W Han , M Kamber . *Data Mining Concepts and Techniques*, 2001. Morgan Kaufmann
151 Publishers.
152 [Quinlan ()] ‘Improved Use of Continuous Attributes in C4.5’. J R Quinlan . *Journal of Artificial Intelligence*
153 *Research* 1996. p. .
154 [Quinlan ()] ‘Induction of Decision Trees’. J R Quinlan . *Machine Learning, vol1*, 1986. p. .
155 [International Conference on (2010)] *International Conference on*, Sept. 2010. 2 p. .
156 [Miao and Wang ()] D Q Miao , G Y Wang . *Granular Computing: Past, Now, Future[M]*, Science publishing
157 house, 2007.
158 [Qing Lin] *Qing Lin*, (Zongzhuan Ding)
159 [Pawlak ()] ‘Rough Sets’. Z Pawlak . *International Journal of Information and Computer Science, volII* 1982.
160 p. .
161 [Seifi et al.] *Twins Decision Tree Classification: A Sophisticated Approach to Decision Tree*, F Seifi , H Ahmadi
162 , M Kangavari .