# Application of Decision Tree in Analysis of Intra College Festival Data set

**[1]Sumit Srivastava (Dr.)**
**[1]Department of Information Technology,**
**Poornima College of Engineering, Jaipur, Rajasthan-302022(India)**

**[2]Megha Yadav**
**[2]Department of Computer Science,**
**Poornima College of Engineering, Jaipur, Rajasthan-302022(India)**

**[3]Sandeep Devra**
**[3]Department of Computer Science,**
**Jagannath Institute of Engineering & Technology, Jaipur, Rajasthan-302022(India)**

**Abstract -** Data mining has been used very frequently to extract hidden information from large databases. The paper suggested the use of data mining technique named decision trees along with its algorithm for continuously monitoring the College fest. In this way the data extracted could be used to teach the students on the real time problem scenario application to be monitored at college level. Also the model can be used for the future planning of such fest at college level. The algorithm used is CHAID (Chi-squared Automatic Interaction Detection) using PASW18 as software.
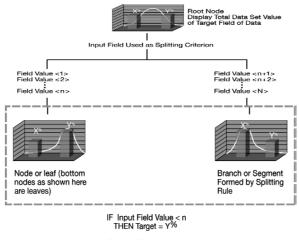
*Keywords-* Data Mining, ID3, C4.5, CART, CHAID

## 1. Introduction

It is the process of extracting patterns from large data sets by combining methods from statistics and artificial intelligence with database management. It is an increasingly important tool by modern business to transform data into business intelligence. It is currently used in a wide range of profiling practices, such as marketing, surveillance, fraud detection, and scientific discovery. In this paper, the effective use of the information available by using the decision trees on a festival data set. The problem is to classify the information available on the basis of year of student participated in the various events, participation on the basis of the gender, particiaption on the basis of type of event(can be clutural, club, sports and tehnical), branch wise information. The algorithm available to be used in Decision trees ID3, C4.5, CART, CHAID to be used

with software tools like PASW18 onwards. Also the comparison with the other algorithms is being done.

## 2. Decision Trees

Decision trees are often used in classification and prediction. It is simple yet a powerful way of knowledge Representation [1]. The models produced by decision trees are represented in the form of tree structure. A leaf node indicates the class of the examples. The instances are classified by sorting them down the tree from the root node to some leaf node.



**FIG 2.1: Decision tree reflecting both continuous and categorical object of analysis.**

The reason of using Decision trees are

1. Decision trees are white boxes means they generate simple, understandable rules
2. Decision trees are non-parametric means no specific data distribution is necessary.

3. Decision trees handle missing values as easily as any normal value of the variable.
4. In decision trees elegant tweaking is possible.
5. Decision trees identify subgroups. Each terminal or intermediate leave in a decision tree can be seen as a subgroup/segment of your population.
6. Decision trees run fast even with lots of observations and variables
7. Decision trees can be used for supervised and unsupervised learning.
8. Decision trees can easily handle unbalanced datasets.
9. Versatility for a wide variety of data mining tasks, such as classification, regression, clustering and feature selection

while some of the disadvantage lies as following

1. While decision trees classify quickly, the time for building a tree may be higher than another type of classifier.
2. Decision trees suffer from a problem of error propagation
3. Disjunctive descriptions may be required
   a. naturally represent disjunctive expressions

the algorithms being used while decision tree analysis is as follows.

**2.1 ID3 Algorithm:-** it is a greedy search technique, for which top-down recursive divide and conquer manner. Generally the question relates to classifying the information. It generally use information gain as a quantitative measure of an attribute[3]. The algorithm for attribute selection measure

- Select the attribute with the highest information gain
- Let $p_i$ be the probability that an arbitrary tuple in D belongs to class $C_i$, estimated by $|C_{i, D}|/|D|$
- Expected information (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^{m} p_i \log_2 (p_i)$$
……(2.1)

- Information needed (after using A to split D into v partitions) to classify D:
  o Information gained by branching on attribute A (If A is discrete valued):

$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times I(D_j)$$

……..(2.2)
- Information gained by branching on attribute A (If A is discrete valued):

$$Gain(A) = Info(D) - Info_A(D)$$
……..(2.3)

**2.2 C4.5 Algorithm:-** C4.5 (a successor of ID3) uses gain ratio to overcome the problem. Here for Computing the information- gain for continuous value attributes the best spilt-point is being calculated [2][4]. The algorithm for calculating the same is:

- Let attribute A be a continuous-valued attribute
- Must determine the *best split point* for A
- Sort the value A in increasing order
- Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*
- $(a_i+a_{i+1})/2$ is the midpoint between the values of $a_i$ and $a_{i+1}$
- The point with the *minimum expected information requirement* for A is selected as the split-point for A

and the SPLIT is calculated as

- D1 is the set of tuples in D satisfying A ≤ split-point, and D2 is the set of tuples in D satisfying A > split-point

It uses the kind of normalization to information gain using a "Split Information" value defined analogously with Info(D) as:

$$SplitINFO_A(D) = -\sum_{i=1}^{k} \frac{D_i}{D} \log \frac{D_i}{D}$$
……..(2.4)

and the gain ratio defined as

$$GainRATIO\ (A) = \frac{GAIN\ (A)}{SplitINFO\ (A)}$$
…………(2.5)

**2.3 CART Algorithm:** Classification and Regression Trees (CART) is a non-parametric technique that produces either classification or regression trees, depending on whether the dependent variable is categorical or numeric, respectively. The Gini index is used in CART. Gini index measures the impurity of *D, a data partition or set of training tuples, as*

$$gini\ (D) = 1 - \sum_{j=1}^{n} p^2_j \qquad …(2.6)$$

IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 1, No 1, January 2012
ISSN (Online): 1694-0814
www.IJCSI.org

216

where $p_j$ is the relative frequency of class $j$ in $D$. The Gini Index consider the binary split of the each of the attribute. If the data set D is split into subset $D_1$ and $D_2$, the gini-index can be defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2) \quad ...(2.7)$$

and Reduction in Impurity (that would be incurred by a binary split on a discrete- or continuous-valued attribute $A$) is:

$$\Delta gini(A) = gini(D) - gini_A(D) \quad ....(2.8)$$

**2.5 CHAID Algorithm:-** Chi-squared Automatic Interaction Detection is a non-binary decision tree algorithm, where split search is designed for categorical values. CHAID uses a Chi-squared test using contingency tables. The advantage for using CHAID as method of evaluation is:

- Used in fields of marketing research and public health
- Can handle a large number of variables
- Designed to identify potentially meaningful patterns in a dataset
- Easy to use and understand

## 3. Decision Trees

The application data set consists of the cultural fest named "AAROHAN2k10" organized as annual fest in one of the private engineering colleges. It includes nearly 22 technical events, 13 club events, 7-8 sports events, 9-10 cultural events with 5500 approx student participating for the competition. For each student participation in activity involving Technical and club is a must with total of 4 events per student has to be submitted. So with almost 5500*4= 22,000 entries as input the event has been successfully completed. The snapshot of the data taken from AAROHAN2k10 is shown below



**Fig 3.1: Training Data Set (taken from Aahoran database)**

During the result analysis, it was asked to predict the participation in each types of events by each branch of each college with Gender of the students can be as the dependent variable for the analysis. Also if the event is supposed to a group then classification of the information on the basis of college, year, branch, team-name, event-name is requested. For such problem CHAID is selected as the best technique for the solution as information can be branched using non-binary variable. Some of the results snapshots calculated using PASW 18 beta version as software is shown below.
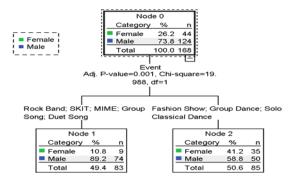


**Fig 3.2:- When the dependant variable is taken to be Gender and independent variables are: Event, College,**

IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 1, No 1, January 2012
ISSN (Online): 1694-0814
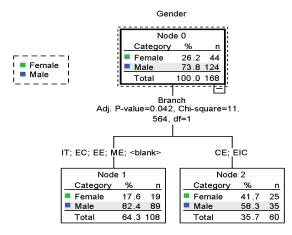www.IJCSI.org

217

Branch)



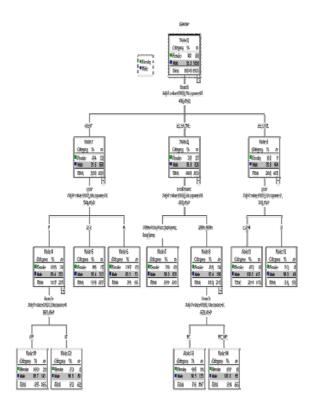Fig 3.3:- When the dependant variable is taken to be Gender and independent variables are: College, Branch)



Fig 3.4:- When the dependant variable is taken to be Gender and independent variables are: college, year, branch, team_name, event_name

## 4. Conclusion

The inference from the various result analyses can be shown as.

- It was predicted that 89.2% out of 73.9% of the males have participated in events like Rock Band, SKIT, MIME, Group Song and Duet Song.
- On the basis of the data studied it can be concluded that 58.8% out of 73.9% of the males have participated in events Cultural and Club Events
- it can be further concluded that 58.3% out of 73.8% of the males are from CE, EIC branches who have participated in cultural events
- On the basis of the data studied it can be concluded that 82.4% out of 73.8% of the males are from IT, ME, EC, EE branches who have participated in cultural events, sports and technical events

With fast analysis of large databases and classification of information using Decision trees, CHAID is founded as the one among the best algorithm for the large data analysis.

## 5. Acknowledgement

## 6. Bibliographies and References

[1] T. Mitchell, *"Decision Tree Learning*", in T. Mitchell, Machine
Learning, the McGraw-Hill Companies, Inc., 1997, pp. 52-78.
[2] P. Winston, "*Learning by Building Identification Trees*", in P. Winston, Artificial Intelligence, Addison-Wesley Publishing Company, 1992, pp. 423-442.
[3] Howard J. Hamilton's CS Course: Knowledge Discovery in Databases. Accessed 06/06/12.
[4] http://www.cs.waikato.ac.nz/ml/weka/, accessed 06/05/21.
[5]http://grb.mnsu.edu/grbts/doc/manual/J48_Decision_Trees.html, accessed 06/06/12.
[6] Quinlan, J.R.: C4.5: Programs for Machine Learning. Morgan Kauffman, 1993.

**[1]Dr. Sumit Srivastava**
**Professor,**
**Department of Information Technology,**
**Poornima College of Engineering,**
**Jaipur-302022(raj)**


**[2]Mrs. Megha Yadav**
**M.Tech, Research Scholar,**
**Department of Computer Science Engineering,**
**Poornima College of Engineering,**
**Jaipur-302022(raj)**


**[3]Mr.Sandeep Devra**
**Associate Professor,**
**Department of Computer Engineering,**
**JNIT, Jaipur-302022(Raj)**