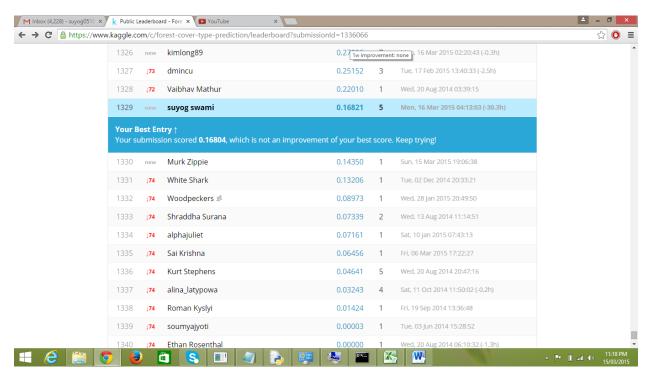## CSE 4334/5334 – Data Mining Spring 2015 – Programming Assignment 2

Kaggle Username: suyog swami

Kaggle Rank: 1329

Kaggle Score: .16821



Design and Implementation:

Classifier: Decision tree - ID3.

IDE : Spyder 2.3.1

Python version: 3.4.2

Details: In this programming assignment I have used decision tree ID3 classifier to predict the forest Cover Type results. In this program there is a single python code file containing 9 functions as follows:

1.  main(): where I accept the two csv files and give calls to the createTree and predict fumction. The function is also used to get the attributes and data of both the files. The operation to write into the result.csv file is also done in this function.
2.  createTree(): In this function the main recursive functionality of the program is placed. It contains function calls to five other functions. A) mostCommonAttriVal() B) chooseBestAttri() C)

getColVal() and D) getCurrentVal E) createTree()[This works recursively]. Also the stop condition for the recursion is check in the beginning of the createTree() function.

3.  mostCommonAttriVal(): This function is used to calculate the most common value for the attribute. The value of this function is used to check the error condition if the data set is empty or there is only one attribute in the dataset.

4.  chooseBestAtrribute(): this function is used to decide which is the best attribute to perform the create three operation in the specific run. The information gain and entropy are calculated from this function which is explained further.

5.  getColVal(): This function is used to get the column values for the best attribute decided in the chooseBestAttribute function.

6.  getCurrentVal(): This function is used to get the dataset from the training dataset which doesn't include records of the best attribute. This dataset is used for next round of recursion.

7.  infoGain(): Information gain which is Entropy of the decision attribute – (freq of best attribute)* entropy of the bestattribute. This function is called in a for loop from the chooseBestAttribute function. It returns the gain to the chooseBestAtrribute function. A call to calculate Entropy is given from this function.

8.  calculateEntropy(): This function calculates entropy using the entropy function discussed in the class.

9.  predict(): After generating the tree. Predict function is used to predict the output of the test.csv's Cover_Type. Here I have used a random function to fill the places whose Cover_Type values the classifier can't predict. Also a hard coding is done to check if the value of testdata attributes is nearest to the values in train and predict the result.

It was an overly exhaustive assignment where I feel I put more efforts than in any programming assignment done before. But I agree looking at the accuracy I wasn't able to optimize the result efficiently. Hardcoding was the last thing I opted for. I tried many variations but was unsuccessful in doing so. One variation I tried was to place the tree nodes as ranges of the attributes and match the ranges with the test data attributes, but it became a more complex process to handle the string and integers interchangeably.