Name: Yeung Wai Kit

Student ID: 20608456

Course: COMP4331

Title: Report for Decision Tree Mini Project

## Preprocessing

Before constructing the decision tree model, data should be preprocessed.

To prevent unexpected behavior in the decision tree model, I first fill all the null values in the dataset with corresponding attribute mean. Then all columns with identical values are dropped because those columns are meaningless for decision tree splits. The process is then followed by standardization to range (-1,1). This transforms data in each column to data with mean=0 and standard deviation =1 and helps unifying data units. Finally, to ensure model training speed we do PCA to reduce the number of features. I choose to preserve 90% portion of explained variance in this model.

## Model creation

My final goal is to build a random forest model. I will first do hyperparameter training to select best parameters for my forest. Then I will prune the DecisionTreeClassifier object with best parameter using cost complexity pruning to avoid overfitting. After all these I should be able to get the best parameters for a non-overfitting tree for my model. Finally, I will rebuild the random forest model using the list of best parameters and pruning parameter.

## Random Forest Model

The reason I use random forest model is because it is a more reliable compared to single decision tree model. Random forest model creates different decision trees using different randomly split datasets and use a majority vote method for final prediction. Using the hyperparameter tuning, every tree inside the forest has the best parameters and the model should be performing the best. To create a random forest model I will use RandomForestClassifier from sklearn.ensemble. This creates a random forest directly with parameters similar as a Decision Tree Model.

## Hyperparameter Tuning

For this step I will use GridSearchCV class. This class basically provide functions that allows me to loop over different parameters of my input model and select the best performing sets automatically. The looping parameter set I will use is:

{'criterion':['gini','entropy'],'max_depth':[250,270,300],'min_samples_split':[3,4,5],'max_leaf_nodes':[150,170,200],'n_estimators' : [500,1000]}

***The best parameters are

{'criterion': 'entropy', 'max_depth': 270, 'max_leaf_nodes': 170, 'min_samples_split': 3, 'n_estimators': 1000}

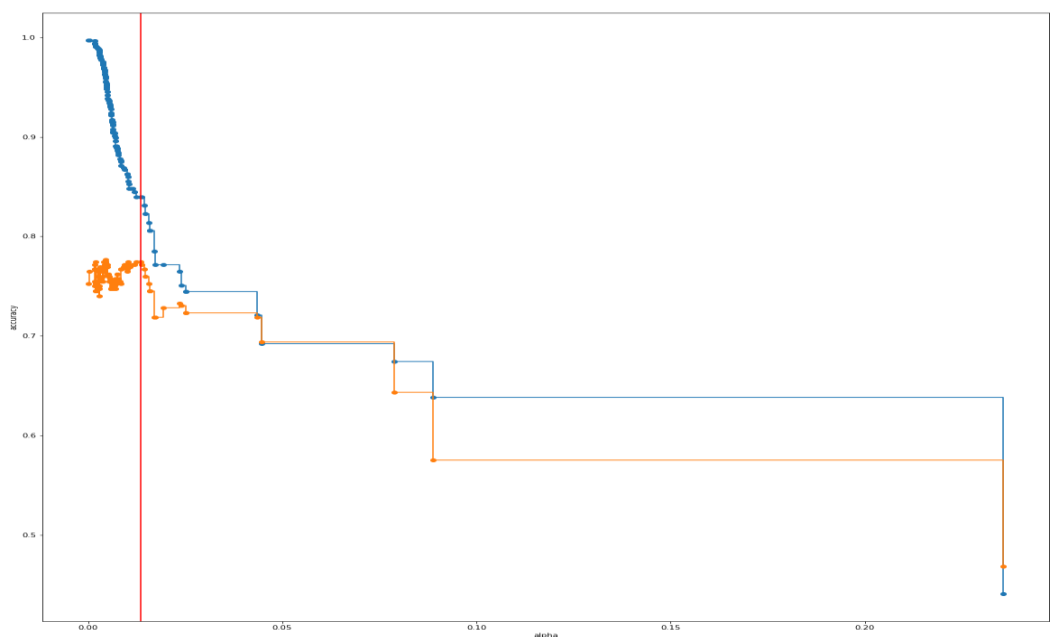 for the sake of computational speed I will just use the best parameters to do GrindSearchCV in my python notebook***

For the parameter cv in GrideSearchCV, I will set cv = 10 to do a 10 fold stratified cross validation on my model.

I will do the search in the function selection(tree,x,y) which inputs the decision tree classifier tree, data x and label data y. This function should return a dictionary of best performing parameters.

## **Pruning**

After selecting the best parameters, I will use it to construct a decision tree model. The problem of overfitting may occur despite the best parameters. To reduce overfitting, I will use cost complexity pruning method to post prune my decision tree model.

Cost complexity pruning method involves a cost complexity pruning parameter, ccp_alpha. By adjusting this parameter, we can get different scores on test sets and training sets. The best value for this parameter is the point where the testing data get the best score while training data gets the lowest score given the highest test score. Details of finding alpha is shown below



From the graph, best cc_alpha value is 0.0135.(This file is also attached as prun_parameter.png)

Now I will explain my procedure in python notebook.

First, I will split the dataset with 80:20 training to test ratio. Then I will use cost_complexity_pruning_path function to find the alphas (cost complexity pruning parameter) and corresponding impurities. Then, similar to hyperparameter tuning, I will use a loop to construct and score different decision tree models with different cost

complexity parameters. Then I will plot a graph to see what value of cost complexity parameters gives the best performance in testing set with lowest performance in training set. This cost complexity parameter will be used to construct the random forest model.
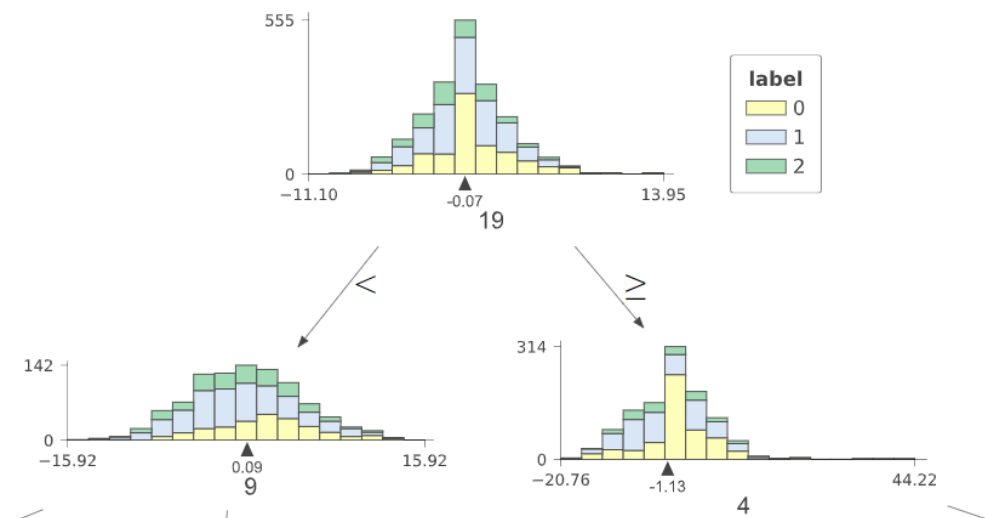
## K-Fold Cross Validation

I will use cross_val_score function to do a 10-fold stratified cross validation on my model. The scores should represent the accuracy of my model. As there should be 10 different scores, I will take the average to represent the overall score of my model.
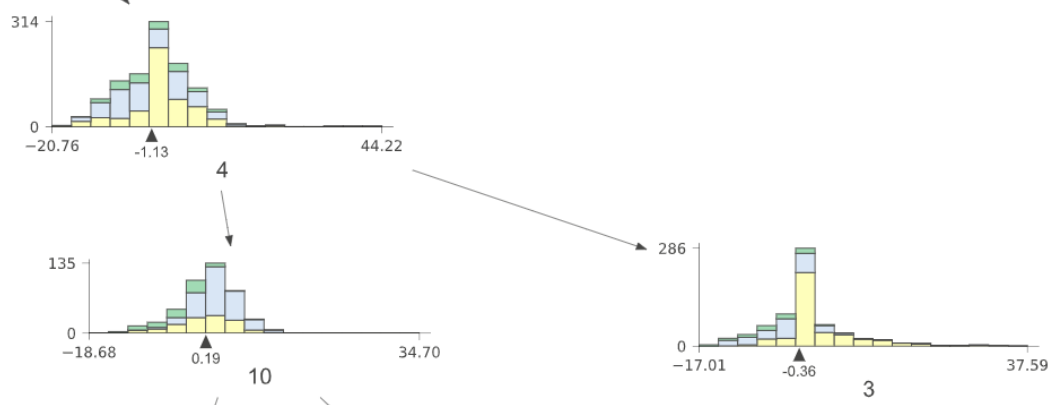
## Three splits from my decision tree

As random forest technique is used, I will choose three splits from the first tree in my forest to display.

1<sup>st</sup> split



2<sup>nd</sup> split



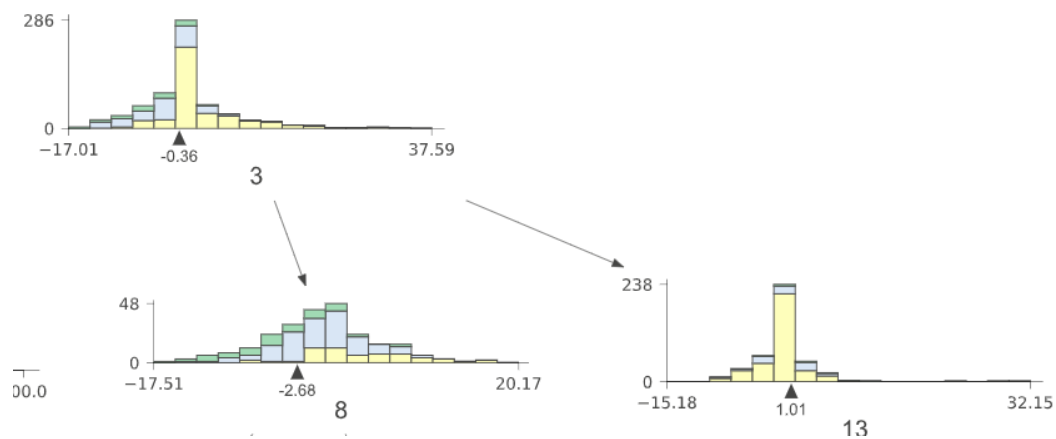3<sup>rd</sup> split

## Final Model

The final output is also attached as decision_tree.png