# Question 3 : Decision Trees

We have performed a grid search on the **_max_depth_** parameter of the decision trees in the range **[1,100)**.

## Accuracy Scores for different parameter settings:

a. **Abalone dataset:**
Raw dataset: The best accuracy of 26.26%  is obtained at a maximum depth of 4
PCA dataset: The best accuracy of 26.33%  is obtained at a maximum depth of 4
LDA dataset: The best accuracy of 25.8%  is obtained at a maximum depth of 3

b. **Wine dataset:**
Raw dataset: The best accuracy of 49.97%  is obtained at a maximum depth of 4
PCA dataset: The best accuracy of 45.57%  is obtained at a maximum depth of 6
LDA dataset: The best accuracy of 53.19%  is obtained at a maximum depth of 1

## Observations :

**Fig 3.1 to 3.6** demonstrate the **Mean Test Accuracy score** vs **Maximum Depth** plots for the 6 datasets. It can be observed from the Accuracy vs Depth plot for all the 6 datasets that the performance of the model first improves with the increase in the maximum depth but degrades rapidly after achieving the peak accuracy score. Thereafter, the accuracy seems to fluctuate in the lower range. This indicates that overfitting may be occurring with the increase in depth. Very low values of depth (in the range of 1-10) seems to be the best choice for all the datasets.
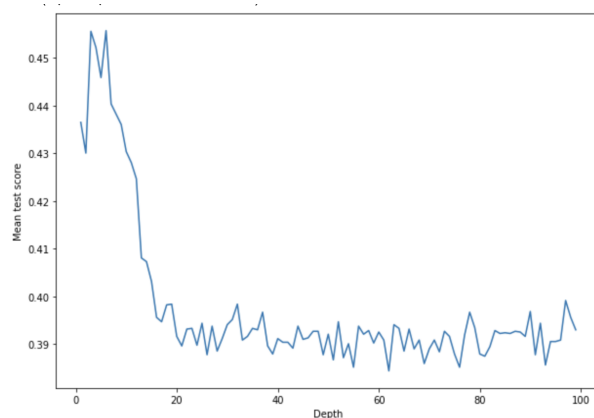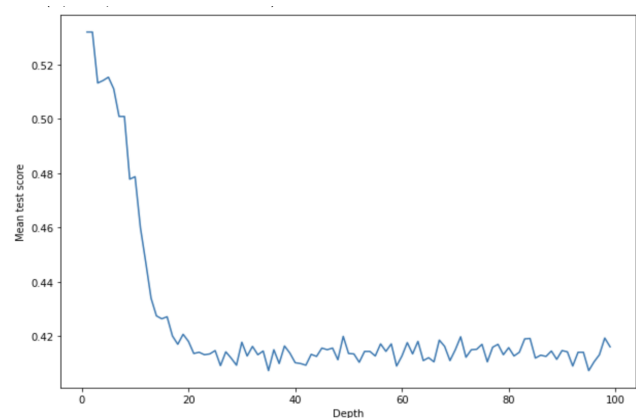
**Fig 3.1: Wine - PCA**

**Fig 3.2: Wine-LDA**
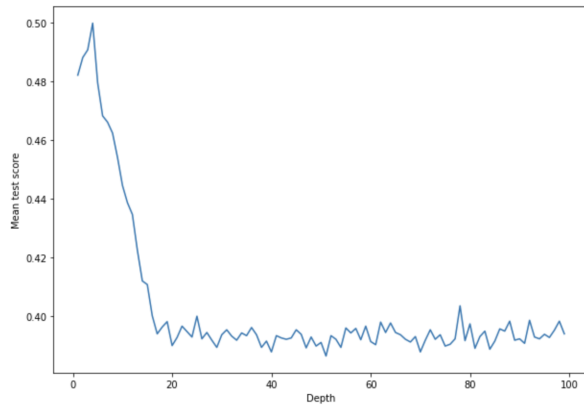


**Fig 3.3: Wine-raw**

**Fig 3.4: Abalone-PCA**

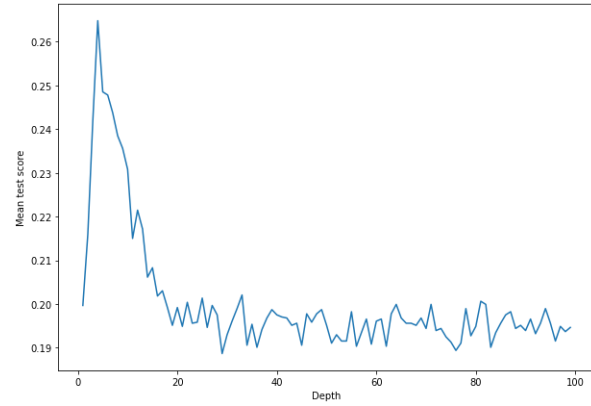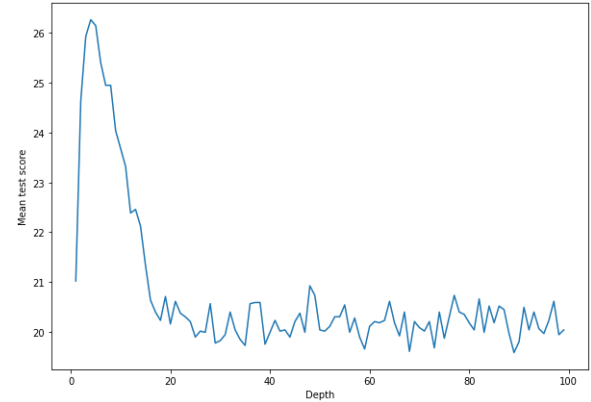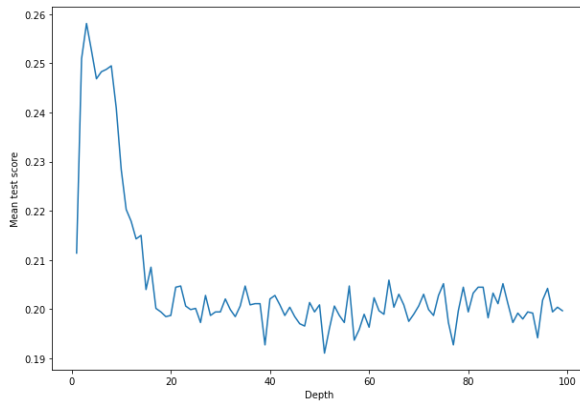**Fig 3.5: Abalone-LDA**



**Fig 3.6:Abalone-raw**





## Interpretability:

### a. Abalone dataset:

**Fig 3.5(a) and (b)** below show the best decision tree for abalone and its root node respectively. It shows the first partitioning feature used is x[6] i.e *'Shell Weight'* which means this feature provided the most information gain and was the most discriminating feature. It should be noted that *'Shell Weight'* was also observed to have a high correlation with the target variable i.e *'Rings'* as per the analysis performed in Assignment1 and hence corroborates the finding that it contains high information gain. Also, some features like x[2] i.e *'Height'* were not included in making any split decision.



```
            x[6] <= -0.678
            gini = 0.895
           samples = 4177
value = [1, 1, 15, 57, 115, 259, 391, 568, 689, 634, 487
    267, 203, 126, 103, 67, 58, 42, 32, 26, 14, 6, 9
             2, 1, 1, 2, 1]
```

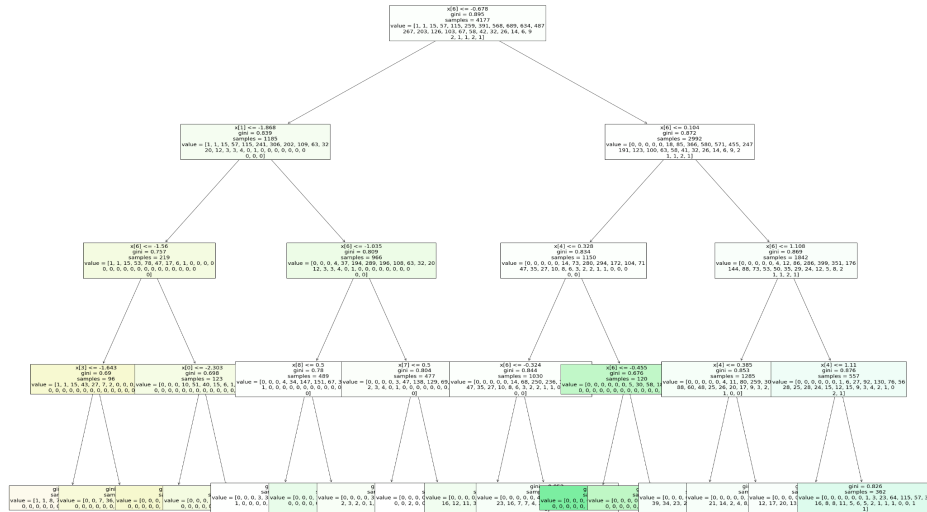**Fig 3.7 (b): Root Node(abalone-raw)**

Fig 3.7(a) : Best performing Decision Tree on Abalone Dataset

## b.  Wine dataset:

**Fig 3.8 (a) and (b)** below show the best decision tree for wine and its root node respectively. It shows the partitioning criterion and the feature used to partition it. As we can see, the topmost partition feature is x[10] i.e. **"Alcohol"** column in the dataset, which means this feature provided the most information gain and was the most discriminating feature. Another interesting observation is that some features, namely - **"Citric Acid"**, **"Total Sulfur Dioxide Density", "pH", "Type"**, weren't even used once in the whole decision tree for splitting, which means those features contributed the least in the decision making process and might be the most irrelevant features. Removal of such features could result in better performance.
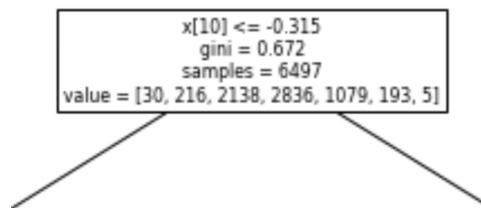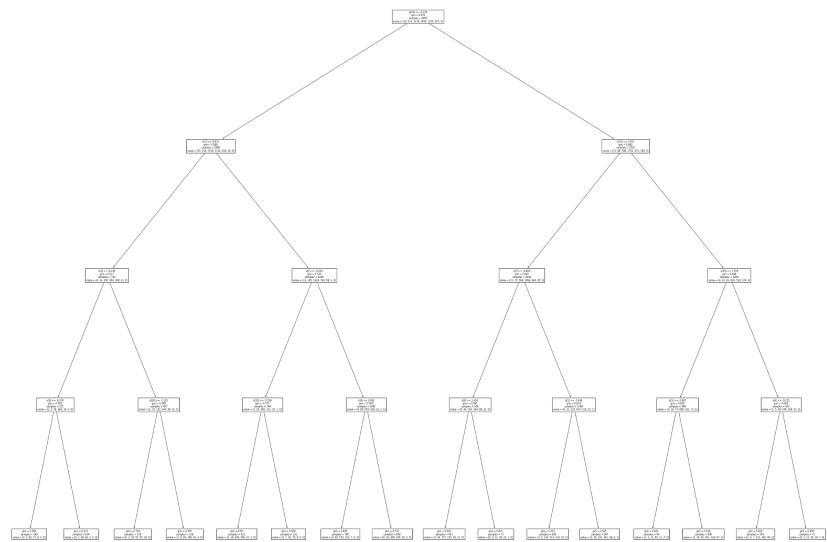


Fig 3.8 (b): Root Node(wine-raw)

**Fig 3.7(a) : Best performing Decision Tree on Wine Dataset**