

Q.3: Decision Tree

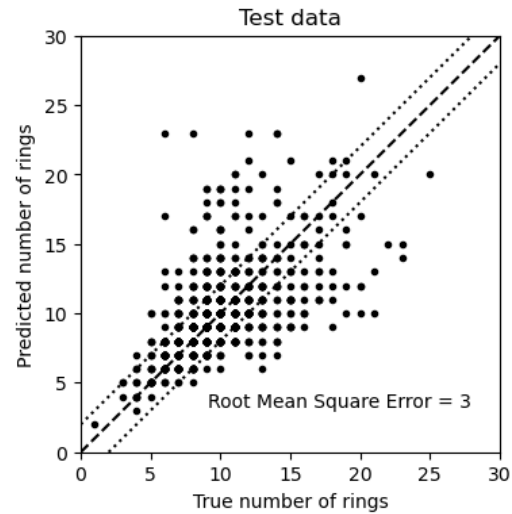
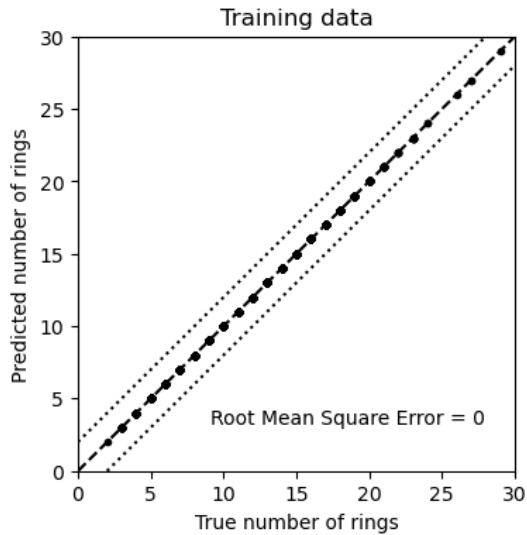
Decision Tree on Abalone dataset

The **DecisionTreeRegressor** is an algorithm used to estimate a continuous variable instead of a discrete one.

Testing score: [-0.05001266528192927, 0.20131855871953352, 0.06866661803868834, 0.14899358230388193, 0.12614396022862273]

Training score: [1.0, 1.0, 1.0, 1.0, 1.0]

This model **overfits** the dataset and that is why, validation error is very high.



The Decision Tree overfits the training set, i.e. its parameters are fine tuned to reproduce the results of the training set but generalized badly to data not seen previously.

GridSearchCV on RAW Abalone data

Using decision tree, with `max_depth` of 4, we got the accuracy of 0.26238260321462337

printing the decision tree using graphviz for Raw abalone data

```
graph TD
    Node0[Shell_weight <= 0.14]
    Node1[Shell_weight <= 0.02]
    Node2[Whole_weight <= 0.02]
    Node3[class: 3]
    Node4[Whole_weight > 0.02]
    Node5[class: 4]
    Node6[Shell_weight > 0.02]
    Node7[Length <= 0.25]
    Node8[class: 4]
    Node9[Length > 0.25]
    Node10[class: 5]
    Node11[Diameter > 0.22]
    Node12[Shell_weight <= 0.09]
    Node13[Sex <= 1.50]
    Node14[class: 7]
    Node15[Sex > 1.50]
    Node16[class: 9]
    Node17[Shell_weight > 0.09]
    Node18[Sex <= 0.50]
    Node19[class: 8]
    Node20[Sex > 0.50]
    Node21[class: 7]
    Node22[Shell_weight > 0.14]
    Node23[Shell_weight <= 0.25]
    Node24[Sucked_weight <= 0.43]
    Node25[Shell_weight <= 0.19]
    Node26[class: 8]
    Node27[Shell_weight > 0.19]
    Node28[class: 9]
    Node29[Sucked_weight > 0.43]
    Node30[Shell_weight <= 0.18]
    Node31[class: 10]
    Node32[Shell_weight > 0.18]
    Node33[class: 9]
    Node34[Shell_weight > 0.25]
    Node35[Shell_weight <= 0.39]
    Node36[Sucked_weight <= 0.44]
    Node37[class: 10]
    Node38[Sucked_weight > 0.44]
    Node39[class: 10]
    Node40[Shell_weight > 0.39]
    Node41[Sucked_weight <= 0.61]
    Node42[class: 10]
    Node43[Sucked_weight > 0.61]
    Node44[class: 11]

    Node0 --> Node1
    Node0 --> Node11
    Node1 --> Node2
    Node1 --> Node4
    Node2 --> Node3
    Node4 --> Node5
    Node11 --> Node12
    Node11 --> Node17
    Node12 --> Node13
    Node12 --> Node15
    Node13 --> Node14
    Node15 --> Node16
    Node17 --> Node18
    Node17 --> Node20
    Node18 --> Node19
    Node20 --> Node21
    Node22 --> Node23
    Node22 --> Node29
    Node23 --> Node24
    Node23 --> Node29
    Node24 --> Node25
    Node24 --> Node27
    Node25 --> Node26
    Node27 --> Node28
    Node29 --> Node30
    Node29 --> Node32
    Node30 --> Node31
    Node32 --> Node33
    Node34 --> Node35
    Node34 --> Node40
    Node35 --> Node36
    Node35 --> Node38
    Node36 --> Node37
    Node38 --> Node39
    Node40 --> Node41
    Node40 --> Node43
    Node41 --> Node42
    Node43 --> Node44
```

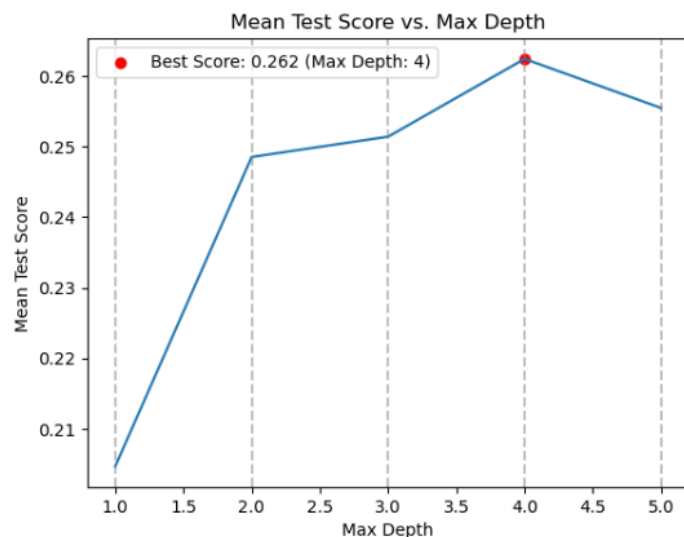
Summary of the decision tree

The decision tree for the abalone dataset has a total of 15 nodes and a maximum depth of 5. The first split is based on the "shell weight" feature, with a threshold of 0.14 g. If the shell weight is less than or equal to 0.14 g, the tree continues to split based on the "diameter" feature, with a threshold of 0.22 cm.

If the diameter is less than or equal to 0.22 cm, the tree further splits on "shell weight" and "whole weight" features. If the shell weight is less than or equal to 0.02 g and the whole weight is less than or equal to 0.02 g, the predicted age of the abalone is 3 years. If the shell weight is less than or equal to 0.02 g and the whole weight is greater than 0.02 g, the predicted age is 4 years. If the shell weight is greater than 0.02 g and the length is less than or equal to 0.25 cm, the predicted age is 4 years. If the length is greater than 0.25 cm, the predicted age is 5 years.

further tree continues to use the feature's threshold and classifies the abalones based on their ring classes.

Plotting the Max depth vs Mean test score for RAW abalone dataset



[illegible]

Tuning the hyperparameter and finding the best hyper-parameter that maximizes the accuracy for Raw abalone dataset

Best parameters: {'max_depth': 5, 'max_features': 'log2', 'max_leaf_nodes': 20, 'min_samples_leaf': 4, 'min_weight_fraction_leaf': 0.0, 'splitter': 'best'}

Best accuracy score: 0.267408818726184

Best parameters: {'max_depth': 3}

Best accuracy score: 0.2542440477895883

printing the decision tree using graphviz for PCA Pre-processed abalone data and the Max depth vs Mean test score for PCA pre-processed abalone dataset

Mean Test Score vs. Max Depth

Max Depth	Mean Test Score
1.0	0.205
2.0	0.230
3.0	0.254
4.0	0.254
5.0	0.248

Best Score: 0.254 (Max Depth: 3)

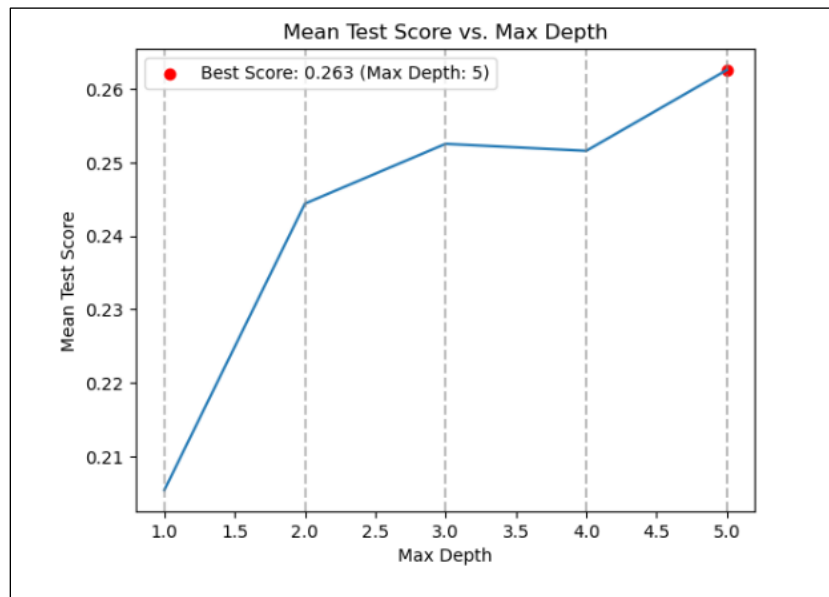
Best parameters: {'max_depth': 5}
Best accuracy score: 0.2626232702059995

Printing the decision tree using graphviz using LDA preprocessed Abalone dataset and the Max depth vs Mean test score for LDA pre-processed abalone dataset

```

--- LDA1 <= -0.60
--- LDA1 <= -3.28
--- LDA1 <= -4.13
--- LDA2 <= -2.77
--- LDA2 <= -3.89
--- class: 1
--- LDA2 > -3.89
--- class: 3
--- LDA2 > -2.77
--- LDA3 <= -0.49
--- class: 4
--- LDA3 > -0.49
--- class: 4
--- LDA1 > -4.13
--- LDA3 <= -0.17
--- LDA1 <= -3.63
--- class: 5
--- LDA1 > -3.63
--- class: 6
--- LDA3 > -0.17
--- LDA1 <= -3.64
--- class: 7
--- LDA1 > -3.64
--- class: 5
--- LDA1 > -3.28
--- LDA1 <= -1.45
--- LDA1 <= -2.01
--- LDA1 <= -3.25
--- class: 7
--- LDA1 > -3.25
--- class: 6
--- LDA1 > -2.01
--- LDA2 <= -0.08
--- class: 7
--- LDA2 > -0.08
--- class: 7
--- LDA1 > -1.45
--- LDA2 <= 0.28
--- LDA3 <= -1.12
--- class: 9
--- LDA3 > -1.12
--- class: 8
--- LDA2 > 0.28
--- LDA3 <= -1.32
--- class: 9
--- LDA3 > -1.32
--- class: 7
--- LDA1 > -0.60
--- LDA2 <= 0.09
--- LDA2 <= -1.45
--- LDA3 <= -0.02
--- LDA2 <= -5.73
--- class: 17
--- LDA2 > -5.73
--- class: 11
--- LDA3 > -0.02
--- LDA1 <= 1.20
--- class: 13
--- LDA1 > 1.20
--- class: 16
--- LDA2 > 1.45
--- LDA3 <= -0.31
--- LDA1 <= 0.57
--- class: 10
--- LDA1 > 0.57
--- class: 11
--- LDA3 > -0.31
--- LDA2 <= -0.19
--- class: 13
--- LDA2 > -0.19
--- class: 10
--- LDA2 > 0.09
--- LDA1 <= 0.23
--- LDA2 <= 0.63
--- LDA2 <= 0.58
--- class: 8
--- LDA2 > 0.58
--- class: 10
--- LDA2 > 0.63
--- LDA3 <= 0.23
--- class: 9

```



Tuning the hyperparameter and finding the best hyper-parameter that maximizes the accuracy for LDA pre-processed abalone dataset

Best parameters: {'max_depth': 5,
'max_features': 'log2',
'max_leaf_nodes': 20,
'min_samples_leaf': 5,
'min_weight_fraction_leaf': 0.0,
'splitter': 'best'}

Best accuracy score: 0.2585531329685128

Implementation of Decision trees on wine dataset starts here

Using Raw Wine Data

Decision tree accuracy and best hyperparameter using raw wine data

Best parameters: {'max_depth': 5, 'max_features': 'log2', 'max_leaf_nodes': 20, 'min_samples_leaf': 4, 'min_weight_fraction_leaf': 0.0, 'splitter': 'best'}

Best accuracy score: 0.5320946290045596

Decision tree accuracy and best hyperparameter using PCA pre-processed wine data

Best parameters: {'max_depth': 5, 'max_features': 'log2', 'max_leaf_nodes': 20, 'min_samples_leaf': 1, 'min_weight_fraction_leaf': 0.0, 'splitter': 'best'}

Best accuracy score: 0.5034671640907208

Decision tree accuracy and best hyperparameter using LDA pre-processed wine data

Best parameters: {'max_depth': 5, 'max_features': 'log2', 'max_leaf_nodes': None, 'min_samples_leaf': 2, 'min_weight_fraction_leaf': 0.0, 'splitter': 'best'}

Best accuracy score: 0.5460980635992183

Conclusion:

LDA is performing better than PCA because the goal of the dataset is to classify the abalones into different age groups based on their physical characteristics. LDA takes into account the class information while PCA does not. Therefore, LDA is better suited for this classification problem.

Moreover, the abalone dataset has a low number of features compared to its sample size. This means that the dataset may not have a high degree of redundancy, which is necessary for PCA to work well. LDA is less affected by the degree of redundancy in the data because it explicitly takes into account the class information.

LDA is suited for the abalone dataset because it is specifically designed for classification problems and takes into account the class information, which is essential for the task of predicting the age of abalones based on their physical characteristics.