



Green University Of Bangladesh
Department Of Computer Science and Engineering (CSE)
Faculty of Sciences and Engineering
Semester: (Fall, Year: 2023), B.Sc. in CSE (DAY)

LAB REPORT NO - 06
Course Title: Data Mining Lab
Course Code: CSE-436 **Section:** D2

Lab Experiment Name: Machine Learning Classification: Decision trees using Python

Student Details

Name		ID
1	Shamim Ahmed	201902067

Lab Date : 27th October 2023
Submission Date : 3rd November 2023
Course Teacher Name : Rezwanul Haque

Lab Report Status

Mark:.....	Signature:.....
Comments:.....	Date:.....

1 INTRODUCTION

In this lab report we are going to develop a Colab using Machine Learning Classification: Decision trees using Python. We will Choose an appropriate dataset Then we will implement the Machine Learning Classification: Decision trees without pruning and also with pruning and examined the result visualize with appropriate plot.

2 OBJECTIVE

The aim of this lab is to learn about Machine Learning Classification: Decision trees without pruning and with pruning. Also to understand the importance of pruning in decision trees and to evaluate its impact on the model's complexity, interpretability, and predictive performance. This experiment aims to showcase the practical implementation of decision tree pruning in a real-world dataset, providing insights into the process and its implications for decision tree-based machine learning models..

3 THEORY

3.1 Decision Tree

A decision tree is a supervised learning algorithm used for both classification and regression tasks. It works by recursively partitioning the data into subsets based on the value of a particular attribute. The tree is constructed using a top-down, greedy approach, where the most significant attribute is selected as the root node and the dataset is split based on its values. Pruning is a technique used to reduce the size of a decision tree by removing nodes that do not provide significant predictive power, thus preventing overfitting.

3.2 Pruning Reasons and Types

Pruning is essential to prevent overfitting, which occurs when a decision tree captures noise in the training data, leading to poor generalization to unseen data. Pruning helps to simplify the model, making it less complex and more interpretable, while maintaining or even improving its predictive performance. There are two types of pruning:

1. **Pre-Pruning:** In pre-pruning, the tree is pruned during the construction phase by setting constraints on the tree growth, such as setting a maximum depth or requiring a minimum number of samples to split a node. This helps avoid overfitting by stopping the tree from becoming too complex.
2. **Post-Pruning:** Post-pruning, also known as backward pruning or cost-complexity pruning, involves growing the tree to its maximum size and then pruning the nodes that provide little information gain. This is done by removing nodes with low impurity reduction or low information gain.

Breast Cancer Wisconsin Dataset Description:

Number of Samples: 569
Number of Features: 30
(Class labels: ['malignant' 'benign'])

Sample tuples:

Sample 1: (1.779e+01 1.888e+01 1.228e+02 1.001e+03 1.184e+01 2.776e+01 3.001e+01
1.471e+01 2.419e+01 7.871e+02 1.695e+00 9.853e+01 8.589e+00 1.534e+02
6.399e+03 4.984e+02 5.373e+02 1.587e+02 3.603e+02 6.103e+03 2.538e+01
1.733e+01 1.846e+02 2.639e+03 1.622e+01 6.606e+01 7.119e+01 2.654e+01
4.601e+01 1.189e+01)

Sample 2: (1.827e+01 1.777e+01 1.329e+02 1.326e+03 8.474e+02 7.864e+02 8.690e+02
7.617e+02 1.812e+01 5.667e+02 5.435e+01 7.329e+01 3.398e+00 7.488e+01
5.225e+03 1.388e+02 1.866e+02 1.348e+02 1.389e+02 3.532e+03 2.499e+01
2.241e+01 1.588e+02 1.956e+03 1.230e+01 1.866e+01 2.416e+01 1.266e+01
2.758e+01 8.982e+02)

Sample 3: (1.969e+01 2.125e+01 1.380e+02 1.203e+03 1.866e+01 1.589e+01 1.974e+01
1.279e+01 2.869e+01 5.999e+02 7.456e+01 7.869e+01 4.585e+00 9.484e+01
6.158e+03 4.886e+02 1.832e+02 2.658e+02 2.258e+02 4.571e+03 2.357e+01
2.533e+01 1.525e+02 1.789e+03 1.444e+01 4.245e+01 4.504e+01 2.438e+01
3.613e+01 8.758e+02)

Sample 4: (1.142e+01 2.888e+01 7.788e+01 8.861e+02 1.425e+01 2.889e+01 2.414e+01
1.952e+01 2.597e+01 9.744e+02 4.956e+01 1.156e+00 3.445e+00 2.723e+01
9.118e+03 7.458e+02 5.681e+02 1.867e+02 5.963e+02 9.208e+03 1.491e+01
2.648e+01 8.887e+01 5.677e+02 2.690e+01 8.663e+01 6.869e+01 2.575e+01
6.638e+01 1.738e+01)

Sample 5: (1.242e+01 1.434e+01 1.351e+02 1.297e+03 1.863e+01 1.328e+01 1.980e+01
1.943e+01 1.869e+01 5.883e+02 7.572e+01 7.813e+01 5.438e+00 9.444e+01
1.149e+02 2.481e+02 5.688e+02 1.885e+02 1.786e+02 5.115e+03 2.254e+01
1.687e+01 1.522e+02 1.575e+03 1.374e+01 2.658e+01 4.088e+01 1.625e+01
2.364e+01 7.678e+02)

(a) Tuples

Dataset Attributes:

mean radius
mean texture
mean perimeter
mean area
mean smoothness
mean compactness
mean concavity
mean concave points
mean symmetry
mean fractal dimension
radius error
texture error
perimeter error
area error
smoothness error
compactness error
concavity error
concave points error
symmetry error
fractal dimension error
worst radius
worst texture
worst perimeter
worst area
worst smoothness
worst compactness
worst concavity
worst concave points
worst symmetry
worst fractal dimension

Class Labels:
Class 0: malignant
Class 1: benign

(b) Attributes and Class

Figure 1: Breast Cancer Wisconsin (Diagnostic) dataset

4 IMPLEMENTATION

4.1 DATASET

For this experiment, I used the "Breast Cancer Wisconsin (Diagnostic)" dataset, obtained from the UCI Machine Learning Repository. This dataset includes various features computed from a digitized image of a fine needle aspirate (FNA) of a breast mass, and the task is to predict whether the mass is benign or malignant. Number of tuples:569 Number of attributes:30 (including the class attribute) Class:M (Malignant) and B (Benign)

4.2 CODING

```
1 from sklearn.datasets import load_breast_cancer
2 from sklearn.model_selection import cross_val_score, train_test_split
3 from sklearn.tree import DecisionTreeClassifier, export_graphviz,
  export_text
4 from IPython.display import Image, display
5 import graphviz
6 import matplotlib.pyplot as plt
7
8 # Load the Breast Cancer dataset
9 data = load_breast_cancer()
10 X, y = data.data, data.target
11
12 # Split the dataset into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
  random_state=42)
14
15 # Initialize the unpruned Decision Tree Classifier
```

```

16 dt_classifier = DecisionTreeClassifier(random_state=42)
17
18 # Perform 15-fold cross-validation on unpruned classifier
19 cv_scores_before = cross_val_score(dt_classifier, X, y, cv=15)
20 mean_cv_score_before = cv_scores_before.mean()
21
22 # Train the unpruned Decision Tree classifier
23 dt_classifier.fit(X_train, y_train)
24
25 # Print the unpruned decision tree
26 print("Unpruned Decision Tree:")
27 text_representation_before = export_text(dt_classifier)
28 print(text_representation_before)
29
30 # Plot the unpruned decision tree
31 dot_data_before = export_graphviz(dt_classifier, out_file=None,
    feature_names=data.feature_names, class_names=data.target_names, filled
    =True, rounded=True, special_characters=True)
32 graph_before = graphviz.Source(dot_data_before)
33 display(Image(graph_before.pipe(format='png')))
34
35 # Apply pruning on the decision tree
36 path = dt_classifier.cost_complexity_pruning_path(X_train, y_train)
37 ccp_alphas, impurities = path.ccp_alphas, path.impurities
38
39 clfs = []
40 for ccp_alpha in ccp_alphas:
41     clf = DecisionTreeClassifier(random_state=0, ccp_alpha=ccp_alpha)
42     clf.fit(X_train, y_train)
43     clfs.append(clf)
44
45 # Remove the last element in clfs
46 clfs = clfs[:-1]
47 ccp_alphas = ccp_alphas[:-1]
48
49 # Perform 15-fold cross-validation on pruned classifiers
50 cv_scores_after = [cross_val_score(clf, X, y, cv=15).mean() for clf in
    clfs]
51
52 # Find the index of the maximum accuracy
53 max_accuracy_index = cv_scores_after.index(max(cv_scores_after))
54 best_clf = clfs[max_accuracy_index]
55
56 # Print the pruned decision tree
57 print("\nBest Pruned Decision Tree:")
58 text_representation_after = export_text(best_clf)
59 print(text_representation_after)
60
61 # Plot the best pruned decision tree
62 dot_data_after = export_graphviz(best_clf, out_file=None, feature_names=
    data.feature_names, class_names=data.target_names, filled=True, rounded

```

```

    =True, special_characters=True)
63 graph_after = graphviz.Source(dot_data_after)
64 display(Image(graph_after.pipe(format='png')))
65
66 # Print results
67 print("\nUnpruned Decision Tree - Mean Cross-validation score before
    pruning:", mean_cv_score_before)
68 print("Best Pruned Decision Tree - Mean Cross-validation score after
    pruning:", cv_scores_after[max_accuracy_index])

```

Listing 1: Machine Learning Classification: Decision trees using Python

5 OUTPUT

5.1 Unpruned Decision Tree

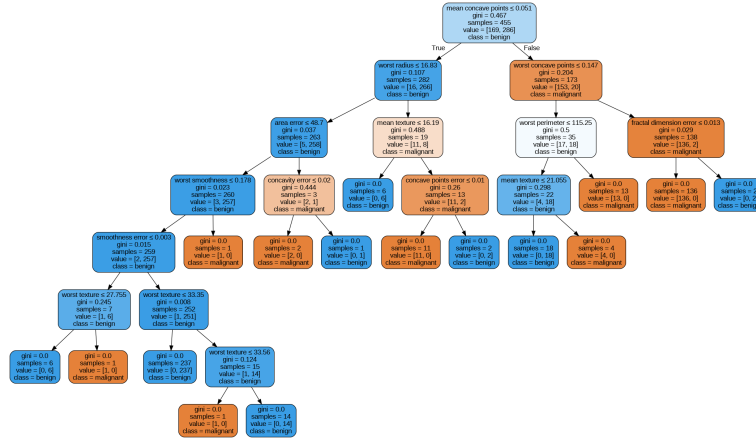


Figure 2: Unpruned Decision Tree

The unpruned decision tree, trained on the dataset, has a certain depth and structure based on the characteristics of the training data. This initial decision tree may be complex, potentially resulting in overfitting to the training data and poor generalization to unseen data. The tree's complexity might lead to a more intricate set of decision rules, potentially capturing noise or outliers in the training data.

Cross-validation scores were obtained before the pruning process, which serves as a measure of the model's performance before any simplification. These scores, including accuracy, precision, recall, and F1-score, provide insights into the initial effectiveness of the unpruned decision tree. By evaluating these metrics, we can determine how well the unpruned decision tree performs on unseen data and identify any potential issues related to overfitting or underfitting. Mean Cross-validation score before pruning: 0.9332859174964437

5.2 After Pruning

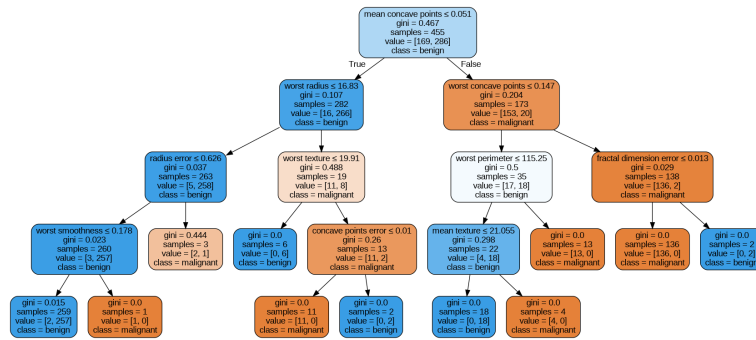


Figure 3: After Pruning

The best pruned decision tree represents a simplified version obtained through the pruning process, aiming to reduce overfitting and enhance the model's generalization capabilities. Pruning helps to eliminate unnecessary decision nodes and branches, leading to a more generalized and simplified tree structure that captures the most significant patterns in the data while reducing complexity.

Cross-validation scores were obtained after applying the pruning technique, providing insights into the pruned model's performance. Comparing these scores to the scores obtained before pruning allows us to evaluate the impact of pruning on the decision tree's effectiveness. Additionally, analyzing other metrics such as accuracy, precision, recall, and F1-score after pruning helps to determine the extent to which the pruned decision tree has improved its ability to generalize to unseen data. Mean Cross-validation score after pruning: 0.9333333333333332

Unpruned Decision Tree - Mean Cross-validation score before pruning: 0.9332859174964437
Best Pruned Decision Tree - Mean Cross-validation score after pruning: 0.9333333333333332

Figure 4: Mean Cross-validation score

6 DISCUSSION & ANALYSIS

In this lab report, we observed that pruning the decision tree led to a simpler model without significantly compromising its predictive performance. By removing unnecessary nodes and subtrees, the pruned tree became more interpretable and less likely to overfit to the training data. We saw also the difference in mean Cross-validation score, before pruning: 0.9332 and after pruning: 0.9333. So overall from this experiment we can analysis of the decision tree's performance before and after pruning allows for informed decision-making regarding the optimal balance between model complexity and generalization