# Assignment-4: Decision Trees and XGBoost

Sajan Kumer Sarker
Id: 2111131642
Email: sajan.sarker@northsouth.edu

Rafsan Jani Chowdhury
Id: 2011424642
Email: rafsan.chowdhury@northsouth.edu

Rosely Mohammad
Id: 2014219642
Email: rosely.mohammad@northsouth.edu

November 9, 2024

## 1. INTRODUCTION

In this Assignment report on Decision Tree and XGBoost we will discuss about training and evaluation of the Decision Tree model and XGBoost model on four different datasets: 'Car Evaluation (for Multi-class Classification)', 'DT-BrainCancer (for Binary Classification)', 'DT-Wage (for Regression)', 'DT-Credits (for Regression)'. For these models we'll report Accuracy, Precision, Average Precision, Recall, Average Recall, Confusion Matrix, F-Score, Precision-Recall Curve and Mean-Squared Error. Here, we split this assignment into 4 parts: 1. Multi-class Classification on 'Car Evaluation' dataset, 2. Binary Classification on 'DT-BrainCancer' dataset, 3. Regression on 'DT-Wage' dataset, 4. Regression on 'DT-Credits' dataset. Each part is handled in a separate notebook file, with each file focused on a single dataset. To implement these models we've follow multiples resources. [1] [2] [3] [4]

## 2. METHODS

**Necessary Imports:**

```
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import math
5  import seaborn as sns
6  import warnings
7  warnings.filterwarnings('ignore')
8  %matplotlib inline
9  from sklearn.model_selection import train_test_split, GridSearchCV
10 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
11 from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
12 from xgboost import XGBClassifier, XGBRegressor
```

We have imported all necessary libraries and modules to be used in each notebook.

## 2.1 Multi-class Classification on 'Car Evaluation' dataset

**Loading Dataset:**

```
1  car_evaluation = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/Assignment-4/dataset/ca
2
3  car_evaluation.head()
```

Load the Dataset using pandas library and check first few data from the dataset.

**Assign Column name:**

```
1  column_names = ['buying', 'maint', 'doors', 'persons', 'lug_boot', 'safety', 'class']
2  car_evaluation.columns = column_names
3  car_evaluation.head()
```

Assign each columns with a name, since the original dataset didn't have any headings.

**Checking Missing Values:**

```
1  print('Number_of_missing_values:_')
2  print(car_evaluation.isnull().sum())
```

Checking if any missing values are founds in any features.

**Mapping Categorical Data:**

```
1  car_evaluation['buying'] = car_evaluation['buying'].map({'vhigh':3, 'high':2, 'med':1, 'low':0
2  car_evaluation['maint'] = car_evaluation['maint'].map({'vhigh':3, 'high':2, 'med':1, 'low':0})
3  car_evaluation['doors'] = car_evaluation['doors'].map({'2':0, '3':1, '4':2, '5more':3})
4  car_evaluation['persons'] = car_evaluation['persons'].map({'2':0, '4':1, 'more':2})
5  car_evaluation['lug_boot'] = car_evaluation['lug_boot'].map({'small':0, 'med':1, 'big': 2})
6  car_evaluation['safety'] = car_evaluation['safety'].map({'low':0, 'med':1, 'high': 2})
7
8  car_evaluation['class'] = car_evaluation['class'].map({'unacc':0, 'acc':1, 'good': 2, 'vgood':
9  car_evaluation.head()
```

Since there is no complex relationship between the categories, so we just mapped them to values ranging from 0 to 4. Then each categorical columns are encoded in 0 to 4.

**Split Data into Training, Validation and Test Sets:**

```
1  X = car_evaluation.drop(columns='class')
2  y = car_evaluation['class']
3  X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
4  X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42
5
6  print(f'Train_Data_Shape_(X,_y):_{X_train.shape,_y_train.shape}')
7  print(f'Validation_Data_Shape_(X,_y):_{X_val.shape,_y_val.shape}')
8  print(f'Test_Data_Shape_(X,_y):_{X_test.shape,_y_test.shape}')
```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing.

**Initialize the Classifier Models and Hyper-parameters for GridSearchCV:**

```
1  # initialize the classifier models
2  dt_classifier = DecisionTreeClassifier(random_state=42)
3  xgb_classifier = XGBClassifier(random_state=42)
4
```

```
5   # initialize hyperparameters for gridsearchCV
6   dt_param_grid = {
7       'max_depth': [2,4,8,10],
8       'min_samples_split': [2,4,5,10],
9       'min_samples_leaf': [1,3,5,8,10]
10  }
11
12  xgb_param_grid = {
13      'max_depth': [3,5,7],
14      'learning_rate': [0.1, 0.01, 0.001],
15      'n_estimators': [100, 200, 300],
16      'subsample': [0.8, 1.0]
17  }
```

Here, we initialize the Decision Tree and XGBoost classifier models. Then, we define the hyperparameter grids for each model to fine-tune the key parameters. For the Decision Tree classifier, we selected a specific range of values for max_depth, min_samples_split, and min_samples_leaf, allowing the model to explore different tree structures. We do the same for the XGBoost model, where we choose a specific range of values for max_depth, learning_rate, n_estimators, and subsample. [1]

**Training Decision Tree Model:**

```
1   # Decision Tree:
2   dt_grid_search = GridSearchCV(estimator=dt_classifier, param_grid=dt_param_grid, cv=5, scoring=
3   dt_grid_search.fit(X_train, y_train)
4
5   dt_best_params = dt_grid_search.best_params_
6   dt_best_score = dt_grid_search.best_score_
7   print(f"Decision Tree Best Parameters: {dt_best_params}")
8   print(f"Decision Tree Best Score: {dt_best_score}")
```

Here, we perform hyperparameter tuning for the Decision Tree model using GridSearchCV by testing different combinations of parameters with 5-fold cross-validation to evaluate the models performance. Then we fit it to the training data and find the best parameters and best scores. [1]

**Training XGBoost Model:**

```
1   # XGBoost:
2   xgb_grid_search = GridSearchCV(estimator=xgb_classifier, param_grid=xgb_param_grid, cv=5, scor
3   xgb_grid_search.fit(X_train, y_train)
4
5   xgb_best_params = xgb_grid_search.best_params_
6   xgb_best_score = xgb_grid_search.best_score_
7   print(f"XGBoost Best Parameters: {xgb_best_params}")
8   print(f"XGBoost Best Score: {xgb_best_score}")
```

There we do the same thing for XGBoost model. We perform hyperparameter tuning with 5-fold cross-validation. Then fit it with training data and find the best parameters. [1]

**Retraining the Model with Validation and Training Set:**

```
1   # now retrain the model on the combined training and validation set with the best hyperparamet
2   X_train_val = np.concatenate((X_train, X_val))
3   y_train_val = np.concatenate((y_train, y_val))
4
5   dt_final_model = DecisionTreeClassifier(**dt_best_params, random_state=42)
6   xgb_final_model = XGBClassifier(**xgb_best_params, random_state=42)
7
```

```
8   # retrain decision tree final model
9   dt_final_model.fit(X_train_val, y_train_val)
10  print(f"Decision_Tree_Best_Score:_{dt_best_score}")
11
12  # retrain xgboost final model
13  xgb_final_model.fit(X_train_val, y_train_val)
14  print(f"XGBoost_Best_Score:_{xgb_best_score}")
```

Here we retrain both models with the combined training and vaalidation data to make the best use of available data with the optimal best hyper parameters.

**Evaluate the Final Model on Test Set:**

```
1   # evaluate the final model on test set
2   y_test_pred_dt = dt_final_model.predict(X_test)
3   y_test_pred_xgb = xgb_final_model.predict(X_test)
```

This step, we evaluate the final model's performance on test set.

**User Define Evaluation Metrics Function:**

```
1   # functions to calculate: accuracy, average precision, average recall, average f-score.
2
3   # Accuracy
4   def accuracy_score(y_actual, y_predict):
5     correct = 0
6     total_samples = len(y_actual)
7
8     for predict, actual in zip(y_predict, y_actual):
9       if predict == actual:
10        correct += 1
11
12    return (correct/ total_samples)
13
14  # Precision
15  def precision_score(y_actual, y_predict):
16    precision_scores = []
17
18    for class_label in set(y_actual):
19      true_positive = 0
20      false_positive = 0
21
22      for actual, predict in zip(y_actual, y_predict):
23        if predict == class_label:
24          if actual == class_label:
25            true_positive += 1
26          else:
27            false_positive += 1
28
29      if true_positive+false_positive ==0:
30        class_precision = 0.0
31      else:
32        class_precision = true_positive / (true_positive + false_positive)
33      precision_scores.append(class_precision)
34    return precision_scores
35
36  # average precision
37  def average_precision_score_macro(precision_score, y_actual):
```

```python
38      return sum(precision_score)/len(set(y_actual))
39
40  # recall
41  def recall_score(y_actual, y_predict):
42      recall_scores = []
43
44      for class_label in set(y_actual):
45          true_positive = 0
46          false_negative = 0
47
48          for actual, predict in zip(y_actual, y_predict):
49              if actual == class_label:
50                  if predict == class_label:
51                      true_positive += 1
52                  else:
53                      false_negative += 1
54
55          if true_positive+false_negative ==0:
56              class_recall = 0.0
57          else:
58              class_recall = true_positive / (true_positive + false_negative)
59
60          recall_scores.append(class_recall)
61      return recall_scores
62
63  # average recall
64  def average_recall_score_macro(recall_score, y_actual):
65      return sum(recall_score)/len(set(y_actual))
66
67  # f-score
68  def fscore(y_actual,y_predict):
69      f_score_total = []
70
71      for class_label in set(y_actual):
72          true_positive = 0
73          false_positive = 0
74          false_negative = 0
75
76          for actual, predict in zip(y_actual, y_predict):
77              if predict == class_label:
78                  if actual == class_label:
79                      true_positive += 1
80                  else:
81                      false_positive += 1
82              elif actual == class_label:
83                  false_negative += 1
84
85          if true_positive + false_positive == 0:
86              precision = 0.0
87          else:
88              precision = true_positive/(true_positive+false_positive)
89
90          if true_positive + false_negative == 0:
91              precision = 0.0
92          else:
93              recall = true_positive/(true_positive+false_negative)
94
95          if precision + recall ==0:
```

```
96        f_score = 0.0
97      else:
98        f_score = 2*true_positive/((2*true_positive)+false_positive+false_negative)
99
100       f_score_total.append(f_score)
101    return f_score_total
102
103  # fscore average
104  def fscore_average_macro(fscore, y_actual):
105    return sum(fscore)/len(set(y_actual))
```

Here we define some functions manually to report the evaluation matrices (f-score, precision, recall, average precision, average recall) without using evaluation matrices scikit learn's library function.

### Evaluation Metrics:

```
1  print("Decision␣Tree:␣")
2  print(f'Accuracy:␣{accuracy_score(y_test,␣y_test_pred_dt)}')
3  precision = precision_score(y_test, y_test_pred_dt)
4  print(f'Average␣Precision:␣{average_precision_score_macro(precision,␣y_test)}')
5  recall = recall_score(y_test, y_test_pred_dt)
6  print(f'Average␣Recall:␣{average_recall_score_macro(recall,␣y_test)}')
7  f_score = fscore(y_test, y_test_pred_dt)
8  print(f'Average␣F1-Score:␣{fscore_average_macro(f_score,␣y_test)}')
9  print()
10
11  print("XGBoost:␣")
12  print(f'Accuracy:␣{accuracy_score(y_test,␣y_test_pred_xgb)}')
13  precision = precision_score(y_test, y_test_pred_xgb)
14  print(f'Average␣Precision:␣{average_precision_score_macro(precision,␣y_test)}')
15  recall = recall_score(y_test, y_test_pred_xgb)
16  print(f'Average␣Recall:␣{average_recall_score_macro(recall,␣y_test)}')
17  f_score = fscore(y_test, y_test_pred_xgb)
18  print(f'Aveerage␣F1-Score:␣{fscore_average_macro(f_score,␣y_test)}')
```

We print the evaluation metrics using the defined functions.

### Confusion Matrix:

```
1  num_class = len(set(y_test))
2  fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(20, 8))
3
4  # confusion matrix for Decision Tree
5  conf_matrix_dt = np.zeros((num_class, num_class), dtype=int)
6  for actual, predict in zip(y_test, y_test_pred_dt):
7      conf_matrix_dt[actual, predict] += 1
8
9  # plotting
10  ax1.imshow(conf_matrix_dt, cmap='Blues', interpolation='nearest')
11  ax1.set_title('Decision␣Tree', fontsize=16)
12  ax1.set_xlabel('Predicted␣Label', fontsize=14)
13  ax1.set_ylabel('Actual␣Label', fontsize=14)
14  ax1.set_xticklabels(['unacc', 'acc', 'good', 'vgood'], fontsize=12)
15  ax1.set_yticklabels(['unacc', 'acc', 'good', 'vgood'], fontsize=12)
16
17  for i in range(num_class):
18      for j in range(num_class):
19          ax1.text(j, i, str(conf_matrix_dt[i, j]), ha='center', va='center', color='red', fonts
```

```
20
21 # confusion matrix for XGBoost
22 conf_matrix_xgb = np.zeros((num_class, num_class), dtype=int)
23 for actual, predict in zip(y_test, y_test_pred_xgb):
24     conf_matrix_xgb[actual, predict] += 1
25
26 # Plotting
27 ax2.imshow(conf_matrix_xgb, cmap='Blues', interpolation='nearest')
28 ax2.set_title('XGBoost', fontsize=16)
29 ax2.set_xlabel('Predicted Label', fontsize=14)
30 ax2.set_ylabel('Actual Label', fontsize=14)
31 ax2.set_xticklabels(['unacc', 'acc', 'good', 'vgood'], fontsize=12)
32 ax2.set_yticklabels(['unacc', 'acc', 'good', 'vgood'], fontsize=12)
33
34 for i in range(num_class):
35     for j in range(num_class):
36         ax2.text(j, i, str(conf_matrix_xgb[i, j]), ha='center', va='center', color='red', font
37
38 # common properties
39 for ax in (ax1, ax2):
40     ax.set_xticks(np.arange(num_class))
41     ax.set_yticks(np.arange(num_class))
42     ax.tick_params(axis='both', which='major', labelsize=12)
43
44 plt.xticks([0, 1, 2, 3], ['unacc', 'acc', 'goood', 'vgood'])
45 plt.yticks([0, 1, 2, 3], ['unacc', 'acc', 'goood', 'vgood'])
46 plt.suptitle('Confusion Matrices for Decision Tree and XGBoost', fontsize=18)
47 plt.show()
```

There we plot the Confusion Matrix for both model based on the predictions and actual values.

## 2.2 Binary Classification on 'DT-BrainCancer' dataset

**Loading Dataset:**

```
1 df = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/Assignment-4/dataset/DT-BrainCancer
2 print(df.shape)
3 print(df.head())
```

Load the Dataset using pandas library and check first few data from the dataset and its Shape.

**Checking Missing Values:**

```
1 print('Number of missing values: ')
2 print(df.isnull().sum())
```

Checking if any missing values are founds in any features.

**Drop Column:**

```
1 df = df.drop(columns=['Unnamed: 0', 'sex'])
2 df = df.dropna()
3 print(df.shape)
4 print(df.head())
```

Drop the unwanted column and remove the missing values row.

**Encoded Categorical Data:**

```
1  le = LabelEncoder()
2  df['diagnosis'] = le.fit_transform(df['diagnosis'])
3  df['loc'] = le.fit_transform(df['loc'])
4  print(df.head(), df.tail())
```

Since there has few complex relationship between the categories, so we encoded them using LabelEnocer. Then we print few samples from the first and few from the last of this encoded dataset.

**Split Data into Training, Validation and Test Sets:**

```
1  X = df.drop(columns=['status'])
2  y = df['status']
3
4  X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
5  X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42
6
7  print(f'Train_Data_Shape_(X,_y):_{X_train.shape,_y_train.shape}')
8  print(f'Validation_Data_Shape_(X,_y):_{X_val.shape,_y_val.shape}')
9  print(f'Test_Data_Shape_(X,_y):_{X_test.shape,_y_test.shape}')
```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing.

**Initialize the Classifier Models and Hyper-parameters for GridSearchCV:**

```
1   # initialize the classifier models
2   dt_classifier = DecisionTreeClassifier(random_state=42)
3   xgb_classifier = XGBClassifier(random_state=42)
4
5   # initialize hyperparameters for gridsearchCV
6   dt_param_grid = {
7       'max_depth': [2,3,4,5,8,10,12,15],
8       'min_samples_split': [1,2,4,5,10],
9       'min_samples_leaf': [3,5,10,12,15,16]
10  }
11
12  xgb_param_grid = {
13      'max_depth': [3,5,7,10,12],
14      'learning_rate': [0.1,0.01,0.001],
15      'n_estimators': [50,100,150,200,300],
16      'subsample': [0.8,1.0]
17  }
```

Here, we did the same thing that we do for Multi-Class Classification dataset. We initialize the Decision Tree and XGBoost classifier models. Then, we define the hyperparameter grids for each model to fine-tune the key parameters. For the Decision Tree classifier, we selected a specific range of values for `max_depth`, `min_samples_split`, and `min_samples_leaf`, allowing the model to explore different tree structures. We do the same for the XGBoost model, where we choose a specific range of values for `max_depth`, `learning_rate`, `n_estimators`, and `subsample`. [1]

**Training Decision Tree Model:**

```
1  # Decision Tree:
2  dt_grid_search = GridSearchCV(estimator=dt_classifier, param_grid=dt_param_grid, cv=5, scoring
3  dt_grid_search.fit(X_train, y_train)
```

```
4
5  dt_best_params = dt_grid_search.best_params_
6  dt_best_score = dt_grid_search.best_score_
7  print(f"Decision␣Tree␣Best␣Parameters:␣{dt_best_params}")
8  print(f"Decision␣Tree␣Best␣Score:␣{dt_best_score}")
```

Here, we perform hyperparameter tuning for the Decision Tree model using GridSearchCV by testing different combinations of parameters with 5-fold cross-validation to evaluate the models performance. Then we fit it to the training data and find the best parameters and best scores. [1]

**Training XGBoost Model:**

```
1  # XGBoost:
2  xgb_grid_search = GridSearchCV(estimator=xgb_classifier, param_grid=xgb_param_grid, cv=5, scor
3  xgb_grid_search.fit(X_train, y_train)
4
5  xgb_best_params = xgb_grid_search.best_params_
6  xgb_best_score = xgb_grid_search.best_score_
7  print(f"XGBoost␣Best␣Parameters:␣{xgb_best_params}")
8  print(f"XGBoost␣Best␣Score:␣{xgb_best_score}")
```

There we do the same thing for XGBoost model. We perform hyperparameter tuning with 5-fold cross-validation. Then fit it with training data and find the best parameters. [1]

**Retraining the Model with Validation and Training Set:**

```
1  # now retrain the model on the combined training and validation set with the best hyperparamet
2  X_train_val = np.concatenate((X_train, X_val))
3  y_train_val = np.concatenate((y_train, y_val))
4
5  dt_final_model = DecisionTreeClassifier(**dt_best_params, random_state=42)
6  xgb_final_model = XGBClassifier(**xgb_best_params, random_state=42)
7
8  # retrain decision tree final model
9  dt_final_model.fit(X_train_val, y_train_val)
10 print(f"Decision␣Tree␣Best␣Score:␣{dt_best_score}")
11
12 # retrain xgboost final model
13 xgb_final_model.fit(X_train_val, y_train_val)
14 print(f"XGBoost␣Best␣Score:␣{xgb_best_score}")
```

Here we retrain both models with the combined training and vaalidation data to make the best use of available data with the optimal best hyper parameters.

**Evaluate the Final Model on Test Set:**

```
1  # evaluate the final model on test set
2  y_test_pred_dt = dt_final_model.predict(X_test)
3  y_test_pred_xgb = xgb_final_model.predict(X_test)
```

This step, we evaluate the final model's performance on test set.

**Evaluation Metrics:**

```
1  print("Decision␣Tree:␣")
2  print(f'Accuracy:␣{accuracy_score(y_test,␣y_test_pred_dt)}')
3  print(f'Precision:␣{precision_score(y_test,␣y_test_pred_dt)}')
4  print(f'Recall:␣{recall_score(y_test,␣y_test_pred_dt)}')
```

```
5  print(f'F1-Score:␣{fscore(y_test,␣y_test_pred_dt)}')
6  print()
7
8  print("XGBoost:␣")
9  print(f'Accuracy:␣{accuracy_score(y_test,␣y_test_pred_xgb)}')
10 print(f'Precision:␣{precision_score(y_test,␣y_test_pred_xgb)}')
11 print(f'Recall:␣{recall_score(y_test,␣y_test_pred_xgb)}')
12 print(f'F1-Score:␣{fscore(y_test,␣y_test_pred_xgb)}')
```

We print the evaluation metrics Accuracy, Precision, Recall and F-score using the defined functions.

**Confusion Matrix:**

```
1  num_class = len(set(y_test))
2  fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(20, 8))
3
4  # confusion matrix for Decision Tree
5  conf_matrix_dt = np.zeros((num_class, num_class), dtype=int)
6  for actual, predict in zip(y_test, y_test_pred_dt):
7      conf_matrix_dt[actual, predict] += 1
8
9  # plotting
10 ax1.imshow(conf_matrix_dt, cmap='Blues', interpolation='nearest')
11 ax1.set_title('Decision␣Tree', fontsize=16)
12 ax1.set_xlabel('Predicted␣Label', fontsize=14)
13 ax1.set_ylabel('Actual␣Label', fontsize=14)
14
15 for i in range(num_class):
16     for j in range(num_class):
17         ax1.text(j, i, str(conf_matrix_dt[i, j]), ha='center', va='center', color='red', fonts
18
19 # confusion matrix for XGBoost
20 conf_matrix_xgb = np.zeros((num_class, num_class), dtype=int)
21 for actual, predict in zip(y_test, y_test_pred_xgb):
22     conf_matrix_xgb[actual, predict] += 1
23
24 # Plotting
25 ax2.imshow(conf_matrix_xgb, cmap='Blues', interpolation='nearest')
26 ax2.set_title('XGBoost', fontsize=16)
27 ax2.set_xlabel('Predicted␣Label', fontsize=14)
28 ax2.set_ylabel('Actual␣Label', fontsize=14)
29
30 for i in range(num_class):
31     for j in range(num_class):
32         ax2.text(j, i, str(conf_matrix_xgb[i, j]), ha='center', va='center', color='red', font
33
34 # common properties
35 for ax in (ax1, ax2):
36     ax.set_xticks(np.arange(num_class))
37     ax.set_yticks(np.arange(num_class))
38     ax.tick_params(axis='both', which='major', labelsize=12)
39
40 plt.suptitle('Confusion␣Matrices␣for␣Decision␣Tree␣and␣XGBoost', fontsize=18)
41 plt.show()
```

There we plot the Confusion Matrix for both model based on the predictions and actual values.

**Precision-Recall Curve :**

```
1  from sklearn.metrics import precision_recall_curve, average_precision_score
2  y_scores = dt_final_model.predict_proba(X_test)[:, 1]
3
4  precision, recall, thresholds = precision_recall_curve(y_test, y_scores)
5  average_precision = average_precision_score(y_test, y_scores)
6
7  plt.figure(figsize=(8, 6))
8  plt.plot(recall, precision, label=f'Precision-Recall curve (AP = {average_precision:.2f})')
9  plt.xlabel('Recall', fontsize=14)
10 plt.ylabel('Precision', fontsize=14)
11 plt.title('Precision-Recall Curve', fontsize=16)
12 plt.legend(loc='best')
13 plt.grid()
14 plt.show()
15
16 y_scores = xgb_final_model.predict_proba(X_test)[:, 1]
17
18 precision, recall, thresholds = precision_recall_curve(y_test, y_scores)
19 average_precision = average_precision_score(y_test, y_scores)
20
21 plt.figure(figsize=(8, 6))
22 plt.plot(recall, precision, label=f'Precision-Recall curve (AP = {average_precision:.2f})')
23 plt.xlabel('Recall', fontsize=14)
24 plt.ylabel('Precision', fontsize=14)
25 plt.title('Precision-Recall Curve', fontsize=16)
26 plt.legend(loc='best')
27 plt.grid()
28 plt.show()
```

Here, we plot the Precision-Recall Curve for both models using the built-in function from the scikit-learn module, as defining it manually is complex. This Precision-Recall Curve actually show the trade-off between precision and recall for a model at different thresholds.

## 2.3  Regression on 'DT-Wage' dataset

**Loading Dataset:**

```
1  df = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/Assignment-4/dataset/DT-Wage.csv')
2  print(df.shape)
3  print(df.head())
```

Load the Dataset using pandas library and check first few data from the dataset and its Shape.

**Checking Missing Values:**

```
1  print('Number of missing values: ')
2  print(df.isnull().sum())
```

Checking if any missing values are founds in any features.

**Encoded Categorical Data:**

```
1  df = pd.get_dummies(df, columns=['maritl', 'race', 'education', 'region', 'jobclass', 'health'
2  df.columns = df.columns.str.replace('[', '_').str.replace(']', '_').str.replace('<', '_')
3  print(df.shape)
4  print(df.head())
```

Since there has few complex relationship between the categories, so we used `pd.get_dummies` to convert the categorical columns into one-hot encoded columns Then, we replace special characters ([, ], ¡) in the column names with underscores. Then we print few samples from this encoded dataset.

**Split Data into Training, Validation and Test Sets:**

```
1  X = df.drop(columns=['wage'])
2  y = df['wage']
3
4  X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
5  X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42
6
7  print(f'Train_Data_Shape_(X,_y):_{X_train.shape,_y_train.shape}')
8  print(f'Validation_Data_Shape_(X,_y):_{X_val.shape,_y_val.shape}')
9  print(f'Test_Data_Shape_(X,_y):_{X_test.shape,_y_test.shape}')
```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing.

**Initialize the Classifier Models and Hyper-parameters for GridSearchCV:**

```
1  # initialize the regressor models
2  dt_regressor = DecisionTreeRegressor(random_state=42)
3  xgb_regressor = XGBRegressor(random_state=42)
4
5  # initialize hyperparameters for gridsearchCV
6  dt_param_grid = {
7      'max_depth': [2,3,4,5,8,10],
8      'min_samples_split': [2,4,5,10],
9      'min_samples_leaf': [1,3,5,8,10,12]
10 }
11
12 xgb_param_grid = {
13     'max_depth': [1,3,5,7,10],
14     'learning_rate': [0.1,0.01,0.001],
15     'n_estimators': [10000,15000,16000],
16     'subsample': [0.8,1.0]
17 }
```

Here, we did the same thing that we do for Multi-Class Classification dataset and Binary Classification dataset. We initialize the Decision Tree and XGBoost classifier models. Then, we define the hyperparameter grids for each model to fine-tune the key parameters. For the Decision Tree classifier, we selected a specific range of values for `max_depth`, `min_samples_split`, and `min_samples_leaf`, allowing the model to explore different tree structures. We do the same for the XGBoost model, where we choose a specific range of values for `max_depth`, `learning_rate`, `n_estimators`, and `subsample`. [1]

**Training Decision Tree Model:**

```
1  # Decision Tree:
2  dt_grid_search = GridSearchCV(estimator=dt_regressor, param_grid=dt_param_grid, cv=3, scoring=
3  dt_grid_search.fit(X_train, y_train)
4
5  dt_best_params = dt_grid_search.best_params_
6  print(f"Decision_Tree_Best_Parameters:_{dt_best_params}")
```

Here, we perform hyperparameter tuning for the Decision Tree model using GridSearchCV by testing different combinations of parameters with 3-fold cross-validation to evaluate the models performance. Then we fit it to the training data and find the best parameters and best scores. [1]

**Training XGBoost Model:**

```
1  # XGBoost:
2  xgb_grid_search = GridSearchCV(estimator=xgb_regressor, param_grid=xgb_param_grid, cv=3, scori
3  xgb_grid_search.fit(X_train, y_train)
4
5  xgb_best_params = xgb_grid_search.best_params_
6  print(f"XGBoost␣Best␣Parameters:␣{xgb_best_params}")
```

There we do the same thing for XGBoost model. We perform hyperparameter tuning with 3-fold cross-validation. Then fit it with training data and find the best parameters. [1]

**Retraining the Model with Validation and Training Set:**

```
1   # now retrain the model on the combined training and validation set with the best hyperparamet
2   X_train_val = np.concatenate((X_train, X_val))
3   y_train_val = np.concatenate((y_train, y_val))
4
5   dt_final_model = DecisionTreeClassifier(**dt_best_params, random_state=42)
6   xgb_final_model = XGBClassifier(**xgb_best_params, random_state=42)
7
8   # retrain decision tree final model
9   dt_final_model.fit(X_train_val, y_train_val)
10
11  # retrain xgboost final model
12  print(f"XGBoost␣Best␣Score:␣{xgb_best_score}")
```

Here we retrain both models with the combined training and vaalidation data to make the best use of available data with the optimal best hyper parameters.

**Evaluate the Final Model on Test Set:**

```
1  # evaluate the final model on test set
2  y_test_pred_dt = dt_final_model.predict(X_test)
3  y_test_pred_xgb = xgb_final_model.predict(X_test)
```

This step, we evaluate the final model's performance on test set.

**Calculating MSE:**

```
1   def mse(y_actual, y_pred):
2       error = 0.0
3
4       for actual, predicted in zip(y_actual, y_pred):
5           error += (actual - predicted) ** 2
6       mse = error / len(y_actual)
7       return mse
8
9   print(mse(y_test, y_test_pred_dt))
10  print(mse(y_test, y_test_pred_xgb))
```

Calculates the MSE with user defined function and print it for Decision Tree and XGBoost models.

## 2.4 Regression on 'DT-Credit' dataset

**Loading Dataset:**

```
1 df = pd.read_csv('/content/drive/MyDrive/CSE445-Assignment/Assignment-4/dataset/DT-Credit.csv'
2 print(df.shape)
3 print(df.head())
```

Load the Dataset using pandas library and check first few data from the dataset and its Shape.

**Checking Missing Values:**

```
1 print('Number_of_missing_values:_')
2 print(df.isnull().sum())
```

Checking if any missing values are founds in any features.

**Encoded Categorical Data:**

```
1 df = pd.get_dummies(df, columns=['Own', 'Student', 'Married', 'Region'])
2 print(df.shape)
3 print(df.head())
```

Since there has few complex relationship between the categories, so we used `pd.get_dummies` to convert the categorical columns into one-hot encoded columns. Then we print few samples from this encoded dataset.

**Split Data into Training, Validation and Test Sets:**

```
1 X = df.drop(columns=['Balance'])
2 y = df['Balance']
3
4 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.3, random_state=42)
5 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42
6
7 print(f'Train_Data_Shape_(X,_y):_{X_train.shape,_y_train.shape}')
8 print(f'Validation_Data_Shape_(X,_y):_{X_val.shape,_y_val.shape}')
9 print(f'Test_Data_Shape_(X,_y):_{X_test.shape,_y_test.shape}')
```

Splitting 70 percent data for training, 15 percent for validation and rest 15 percent for testing.

**Initialize the Classifier Models and Hyper-parameters for GridSearchCV:**

```
1  # initialize the Regressor models
2  dt_regressor = DecisionTreeRegressor(random_state=42)
3  xgb_regressor = XGBRegressor(random_state=42)
4
5  # initialize hyperparameters for gridsearchCV
6  dt_param_grid = {
7      'max_depth': [2,3,4,5,8,10],
8      'min_samples_split': [2,4,5,10],
9      'min_samples_leaf': [1,3,5,8,10,12]
10 }
11
12 xgb_param_grid = {
13     'max_depth': [1,3,5,7,10],
14     'learning_rate': [0.1,0.01,0.001],
15     'n_estimators': [10000,15000,16000],
16     'subsample': [0.8,1.0]
17 }
```

Here, we did the same thing that we do for Multi-Class Classification dataset and Binary Classification dataset. We initialize the Decision Tree and XGBoost classifier models. Then, we define the hyperparameter grids for each model to fine-tune the key parameters. For the Decision Tree classifier, we selected a specific range of values for `max_depth`, `min_samples_split`, and `min_samples_leaf`, allowing the model to explore different tree structures. We do the same for the XGBoost model, where we choose a specific range of values for `max_depth`, `learning_rate`, `n_estimators`, and `subsample`. [1]

**Training Decision Tree Model:**

```
1  # Decision Tree:
2  dt_grid_search = GridSearchCV(estimator=dt_regressor, param_grid=dt_param_grid, cv=3, scoring=
3  dt_grid_search.fit(X_train, y_train)
4
5  dt_best_params = dt_grid_search.best_params_
6  print(f"Decision_Tree_Best_Parameters:_{dt_best_params}")
```

Here, we perform hyperparameter tuning for the Decision Tree model using GridSearchCV by testing different combinations of parameters with 3-fold cross-validation to evaluate the models performance. Then we fit it to the training data and find the best parameters and best scores. [1]

**Training XGBoost Model:**

```
1  # XGBoost:
2  xgb_grid_search = GridSearchCV(estimator=xgb_regressor, param_grid=xgb_param_grid, cv=3, scori
3  xgb_grid_search.fit(X_train, y_train)
4
5  xgb_best_params = xgb_grid_search.best_params_
6  print(f"XGBoost_Best_Parameters:_{xgb_best_params}")
```

There we do the same thing for XGBoost model. We perform hyperparameter tuning with 3-fold cross-validation. Then fit it with training data and find the best parameters. [1]

**Retraining the Model with Validation and Training Set:**

```
1   # now retrain the model on the combined training and validation set with the best hyperparamet
2   X_train_val = np.concatenate((X_train, X_val))
3   y_train_val = np.concatenate((y_train, y_val))
4
5   dt_final_model = DecisionTreeClassifier(**dt_best_params, random_state=42)
6   xgb_final_model = XGBClassifier(**xgb_best_params, random_state=42)
7
8   # retrain decision tree final model
9   dt_final_model.fit(X_train_val, y_train_val)
10
11  # retrain xgboost final model
12  print(f"XGBoost_Best_Score:_{xgb_best_score}")
```

Here we retrain both models with the combined training and vaalidation data to make the best use of available data with the optimal best hyper parameters.

**Evaluate the Final Model on Test Set:**

```
1  # evaluate the final model on test set
2  y_test_pred_dt = dt_final_model.predict(X_test)
3  y_test_pred_xgb = xgb_final_model.predict(X_test)
```

This step, we evaluate the final model's performance on test set.

**Calculating MSE:**

```
1  def mse(y_actual, y_pred):
2      error = 0.0
3
4      for actual, predicted in zip(y_actual, y_pred):
5          error += (actual - predicted) ** 2
6      mse = error / len(y_actual)
7      return mse
8
9  print(mse(y_test, y_test_pred_dt))
10 print(mse(y_test, y_test_pred_xgb))
```

Calculates the MSE with user defined function and print it for Decision Tree and XGBoost models.

## 3. EXPERIMENT RESULTS

### 3.1 Decision Tree vs XGBoost in Multi-class Classification on 'Car Evaluation' dataset

**Accuracy:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 97.31% | 98.07% |

The XGBoost Classifier Model achieves higher accuracy (98.07%) compared to the Decision Tree (97.31%).

**Average Precision:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 0.8985 | 0.9224 |

The XGBoost Classifier Model has higher Average Precision (0.9224) compared to Decision Tree (0.8985).
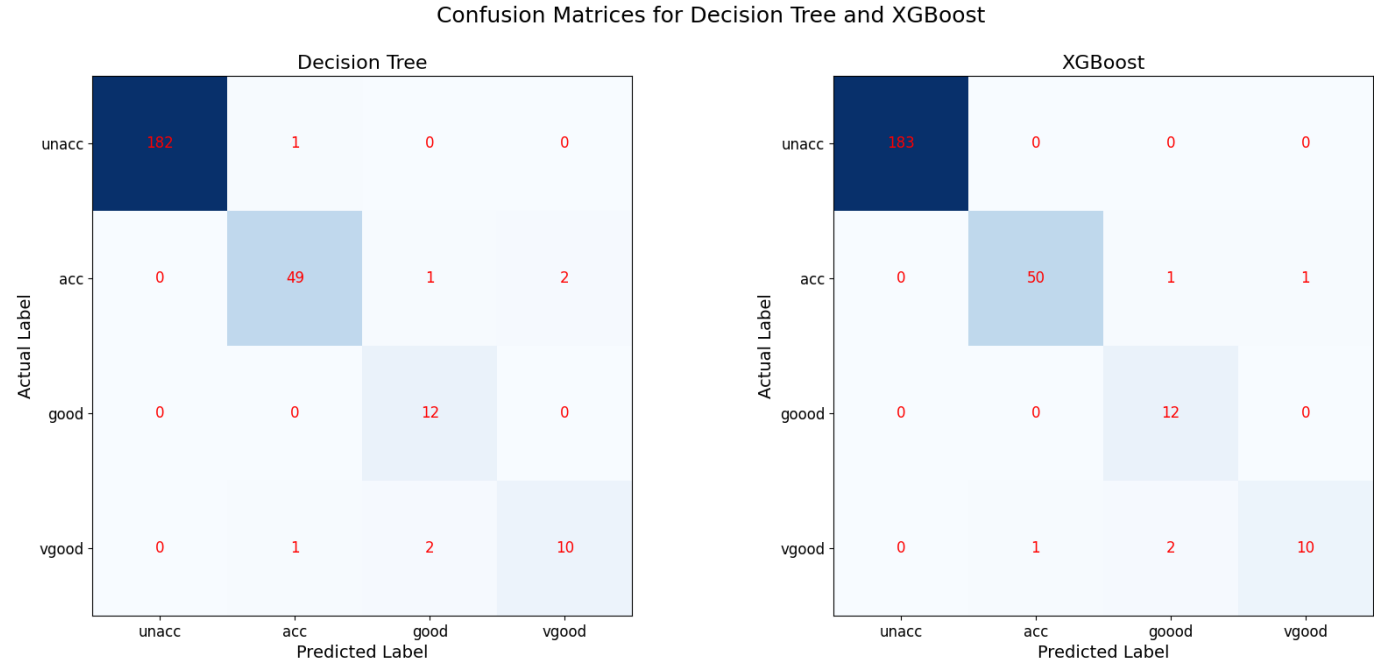
**Average Recall:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 0.9265 | 0.9327 |

The XGBoost Classifier Model has higher Average Precision (0.9224) compared to Decision Tree (0.8985).

**F-Score:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 0.9094 | 0.9233 |

The XGBoost Classifier Model has higher Average Precision (0.9224) compared to Decision Tree (0.8985).

**Confusion Matrix of Decision Tree vs XGBoost:**

Confusion Matrices for Decision Tree and XGBoost



## 3.2 Decision Tree vs XGBoost in Binary Classification on 'DT-BrainCancer' dataset

**Accuracy:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 78.57% | 92.85% |

The XGBoost Classifier Model achieves higher accuracy (92.85%) compared to the Decision Tree (78.57%).

**Precision:**

| Decision Tree | | XGBoost | |
|:---:|:---:|:---:|:---:|
| Class 0 | Class 1 | CLass 0 | Class 1 |
| 0.75 | 1.00 | 1.0 | 0.83 |

Overall XGBoost has higher recall compare to Decision Tree.

**Recall:**

| Decision Tree | | XGBoost | |
|:---:|:---:|:---:|:---:|
| Class 0 | Class 1 | CLass 0 | Class 1 |
| 1.00 | 0.40 | 0.88 | 1.00 |

Overall XGBoost has higher precision compare to Decision Tree.

**F-Score:**
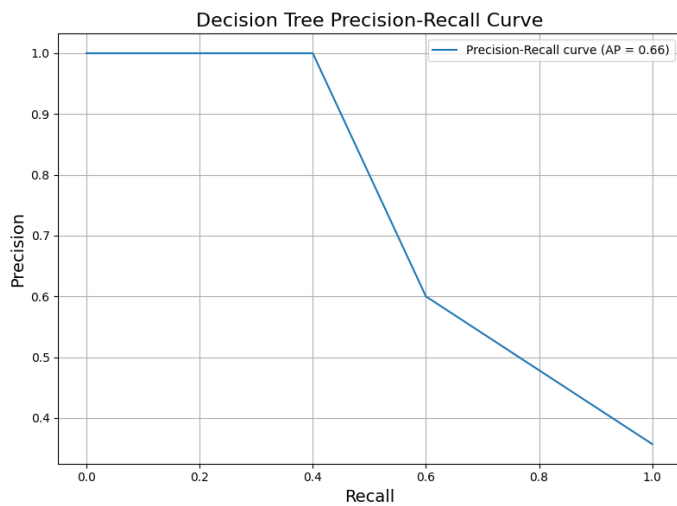
| Decision Tree | | XGBoost | |
|---|---|---|---|
| **Class 0** | **Class 1** | **CLass 0** | **Class 1** |
| 0.8571 | 0.5714 | 0.9411 | 0.9090 |

The XGBoost Classifier demonstrates higher F-score for both classes compared to decision tree.

**Confusion Matrix of Decision Tree vs XGBoost:**



Confusion Matrices for Decision Tree and XGBoost

**Precision-Recall Curve of Decision Tree vs XGBoost:**

### 3.3 Decision Tree vs XGBoost in Regression on 'DT-Wage' dataset

**Mean Squared Error:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 0.6654 | 2.6027 |

The XGBoost has a higher mean squared error 2.6027 compared to Decision Tree which have 0.6654. So the decision tree is a better fit for this data.

### 3.4 Decision Tree vs XGBoost in Regression on 'DT-Credit' dataset

**Mean Squared Error:**

| Decision Tree | XGBoost |
|:---:|:---:|
| 17559.2752 | 3757.1172 |

The XGBoost has a lower mean squared error 3757.11752 compared to Decision Tree which have 17559.2752. So the XGBoost is a better fit for this data though both do not fit the dataset well ude to extremely high MSE.

## 4. DISCUSSION

During this assignment, we learned more about Decision Tree, XGBoost, Cross Validation, Evaluation Metrics. We train the two different model on 4 different datasets. On classification task for both binary classification and multiclass classification XGBoost model perform better. F-score, Precision was almost larger compare to decision tree. For Regression task XGBoost work good on DT-Wage dataset but on DT-Credit dataset both models has high MSE. During this whole assignment we follow multiple resources, whole links are given below on the References section.

## REFERENCES

[1] Fermatsavant. Decision tree high accuracy using gridsearchcv. https://www.kaggle.com/code/fermatsavant/decision-tree-high-acc-using-gridsearchcv, 2020.

[2] T. D. Science. An exhaustive guide to classification using decision trees. https://towardsdatascience.com/an-exhaustive-guide-to-classification-using-decision-trees-8d472e77223f, 2024.

[3] S. SDR. Decision tree classification in python. https://medium.com/@shuv.sdr/decision-tree-classification-in-python-b1e59205949c, 2024.

[4] A. Vidhya. An end-to-end guide to understand the math behind xgboost. https://www.analyticsvidhya.com/blog/2018/09/an-end-to-end-guide-to-understand-the-math-behind-xgboost/, 2018.