# XGBOOST CODES

## CEN 481 – INTRODUCTION TO DATA MINING

## AHMET HAKAN BOZ - 2018556014

```python
1   # Model evaluations
2   from sklearn import metrics
3   from sklearn.metrics import accuracy_score, roc_auc_score, roc_curve, RocCurveDisplay
4   from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
5   from sklearn.metrics import classification_report, f1_score, precision_score, recall_score, average_precision_score
6
7   from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
8   from sklearn.model_selection import KFold, StratifiedKFold
9   from sklearn.model_selection import train_test_split, cross_val_score
10
11  from sklearn.pipeline import make_pipeline
12  from sklearn.pipeline import Pipeline
13  from sklearn import preprocessing
14  from sklearn.preprocessing import scale
15  from sklearn.preprocessing import StandardScaler
16
17  # Exploratory data analysis and plotting libraries
18  import numpy as np
19  import pandas as pd
20  import matplotlib.pyplot as plt
21  import plotly.express as px
22  import seaborn as sns
23  from scipy.io import arff
24  import warnings
25  warnings.simplefilter("ignore")
26
27  # Feature Selection
28  import mlxtend
29  from mlxtend.feature_selection import SequentialFeatureSelector as SFS
30
31  # Models from Scikit-Learn
32  from xgboost import XGBClassifier
33
34  # For fixing random_state parameters
35  seed = 20
36
37  df = pd.read_csv("/content/Acoustic Features.csv")
38  o_df = df.copy()
39
40  df.isna().sum()
41
42  df.describe().T
43
44  song_types = df["Class"].value_counts()
45  song_types_df = pd.DataFrame(song_types)
46  song_types_df = song_types.reset_index(level = 0)
47  song_types_df
48
49  # "relax": 0,
50  # "happy": 1,
51  # "sad": 2,
52  # "angry": 3
53  emotion_map = {"relax": 0, "happy": 1, "sad": 2, "angry": 3}
54  df["Class"] = df["Class"].map(emotion_map)
55  df
```

```python
56
57  # Split data into features and target
58  X = df.drop(["Class"], axis = 1)
59  y = df["Class"]
60
61  # Split into train and test set
62  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, stratify = y, random_state = seed)
63
64  scaler = preprocessing.StandardScaler()
65  scaler.fit(X_train)
66  X_trainStandart = scaler.transform(X_train)
67  X_testStandart = scaler.transform(X_test)
68
69  # Creating model
70  clf = XGBClassifier(random_state = seed)
71  # Searching parameters
72  params = {"n_estimators": [50, 100, 500, 1000],
73           "learning_rate": [1, 0.1, 0.01, 0.001]
74          }
75  # Creating grid
76  xg_clf_grid = RandomizedSearchCV(estimator = clf,
77                                   param_distributions = params,
78                                   cv = StratifiedKFold(n_splits = 10,
79                                                        shuffle = True,
80                                                        random_state = seed),
81                                   n_iter = 10,
82                                   verbose = 2,
83                                   scoring = "accuracy",
84                                   n_jobs = -1)
85  # Fit the model
86  xg_model = xg_clf_grid.fit(X_train_select, y_train)
```

XGBClassifier(base_score=None, booster=None, callbacks=None, colsample_bylevel=None, colsample_bynode=None, colsample_bytree=None, early_stopping_rounds=None, enable_categorical=False, eval_metric=None, feature_types=None, gamma=None, gpu_id=None, grow_policy=None, importance_type=None, interaction_constraints=None, learning_rate=0.01, max_bin=None, max_cat_threshold=None, max_cat_to_onehot=None, max_delta_step=None, max_depth=None, max_leaves=None, min_child_weight=None, missing=nan, monotone_constraints=None, n_estimators=1000, n_jobs=None, num_parallel_tree=None, objective='multi:softprob', predictor=None, ...)

```python
87
88  # Get best parameters
89  print("Best parameters for XGB model: ", xg_model.best_params_)
90  # Best parameters
91  xg_best = pd.DataFrame.from_dict(xg_model.best_params_, orient = "index").rename(columns = {0: "Best"})
92  xg_best
93  xg_clf = XGBClassifier(n_estimators = int(xg_best.iloc[0,0]),
94                         learning_rate = float(xg_best.iloc[1,0]),
95                         random_state = seed)
96  # Fit the model
97  xg_clf.fit(X_train, y_train)
98  #Predictions and model accuracy
99  xg_pred = xg_clf.predict(X_test)
100 xg_acc = accuracy_score(y_test, xg_pred)
101 print("XGB Model Accuracy:", xg_acc)
102 xg_acc_tr = xg_clf.score(X_train, y_train)
103 print("XGB Training Accuracy:", xg_acc_tr)
104
105
```

```
XGB Model Accuracy: 0.7875
XGB Training Accuracy: 1.0
```

```
106    # Classification Report
107    print("XGB Classification Report\n\n", classification_report(y_test, xg_pred))
```

```
XGB Classification Report

               precision    recall  f1-score   support

           0       0.80      0.80      0.80        20
           1       0.86      0.90      0.88        20
           2       0.61      0.55      0.58        20
           3       0.86      0.90      0.88        20

    accuracy                           0.79        80
   macro avg       0.78      0.79      0.78        80
weighted avg       0.78      0.79      0.78        80
```
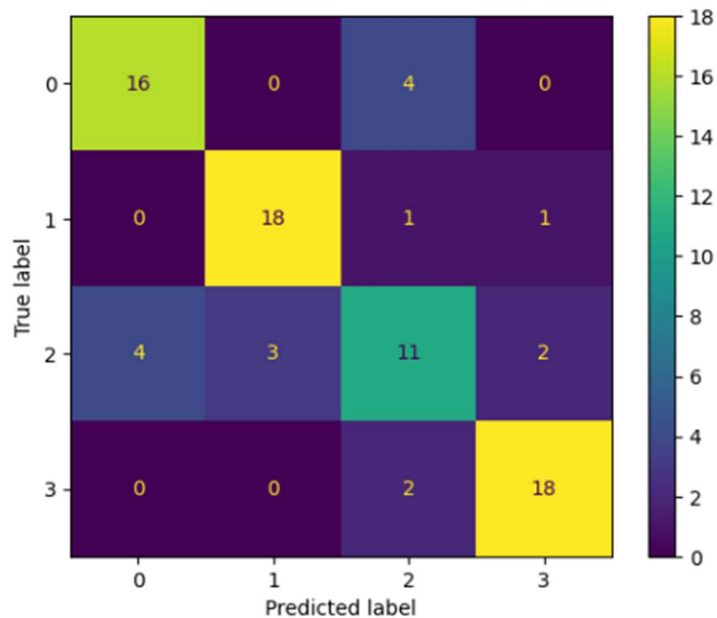
```
109    # For comparison list records
110    xg_recall = recall_score(y_test, xg_pred, average = None)
111    xg_prec = precision_score(y_test, xg_pred, average = None)
112    xg_f1 = f1_score(y_test, xg_pred, average = None)
113    # Confusion matrix
114    xg_cm = confusion_matrix(y_test, xg_pred, labels = xg_clf.classes_)
115    disp = ConfusionMatrixDisplay(confusion_matrix = xg_cm,
116                                  display_labels = xg_clf.classes_)
117    print("XGB Confusion Matrix")
118    disp.plot()
119    plt.show();
120
```

XGB Confusion Matrix



```
120
121    # Train/Test Performance Metrics
122    def calculatePerformance(classifier, X_train, y_train, X_test, y_test):
123        train_pred = classifier.predict(X_train)
124        test_pred = classifier.predict(X_test)
125        scores = {
126            "Train Accuracy": accuracy_score(y_train, train_pred),
127            "Test Accuracy": accuracy_score(y_test, test_pred),
128            "Train Recall": recall_score(y_train, train_pred, average = None),
129            "Test Recall": recall_score(y_test, test_pred, average = None),
130            "Train Precision": precision_score(y_train, train_pred, average = None),
131            "Test Precision": precision_score(y_test, test_pred, average = None),
132            "Train F1": f1_score(y_train, train_pred, average = None),
133            "Test F1": f1_score(y_test, test_pred, average = None)
134        }
135        print("Model Performance Metrics Comparison")
136        return scores
137
138    # Train/Test Performance Metrics
139    xg_pm = pd.DataFrame(calculatePerformance(xg_clf, X_train, y_train, X_test, y_test))*100
140    xg_pm
```

Model Performance Metrics Comparison

| | Train Accuracy | Test Accuracy | Train Recall | Test Recall | Train Precision | Test Precision | Train F1 | Test F1 |
|---|---|---|---|---|---|---|---|---|
| 0 | 100.0 | 78.75 | 100.0 | 80.0 | 100.0 | 80.000000 | 100.0 | 80.000000 |
| 1 | 100.0 | 78.75 | 100.0 | 90.0 | 100.0 | 85.714286 | 100.0 | 87.804878 |
| 2 | 100.0 | 78.75 | 100.0 | 55.0 | 100.0 | 61.111111 | 100.0 | 57.894737 |
| 3 | 100.0 | 78.75 | 100.0 | 90.0 | 100.0 | 85.714286 | 100.0 | 87.804878 |