# Spam Detection HW

**Read complete instructions before starting the HW**

## ∨ Q1: Load the dataset (1 Point)

- For this Hw you will usespam dataset from kaggle which can be found from <u>this</u> link. You can download this data and either upload it in google drive or in colab workspace. Load the data in pandas dataframe.

- There are only two useful columns. These columns are related to (1) label (ham and spam) and the (2) text of email.

- Rename columns as label and message

- Find the % ham and spam in the data.

```
1 #Load Dataset
2
3 import pandas as pd
4
5 df = pd.read_csv('spam.csv', encoding='ISO–8859–1')
6 df.head()
```

```
1 #Rename columns
2
3 df = df[['v1', 'v2']]
4 df.columns = ['label', 'message']
5 df.head()
```

```
1 # % of ham and spam
2
3 import matplotlib.pyplot as plt
4 df['label'].value_counts().plot.pie(autopct='%2.2f%%')
5 plt.show()
```

```
1 #Convert to Binary if spam then 1, else 0 to process ahead
2
3 df['label'] = df['label'].map(lambda x: 1 if x == 'spam' else 0) # converting
```

## Q2 : Provide the metric for evaluating model (1 Point)

As you will notice, the data is highly imbalanced (most messages are labelled as ham and only few are labelled as spam). Always predicting ham will give us very good accuracy (close to 90%). So you need to choose a different metric.

Task: Provde the metric you will choose to evaluate your model. Explain why this is an appropriate metric for this case.

appropriate metric for this case.

**I will be using ROC AUC score as the metric for this case.**

Why:

ROC AUC (Receiver Operating Characteristic Area Under the Curve) is ideal for imbalanced datasets because it evaluates a classifier's ability to rank instances across all possible thresholds, making it insensitive to class imbalance. It provides a comprehensive view of model performance, considering trade-offs between true positive and false positive rates. Additionally, it's interpretable, robust to class skew, and applicable to multi-class classification. Overall, it's a reliable metric for assessing classifier performance in imbalanced datasets.

## ⌄ Q3 : Classification Pipelines (18 Points)

In the previous lectures you learned Data processing, Featurization such as CountVectorizer, TFIDFVectorizer, and also Feature Engineering.

- You will now use folllowing methods to create fearures which you can use in your model.

    1. Sparse Embeddings (TF-IDF) (6 Points)
    2. Feature Engineering (see examples below) (6 Points)
    3. Sparse Embeddings (TF-IDF) + Feature Engineering (6 Points)

**Approach:**

**\*\*Use a smaller subset of dataset (e.g. 5-10 %) to evaluate the three pipelines . Based on your analysis (e.g. model score, learning curves) , choose one pipeline from the three. Provde your rational for choosing the pipleine. Train only the final pipeline on randomly selected larger subset (e.g. 40%) of the data.**

**Requirements:**

1. You can use any ML model (Logistic Regression, XgBoost) for the classification. You will need to tune the **model for imbalanced dataset** (The link on XGBoost tutorial for imbalanced data: https://machinelearningmastery.com/xgboost-for-imbalanced-classification/).

2. For feature engineering, you can choose from the examples below. You do not have to use all of them. You can add other featues as well. Think about what faetures can distinguish a spam from a regular email. Some examples :

    Count of following (Words, characters, digits, exclamation marks, numbers Nouns ProperNouns, AUX, VERBS, Adjectives, named

numbers, Nouns, ProperNouns, AUX, VERBS, Adjectives, named
entities, spelling mistakes (see the link on how to get spelling
mistakes https://pypi.org/project/pyspellchecker/).

3. For Sparse embeddings you will use **tfidf vectorization**. You need to choose appopriate
   parameters e.g. min_df, max_df, max_faetures, n-grams etc.).

4. Think carefully about the pre-processing you will do.

Tip: **Using GridSearch for hyperparameter tuning might take a lot of time. Try using
RandomizedSearch.** You can also explore faster implementation of Gridsearch and
RandomizedSearch in sklearn:

1. Halving Grid Search

2. HalvingRandomSearchCV

## ⌄ Part 1: Random 10% data for testing the 3 pipelines

```
 1 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
 2 from sklearn.model_selection import RepeatedStratifiedKFold
 3 from sklearn.feature_extraction.text import TfidfVectorizer
 4 from sklearn.base import TransformerMixin, BaseEstimator
 5 from sklearn.metrics import roc_auc_score, roc_curve
 6 from sklearn.model_selection import train_test_split
 7 from plot_learning_curve import plot_learning_curve
 8 from sklearn.model_selection import cross_val_score
 9 from sklearn.linear_model import LogisticRegression
10 from sklearn.metrics import classification_report
11 from sklearn.compose import ColumnTransformer
12 from FeaturizerSpacy import ManualFeatures
13 from sklearn.pipeline import Pipeline
14 import CustomPreprocessorSpacy as cp
15 from scipy.sparse import csr_matrix
16 from xgboost import XGBClassifier
17 from pathlib import Path
18 import numpy as np
19 import joblib
```

```
 1 #Small subset sampling
 2 df_ = df.sample(frac=0.1, replace=True, random_state=1)
```

```
 1 #Extracting values of column message into X and label into Y
 2 X, y = df_['message'].values, df_['label'].values
```

```
 1 #Make base directories
```

```
2 base_folder = Path('./')
3 data_folder = base_folder/'datasets/'
4 model_folder = base_folder/'models/'
5
6 model_folder.mkdir(exist_ok=True, parents=True)
7 data_folder.mkdir(exist_ok=True, parents=True)
```

```
1 # Splitting the data into train test splits with test as 20%
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, rando
```

## Approach 1: Sparse Embeddings (TFIDF)

```
1 # Cleaning the text data in X_train and X_test
2 X_train_cleaned_bow = cp.SpacyPreprocessor(model='en_core_web_sm', batch_size=
3                                            remove_url=False, remove_punct=Fals
4 X_test_cleaned_bow = cp.SpacyPreprocessor(model='en_core_web_sm', batch_size=5
5                                            remove_url=False, remove_punct=False
6
7 # save  this to a file
8 file_X_train_cleaned_bow = data_folder / 'x_train_cleaned_bow_small.pkl'
9 joblib.dump(X_train_cleaned_bow, file_X_train_cleaned_bow)
10
11 # save  this to a file
12 file_X_test_cleaned_bow = data_folder / 'x_test_cleaned_bow_small.pkl'
13 joblib.dump(X_test_cleaned_bow, file_X_test_cleaned_bow)
```

```
/content/CustomPreprocessorSpacy.py:83: MarkupResemblesLocatorWarning: The in
  soup = BeautifulSoup(text, "html.parser")
['datasets/x_test_cleaned_bow_small.pkl']
```

```
1 #Using TFIDF and XFBClassifier
2 classifier_1 = Pipeline([
3     ('vectorizer', TfidfVectorizer(analyzer='word', token_pattern=r"[\S]+")),
4     ('classifier', XGBClassifier())])
5
6 # classifier_1 = Pipeline([
7 #     ('vectorizer', TfidfVectorizer(analyzer='word', token_pattern=r"[\S]+"))
8 #     ('classifier', LogisticRegression(max_iter=10000)),])
```

```
1 #Param Grid
2 params_classifier_1 = {
3                         'vectorizer__max_features': [500, 1000, 1500, 2000, 2500
4                         'vectorizer__min_df': [1, 5, 10],  # Minimum number of d
5                         'vectorizer__max_df': [0.5, 0.75, 1.0],  # Maximum propo
6                         'vectorizer__ngram_range': [(1, 1), (1, 2), (2, 2)],  # 
7                         'classifier__n_estimators': [100, 200, 300],  # Number o
8                         'classifier__max_depth': [3, 4, 5]  # Maximum depth of t
9                         }
```

```
-                             ,
10
11 # params_classifier_1 = {'vectorizer__max_features': [1000, 2000, 5000],
12 #                              'classifier__C': [1, 10, 100,]}
```

```
1 #Defining RandomSearch CV Classifier
2 randomsearchCV_classifier_1 = RandomizedSearchCV(estimator=classifier_1, param
```

```
1 #Fit RandomSearch CV on train data
2 %%time
3 randomsearchCV_classifier_1.fit(X_train_cleaned_bow, y_train)
```

```
1 #Validation best scores,parameters and estimators
2 print(f'Best cross-validation score: {randomsearchCV_classifier_1.best_score_:
3 print("\nBest parameters: ", randomsearchCV_classifier_1.best_params_)
4 print("\nBest estimator: ", randomsearchCV_classifier_1.best_estimator_)
```

```
    Best cross-validation score: 0.91

    Best parameters:  {'vectorizer__ngram_range': (1, 2), 'vectorizer__min_df': 5

    Best estimator:  Pipeline(steps=[('vectorizer',
                      TfidfVectorizer(max_df=0.5, max_features=500, min_df=5,
                                      ngram_range=(1, 2), token_pattern='[\\S]+'))
                     ('classifier',
                      XGBClassifier(base_score=None, booster=None, callbacks=None,
                                    colsample_bylevel=None, colsample_bynode=None,
                                    colsample_bytree=None, device=None,
                                    early_stopping_rounds=None,
                                    enable_categorical=False, eval_metric=None,
                                    feat...grow_policy=None,
                                    importance_type=None,
                                    interaction_constraints=None, learning_rate=Nor
                                    max_bin=None, max_cat_threshold=None,
                                    max_cat_to_onehot=None, max_delta_step=None,
                                    max_depth=3, max_leaves=None,
                                    min_child_weight=None, missing=nan,
                                    monotone_constraints=None, multi_strategy=None
                                    n_estimators=100, n_jobs=None,
                                    num_parallel_tree=None, random_state=None, ...
```

```
1 file_best_estimator_pipeline1_round1 = model_folder / \
2      'pipeline1_round1_best_estimator.pkl'
```

```
 2      pipeline1_round1_best_estimator.pkl
 3 file_complete_grid_pipeline1_round1 = model_folder / \
 4     'pipeline1_round1_complete_grid.pkl'
 5
 6 joblib.dump(randomsearchCV_classifier_1.best_estimator_,
 7             file_best_estimator_pipeline1_round1)
 8 joblib.dump(randomsearchCV_classifier_1, file_complete_grid_pipeline1_round1)
 9
10 # load the saved model
11 best_estimator_pipeline1_round1 = joblib.load(file_best_estimator_pipeline1_ro
12 complete_grid_pipeline1_round1 = joblib.load(file_complete_grid_pipeline1_roun
```

```
 1 #Plotting learning curves
 2 %%time
 3 plot_learning_curve(randomsearchCV_classifier_1, 'Learning Curves',
 4                     X_train_cleaned_bow, y_train, n_jobs=-1)
```

```
 1 # let's check the train scores
 2 print(best_estimator_pipeline1_round1.score(X_train_cleaned_bow, y_train))
 3
 4 # let's check the cross validation score
 5 print(complete_grid_pipeline1_round1.best_score_)
```

```
    0.9955056179775281
    0.9148351648351648
```

```
 1 # Final Pipeline
```

```
1 # Final Pipeline
2 def final_pipeline(text):
3     # cleaned_text = cp.SpacyPreprocessor(model='en_core_web_sm').transform(te
4     cleaned_text = joblib.load(file_X_test_cleaned_bow)
5     best_estimator_pipeline1_round1 = joblib.load(
6         file_best_estimator_pipeline1_round1)
7     predictions = best_estimator_pipeline1_round1.predict(cleaned_text)
8     return predictions
```

```
1 # predicted values for Test data set
2 y_test_pred = final_pipeline(X_test)
```

```
1 #Classification report on test data
2 print('\nTest set classification report:\n\n',
3       classification_report(y_test, y_test_pred))
```

```
Test set classification report:

               precision    recall  f1-score   support

           0       0.95      0.98      0.96        97
           1       0.83      0.67      0.74        15

    accuracy                           0.94       112
   macro avg       0.89      0.82      0.85       112
weighted avg       0.93      0.94      0.93       112
```

```
 1 #Printing ROC score results
 2
 3 roc_auc = roc_auc_score(y_test, y_test_pred)
 4 print(f"ROC AUC Score: {roc_auc}")
 5
 6 fpr, tpr, thresholds = roc_curve(y_test, y_test_pred)
 7
 8 # Plot ROC curve
 9 plt.figure(figsize=(4,3))
10 plt.plot(fpr, tpr, label=f'ROC Curve (area = {roc_auc:.2f})')
11 plt.plot([0, 1], [0, 1], 'k--')  # Dashed diagonal
12 plt.xlabel('False Positive Rate')
13 plt.ylabel('True Positive Rate')
14 plt.title('ROC Curve')
15 plt.legend(loc="lower right")
16 plt.show()
```

## ⌄ Approach 2: Feature Engineering

```
 1 # Cleaning the text data in X_train and X_test
 2 X_train_cleaned_basic = cp.SpacyPreprocessor(model='en_core_web_sm',
 3                                             lemmatize=False, lower=False,
 4                                             remove_stop=False, remove_punct=F
 5                                             remove_email=False, remove_url=Fa
 6                                             add_user_mention_prefix=False,
 7                                             basic_clean_only=True).transform(
 8
 9 X_test_cleaned_basic = cp.SpacyPreprocessor(model='en_core_web_sm',
10                                             lemmatize=False, lower=False,
11                                             remove_stop=False, remove_punct=Fa
12                                             remove_email=False, remove_url=Fal
13                                             add_user_mention_prefix=False,
14                                             basic_clean_only=True).transform(X
15
16 # save  this to a file
17 file_X_train_cleaned_basic = data_folder / 'x_train_cleaned_basic_small.pkl'
18 joblib.dump(X_train_cleaned_basic, file_X_train_cleaned_basic)
19
20 # save  this to a file
21 file_X_test_cleaned_basic = data_folder / 'x_test_cleaned_basic_small.pkl'
22 joblib.dump(X_test_cleaned_basic, file_X_test_cleaned_basic)
```

```
    /content/CustomPreprocessorSpacy.py:83: MarkupResemblesLocatorWarning: The inp
      soup = BeautifulSoup(text, "html.parser")
    ['datasets/x_test_cleaned_basic_small.pkl']
```

```
 1 #Using Manual Features
 2 featurizer = ManualFeatures(spacy_model='en_core_web_sm', batch_size = 100)
 3 X_train_cleaned_basic = joblib.load(file_X_train_cleaned_basic)
 4 X_train_features, feature_names = featurizer.fit_transform(X_train_cleaned_bas
```

```
 1 pd.DataFrame(X_train_features, columns=feature_names).head()
```

```
1 #Define classifier and param grid
2 classifier_2 = Pipeline([('classifier', XGBClassifier())])
3 params_classifier_2 = {'classifier__n_estimators': [100, 200, 300],
4                        'classifier__max_depth': [3, 4, 5]}
```

```
1 #Define Random Search CV Classifier
2 randomsearchCV_classifier_2 = RandomizedSearchCV(estimator=classifier_2, param
```

```
1 #Fit on train data
2 randomsearchCV_classifier_2.fit(X_train_features, y_train)
```

```
1 ##Validation best scores,parameters and estimators
2 print(f'Best cross-validation score: {randomsearchCV_classifier_2.best_score_:
3 print("\nBest parameters: ", randomsearchCV_classifier_2.best_params_)
4 print("\nBest estimator: ", randomsearchCV_classifier_2.best_estimator_)
```

```
    Best cross-validation score: 0.97

    Best parameters:  {'classifier__n_estimators': 300, 'classifier__max_depth': :

    Best estimator:  Pipeline(steps=[('classifier',
                     XGBClassifier(base_score=None, booster=None, callbacks=None,
                                   colsample_bylevel=None, colsample_bynode=None,
                                   colsample_bytree=None, device=None,
                                   early_stopping_rounds=None,
                                   enable_categorical=False, eval_metric=None,
                                   feature_types=None, gamma=None, grow_policy=Noi
                                   importance_type=None,
                                   interaction_constraints=None, learning_rate=Noi
                                   max_bin=None, max_cat_threshold=None,
                                   max_cat_to_onehot=None, max_delta_step=None,
```

```
                                  max_cat_to_onehot=None, max_delta_step=None,
                                  max_depth=3, max_leaves=None,
                                  min_child_weight=None, missing=nan,
                                  monotone_constraints=None, multi_strategy=None,
                                  n_estimators=300, n_jobs=None,
                                  num_parallel_tree=None, random_state=None, ...
```

```
 1 file_best_estimator_pipeline2_round1 = model_folder / \
 2     'pipeline2_round1_best_estimator.pkl'
 3 file_complete_grid_pipeline2_round1 = model_folder / \
 4     'pipeline2_round1_complete_grid.pkl'
 5
 6 joblib.dump(randomsearchCV_classifier_2.best_estimator_,
 7             file_best_estimator_pipeline2_round1)
 8 joblib.dump(randomsearchCV_classifier_2, file_complete_grid_pipeline2_round1)
 9
10 # load the saved model
11 best_estimator_pipeline2_round1 = joblib.load(file_best_estimator_pipeline2_ro
12 complete_grid_pipeline2_round1 = joblib.load(file_complete_grid_pipeline2_roun
```

```
 1 #Learning curves for model 2
 2 %%time
 3 plot_learning_curve(best_estimator_pipeline2_round1, 'Learning Curves',
 4                     X_train_features, y_train, n_jobs=-1)
```

```
 1 # let's check the train scores
```

```
2 print(best_estimator_pipeline2_round1.score(X_train_features, y_train))
3
4 # let's check the cross validation score
5 print(complete_grid_pipeline2_round1.best_score_)
```

```
    1.0
    0.9716949716949717
```

```
1 best_estimator_pipeline2_round1
```

```
1 # Final Pipeline
2 def final_pipeline(text):
3     text_cleaned = joblib.load(file_X_test_cleaned_basic)
4     features, feature_names = featurizer.fit_transform(text_cleaned)
5     best_estimator_pipeline2_round1 = joblib.load(
6         file_best_estimator_pipeline2_round1)
7     predictions = best_estimator_pipeline2_round1.predict(features)
8     return predictions
```

```
1 y_test_pred = final_pipeline(X_test)
```

```
1 #Classification report on test data
2 print('\nTest set classification report:\n\n', classification_report(y_test, y
```

```
    Test set classification report:

                  precision    recall  f1-score   support

               0       0.98      0.96      0.97        97
               1       0.76      0.87      0.81        15

        accuracy                           0.95       112
       macro avg       0.87      0.91      0.89       112
    weighted avg       0.95      0.95      0.95       112
```

```
1 ##Printing ROC score results
2
3 roc_auc = roc_auc_score(y_test, y_test_pred)
4 print(f"ROC AUC Score: {roc_auc}")
5
6 fpr, tpr, thresholds = roc_curve(y_test, y_test_pred)
7
8 # Plot ROC curve
9 plt.figure(figsize=(4,3))
```

```
 9  plt.figure(figsize=(4,3))
10  plt.plot(fpr, tpr, label=f'ROC Curve (area = {roc_auc:.2f})')
11  plt.plot([0, 1], [0, 1], 'k--')  # Dashed diagonal
12  plt.xlabel('False Positive Rate')
13  plt.ylabel('True Positive Rate')
14  plt.title('ROC Curve')
15  plt.legend(loc="lower right")
16  plt.show()
```

## ∨ Approach 3: Sparse Embeddings (TF-IDF) + Feature Engineering

```
1  X_train_cleaned_bow = joblib.load(file_X_train_cleaned_bow)
2  X_train_final = pd.concat((pd.DataFrame(X_train_cleaned_bow, columns=['cleaned
3                             pd.DataFrame(X_train_features, columns=feature_name
4
5  X_train_final.head(2)
```

```
1 #Transformer class
2 class SparseTransformer(TransformerMixin, BaseEstimator):
3     def __init__(self):
4         pass
5     def fit(self, X, y=None):
6         return self
7
8     def transform(self, X, y=None):
9         return csr_matrix(X)
```

```
1 #Putting Sparse features and vectorizer in pipeline
2 sparse_features = Pipeline([('sparse', SparseTransformer()), ])
3 vectorizer = Pipeline([('tfidf', TfidfVectorizer(max_features=5)), ])
```

```
1 combined_features = ColumnTransformer(
2     transformers=[
3         ('tfidf', vectorizer, 'cleaned_text'),
4     ], remainder=sparse_features
5 )
```

```
1 #Defining classifier
2 classifier_3 = Pipeline([('combined_features',  combined_features),
3                          ('classifier', XGBClassifier())])
```

```
1 #Defining param grid
2 params_classifier_3 = {'combined_features__tfidf__tfidf__max_features': [500,
3                        'combined_features__tfidf__tfidf__min_df': [1, 2, 3],
4                        'combined_features__tfidf__tfidf__max_df': [0.5, 0.75,
5                        'combined_features__tfidf__tfidf__ngram_range': [(1, 1)
6                        'classifier__n_estimators': [100, 200, 300],
7                        'classifier__max_depth': [3, 4, 5]
8                        }
```

```
1 ##Define Random Search CV Classifier
2 randomsearchCV_classifier_3 = RandomizedSearchCV(estimator=classifier_3, param
```

```
1 ##Fit on train data
2 randomsearchCV_classifier_3.fit(X_train_final, y_train)
```

```
1 #Printing cross valdiation report
2 print("Best cross-validation score: {:.2f}".format(randomsearchCV_classifier_3
3 print("\nBest parameters: ", randomsearchCV_classifier_3.best_params_)
4 print("\nBest estimator: ", randomsearchCV_classifier_3.best_estimator_)
```

```
    Best cross-validation score: 0.97

    Best parameters:  {'combined_features__tfidf__tfidf__ngram_range': (1, 1), 'co

    Best estimator:  Pipeline(steps=[('combined_features',
                      ColumnTransformer(remainder=Pipeline(steps=[('sparse',
                                                                   SparseTransform
                               transformers=[('tfidf',
                                              Pipeline(steps=[('tfidf',
                                                               TfidfVecto
                                              'cleaned_text')])),
                    ('classifier',
                     XGBClassifier(base_score=None, booster=None, callbacks=None,
                                   colsample_bylevel=None, colsample_bynode=None,
                                   colsample_bytr...
                                   feature_types=None, gamma=None, grow_policy=Nor
                                   importance_type=None,
                                   interaction_constraints=None, learning_rate=Nor
                                   max_bin=None, max_cat_threshold=None,
                                   max_cat_to_onehot=None, max_delta_step=None,
                                   max_depth=5, max_leaves=None,
                                   min_child_weight=None, missing=nan,
                                   monotone_constraints=None, multi_strategy=None
                                   n_estimators=200, n_jobs=None,
                                   num_parallel_tree=None, random_state=None, ...
```

```
 1 file_best_estimator_pipeline3_round1 = model_folder / \
 2     'pipeline3_round1_best_estimator.pkl'
 3 file_complete_grid_pipeline3_round1 = model_folder / \
 4     'pipeline3_round1_complete_grid.pkl'
 5
 6 joblib.dump(randomsearchCV_classifier_3.best_estimator_, file_best_estimator_p
 7 joblib.dump(randomsearchCV_classifier_3, file_complete_grid_pipeline3_round1)
 8
 9 # load the saved model
10 best_estimator_pipeline3_round1 = joblib.load(file_best_estimator_pipeline3_ro
11 complete_grid_pipeline3_round1 = joblib.load(file_complete_grid_pipeline3_roun
```

```
1 #learning curve for model 3
2 %%time
3 plot_learning_curve(best_estimator_pipeline3_round1, 'Learning Curves',
4                     X_train_final, y_train, n_jobs=-1)
```

```
1 # let's check the train scores
2 print(best_estimator_pipeline3_round1.score(X_train_final, y_train))
3
4 # let's check the cross validation score
5 print(complete_grid_pipeline3_round1.best_score_)
```

```
    1.0
    0.9727938727938727
```

```
 1 # Final Pipeline
 2 def final_pipeline(text):
 3     # cleaned_text = cp.SpacyPreprocessor(model='en_core_web_sm', batch_size =
 4     cleaned_text = joblib.load(file_X_test_cleaned_bow)
 5     X_features, feature_names = featurizer.fit_transform(text)
 6     X_final = pd.concat((pd.DataFrame(cleaned_text, columns=['cleaned_text']),
 7                          pd.DataFrame(X_features, columns=feature_names)), axi
 8     best_estimator_pipeline3_round1 = joblib.load(
 9         file_best_estimator_pipeline3_round1)
10     predictions = best_estimator_pipeline3_round1.predict(X_final)
11     return predictions
```

```
1 # predicted values for Test data set
2 y_test_pred = final_pipeline(X_test)
```

```
1 #Classification report on test data
2 print('\nTest set classification report:\n\n',
3       classification_report(y_test, y_test_pred))
```

```
    Test set classification report:
```

Test Set Classification Report:

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.98 | 0.95 | 0.96 | 97 |
| 1 | 0.72 | 0.87 | 0.79 | 15 |
| accuracy |  |  | 0.94 | 112 |
| macro avg | 0.85 | 0.91 | 0.88 | 112 |
| weighted avg | 0.94 | 0.94 | 0.94 | 112 |

```
1 #ROC Score on model 3
2 roc_auc = roc_auc_score(y_test, y_test_pred)
3 print(f"ROC AUC Score: {roc_auc}")
4
5 fpr, tpr, thresholds = roc_curve(y_test, y_test_pred)
6
7 # Plot ROC curve
8 plt.figure(figsize=(4,3))
9 plt.plot(fpr, tpr, label=f'ROC Curve (area = {roc_auc:.2f})')
10 plt.plot([0, 1], [0, 1], 'k--')  # Dashed diagonal
11 plt.xlabel('False Positive Rate')
12 plt.ylabel('True Positive Rate')
13 plt.title('ROC Curve')
14 plt.legend(loc="lower right")
15 plt.show()
```

# Running the best pipeline on 40% sampled data

## ⌄ **Required Submissions:**

1. Submit two colab/jupyter notebooks

- (analysis with smaller subset and all three pipelines)
- (analysis with bigger subset and only final pipeline)

2. Pdf version of the notebooks (HWs will not be graded if pdf version is not provided.

3. **The notebooks and pdf files should have the output.**

4. **Name files as follows : FirstName_file1_hw2, FirstName_file2_h2**

```
1 Start coding or generate with AI.
```