# Spam Email Detection Project Documentation

## Monalisa Burma

## Date: 19/02/2024

## Data source

The dataset used in this project is the Spambase dataset, which contains email messages labeled as spam or non-spam (ham). The dataset provides a collection of features extracted from these emails to aid in classification.

Source: https://archive.ics.uci.edu/dataset/94/spambase

## Preprocessing steps

1. **Loading the Data:**

The data was loaded from the spambase.data file using Python's file handling.

2. **Feature Engineering:**

Extracted feature names and data from the loaded file.

Transformed the data into a Pandas DataFrame for further analysis.

3. **Text Vectorization:**

Utilized TF-IDF (Term Frequency-Inverse Document Frequency) vectorization for text-based features.

Processed the email text, applying vectorization to convert it into a numerical format suitable for machine learning.

4. **Train-Test Split:**

Split the dataset into training and testing sets for model evaluation.

# Model Selection

## Support Vector Machine (SVM)

1. **Model Training:**

Implemented a Support Vector Machine (SVM) classifier for spam detection.

Utilized TF-IDF vectorized features for training the SVM model.

2. **Evaluation Metrics:**

Assessed model performance using accuracy, precision, recall, and F1-score.

Conducted a train-test split and evaluated the SVM classifier on the test set.

3. **Hyperparameter Tuning:**

Explored different hyperparameters using GridSearchCV to optimize the SVM model's performance.

# Evaluation Results

## SVM Classifier

- Accuracy: 73.72%

- Classification Report:

```
Accuracy: 0.7372421281216069
Classification Report:
              precision    recall  f1-score   support

         0.0       0.73      0.86      0.79       531
         1.0       0.75      0.57      0.65       390

    accuracy                           0.74       921
   macro avg       0.74      0.72      0.72       921
weighted avg       0.74      0.74      0.73       921
```

- Confusion Matrix:

```
Confusion Matrix:
[[455  76]
 [166 224]]
```

# Model Deployment

## Command Line Interface (CLI):

- Created a simple user interface for users to input an email text.

- Integrated the trained SVM classifier and TF-IDF vectorizer into a joblib-loaded model.

- Implemented a function to predict whether the input email is spam or ham.

# Conclusion

In summary, this project centered around the implementation of a Support Vector Machine (SVM) model for spam email detection, achieving a notable accuracy of 73.72% on the test dataset. Leveraging the TF-IDF vectorization technique, the model exhibited robust performance in distinguishing between spam and ham emails. The classification report provided detailed insights into precision, recall, and F1-score, offering a comprehensive evaluation of the model's effectiveness.

Beyond model training and evaluation, the project extended its utility by incorporating a user-friendly Command Line Interface (CLI). This CLI enables users to input email text, and the SVM classifier, along with the TF-IDF vectorizer, promptly predicts whether the email is spam or ham. This practical deployment aspect enhances the accessibility and applicability of the model, showcasing its potential for real-world scenarios in email filtering and security.

Github Link:

https://github.com/monalisaburma/Coding_Samurai/tree/main/Spam%20Email%20Classifier