## Spam Email Detector Machine Learning Model

This is my project that I made during my free time to demonstrate knowledge in machine learning and cybersecurity. This model has been created to identify spam/phishing emails.

```
1  import pandas as pd
2  import numpy as np
3  from sklearn.model_selection import train_test_split
4  from sklearn.feature_extraction.text import TfidfVectorizer
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
7  import nltk
8  from nltk.corpus import stopwords
9  import re
10 from sklearn.model_selection import train_test_split
```

```
1  nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
True
```

The file that contains around 5000 emails classified into spam or not spam.

```
1  data = pd.read_csv('https://github.com/yushika-j/Playing-with-datasets/raw/
   master/Email%20Spam%20Filtering/emails.csv', encoding='utf-8')
2  data.head()
```

|   | text | spam |
|---|------|------|
| 0 | Subject: naturally irresistible your corporate... | 1 |
| 1 | Subject: the stock trading gunslinger fanny i... | 1 |
| 2 | Subject: unbelievable new homes made easy im ... | 1 |
| 3 | Subject: 4 color printing special request add... | 1 |
| 4 | Subject: do not have money , get software cds ... | 1 |

Next steps:  [ Generate code with `data` ]  [ ◉ View recommended plots ]  [ New interactive sheet ]

## Data Preprocessing

Once loaded, I proceeded to preprocess and vectorize the text data to convert it into a format that a machine learning model can understand. I used TF-IDF or CountVectorizer from scikit-learn to convert the email content into a term-document matrix.

```
1  print(data.isnull().sum())
2
3  # Drop any rows with missing values (if any)
4  data.dropna(inplace=True)
5
6  # Check for duplicates
7  data.drop_duplicates(inplace=True)
```

```
text    0
spam    0
dtype: int64
```

```
1  data['text'] = data['text'].str.lower()
2  #convert to lowercase
```

```
1  data['text'] = data['text'].apply(lambda x: re.sub(r'[^a-z\s]', '', x))
2  # remove special characters
```

```
1  #remove stopwords like 'the' and 'and'
2  stop_words = set(stopwords.words('english'))
3  data['text'] = data['text'].apply(lambda x: ' '.join([word for word in x.split() if word not in (stop_words)]))
```

## Text Vectorization

To convert the cleaned text into a format that can be used by a machine learning model, use TF-IDF or Count Vectorizer to create a term-document matrix.

This will transform each email into a vector of numerical features where each feature will represent a frequency.

```
1  # Initialize the TF-IDF Vectorizer
2  vectorizer = TfidfVectorizer(max_features=1000)  # Limit to 1000 most frequent words
3  X = vectorizer.fit_transform(data['text']).toarray()  # Transforms text data to numerical array
4  # Target variable
5  y = data['spam']
6  #print the data
7  print(X)
```

```
[[0.         0.         0.         ... 0.         0.         0.        ]
 [0.         0.         0.         ... 0.         0.         0.        ]
 [0.         0.         0.         ... 0.         0.         0.        ]
 ...
 [0.         0.         0.         ... 0.         0.         0.        ]
 [0.         0.         0.         ... 0.         0.         0.        ]
 [0.15747477 0.         0.117228   ... 0.         0.         0.        ]]
```

## Train-Test split

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
2
```

## Train a Classification Model

I started with a model like Logistic Regression for simplicity purposes. If the requirements increase, I might change the model to a neural network model.

```
1 from sklearn.linear_model import LogisticRegression
2
3 model = LogisticRegression()
4 model.fit(X_train, y_train)
```

```
▾  LogisticRegression  ⓘ ⑦
LogisticRegression()
```

## Model Evaluation

```
 1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
 2
 3 y_pred = model.predict(X_test)
 4
 5 accuracy = accuracy_score(y_test, y_pred)
 6 precision = precision_score(y_test, y_pred)
 7 recall = recall_score(y_test, y_pred)
 8 f1 = f1_score(y_test, y_pred)
 9
10 print(f'Accuracy: {accuracy}')
11 print(f'Precision: {precision}')
12 print(f'Recall: {recall}')
13 print(f'F1 Score: {f1}')
14
```

```
Accuracy: 0.9754170324846356
Precision: 0.9855072463768116
Recall: 0.918918918918919
F1 Score: 0.951048951048951
```

## Hyperparameter tuning

For optimal performance, I used parameter tuning to ensure the model works perfectly. In this case, I used GridSearchCV.

```
1   from sklearn.model_selection import GridSearchCV
2
3   param_grid = {'C': [0.1, 1, 10, 100]}
4   grid = GridSearchCV(LogisticRegression(), param_grid, refit=True, verbose=2)
5   grid.fit(X_train, y_train)
6
```

```
Fitting 5 folds for each of 4 candidates, totalling 20 fits
[CV] END .................................................C=0.1; total time=   0.1s
[CV] END .................................................C=0.1; total time=   0.1s
[CV] END .................................................C=0.1; total time=   0.1s
[CV] END .................................................C=0.1; total time=   0.1s
[CV] END .................................................C=0.1; total time=   0.1s
[CV] END ...................................................C=1; total time=   0.1s
[CV] END ...................................................C=1; total time=   0.1s
[CV] END ...................................................C=1; total time=   0.1s
[CV] END ...................................................C=1; total time=   0.1s
[CV] END ...................................................C=1; total time=   0.1s
[CV] END ..................................................C=10; total time=   0.1s
[CV] END ..................................................C=10; total time=   0.1s
[CV] END ..................................................C=10; total time=   0.1s
[CV] END ..................................................C=10; total time=   0.1s
[CV] END ..................................................C=10; total time=   0.2s
[CV] END .................................................C=100; total time=   0.1s
[CV] END .................................................C=100; total time=   0.1s
[CV] END .................................................C=100; total time=   0.2s
[CV] END .................................................C=100; total time=   0.1s
[CV] END .................................................C=100; total time=   0.2s
```

```
▸         GridSearchCV          ⓘ ⑦

▸ best_estimator_: LogisticRegression

        ▸ LogisticRegression ⑦
```

The model applies a lower penalty for large coefficients, allowing it to fit the data more closely.

```
1 best_model = grid.best_estimator_
2 y_pred = best_model.predict(X_test)
```

## Retrain on the full dataset

```
1 model = LogisticRegression(C=100)
2 model.fit(X_train, y_train)
```

```
▾  LogisticRegression  ⓘ ⑦
LogisticRegression(C=100)
```

## Cross Validation

Cross-validation is a technique used in machine learning to assess how well a model generalizes to unseen data. Instead of training and testing a model once on a single split of the data, cross-validation divides the data into multiple parts (called "folds") and trains and tests the model multiple times. This process provides a more reliable estimate of model performance and helps prevent overfitting.

```
1   # Convert cv_results_ to a DataFrame for easier viewing
2   cv_results = pd.DataFrame(grid.cv_results_)
3
4   # Display relevant columns
5   print("\nCross-Validation Results:")
6   print(cv_results[['param_C', 'mean_test_score', 'std_test_score']])
7
```

```
Cross-Validation Results:
   param_C  mean_test_score  std_test_score
0      0.1         0.911764        0.010151
1      1.0         0.975199        0.006305
2     10.0         0.984637        0.003249
3    100.0         0.985076        0.003768
```

Cross-validation is a technique used in machine learning to assess how well a model generalizes to unseen data. Instead of training and testing a model once on a single split of the data, cross-validation divides the data into multiple parts (called "folds") and trains and tests the model multiple times. This process provides a more reliable estimate of model performance and helps prevent overfitting.

```
1   # Convert cv_results_ to a DataFrame for easier viewing
2   cv_results = pd.DataFrame(grid.cv_results_)
3
4   # Display relevant columns
5   print("\nCross-Validation Results:")
6   print(cv_results[['param_C', 'mean_test_score', 'std_test_score']])
7
```

```
Cross-Validation Results:
   param_C  mean_test_score  std_test_score
0      0.1         0.911764        0.010151
1      1.0         0.975199        0.006305
2     10.0         0.984637        0.003249
3    100.0         0.985076        0.003768
```