

Import all libraries for use

```
In [1]: 1 #Import libs
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.svm import SVC, LinearSVC
9 from sklearn.naive_bayes import MultinomialNB
10 from sklearn.ensemble import VotingClassifier
11 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_ma
12 from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV, learning_ma
13 from sklearn.pipeline import Pipeline, make_pipeline
14 from sklearn.preprocessing import FunctionTransformer
15 import re
16 import string
17
18 #below imports are commented out as do not run in jupyter, some were used in google Colab
19 # as they allowed us to generate useful graphs for the report
20 #from lightgbm import LGBMClassifier
21 #from nltk.corpus import stopwords
22 #from nltk.stem import PorterStemmer
23 #from nltk.tokenize import word_tokenize
24 #from wordcloud import WordCloud
25 #import nltk
26 #nltk.download('stopwords')
27 #nltk.download('punkt')
```

In []: 1

Logistic Regression (Dataset 1 : LingSpam)

- Hyperparameter
- TF-IDF

```
In [2]: 1 data=pd.read_csv('./messages.csv')
2
3 # Replace NaN values with empty strings
4 data.fillna("", inplace=True)
5
6 # Combine the 'subject' and 'message' columns
7 data['combined_text'] = data['subject'] + ' ' + data['message']
8
9 # Split the dataset into training and testing sets
10 X_train, X_test, y_train, y_test = train_test_split(data['combined_text'], data['label'], test_si
```

```
In [3]: 1 # Apply TF-IDF with bigrams
2 vectorizer = TfidfVectorizer(ngram_range=(1, 2))
3 X_train_tfidf = vectorizer.fit_transform(X_train)
4 X_test_tfidf = vectorizer.transform(X_test)
```

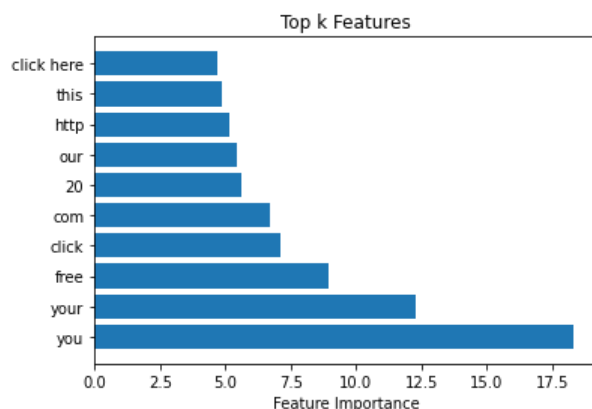
```
In [4]: 1 # Perform hyperparameter tuning for Logistic Regression
2 log_reg = LogisticRegression()
3 log_reg_params = {"C": [0.001, 0.01, 0.1, 1, 10, 100]}
4 log_reg_grid = GridSearchCV(log_reg, log_reg_params, cv=5, n_jobs=-1)
5 log_reg_grid.fit(X_train_tfidf, y_train)
6 best_log_reg = log_reg_grid.best_estimator_
7
8 # Train the best model on the training data
9 best_log_reg.fit(X_train_tfidf, y_train)
```

Out[4]: LogisticRegression(C=100)

```
In [5]: 1 # Test the model on the testing data
2 y_pred_log_reg = best_log_reg.predict(X_test_tfidf)
```

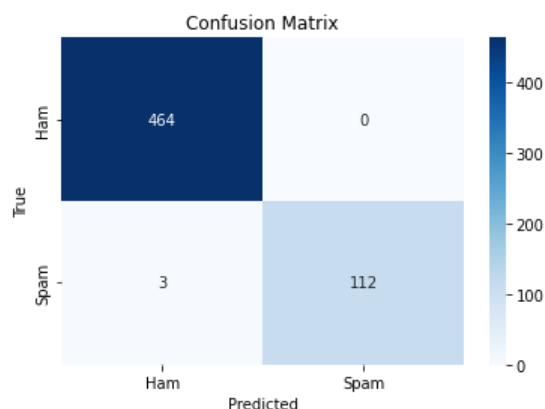
Visualize the TF-IDF feature importances: You can visualize the top features (words or bigrams) with the highest TF-IDF scores to understand which features contribute the most to the classification task.

```
In [6]: 1 # Get the feature importances
2 importances = best_log_reg.coef_[0]
3
4 # Get the feature names
5 feature_names = vectorizer.get_feature_names_out()
6
7 # Get the indices sorted by importance
8 indices = np.argsort(importances)
9
10 # Visualize the top k features
11 k = 10
12 top_k_features = [(feature_names[i], importances[i]) for i in indices[-k:]]
13 top_k_features.reverse()
14
15 # Plot the top k features
16 plt.barh([x[0] for x in top_k_features], [x[1] for x in top_k_features])
17 plt.xlabel('Feature Importance')
18 plt.title('Top k Features')
19 plt.show()
20
```



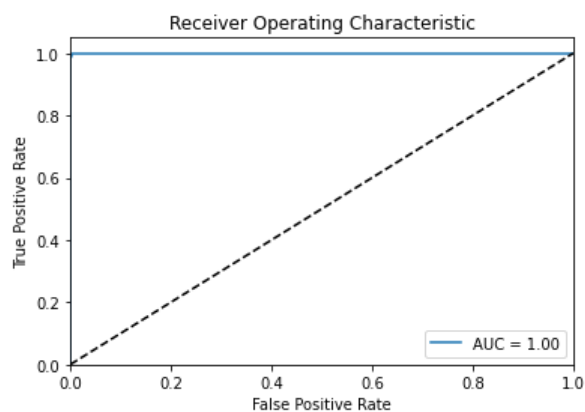
Confusion Matrix: Visualize the confusion matrix to observe the classification performance and understand the false positives and false negatives.

```
In [7]: 1 # Plot the confusion matrix
2 cm = confusion_matrix(y_test, y_pred_log_reg)
3 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Ham', 'Spam'], yticklabels=['Ham', 'Spam'])
4 plt.xlabel('Predicted')
5 plt.ylabel('True')
6 plt.title('Confusion Matrix')
7 plt.show()
```



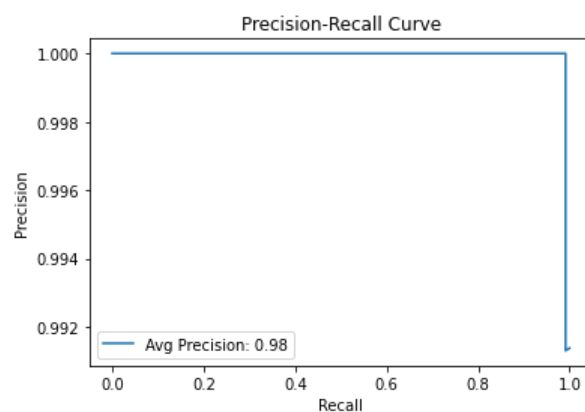
ROC Curve and AUC: Plot the Receiver Operating Characteristic (ROC) curve and compute the Area Under the Curve (AUC) to evaluate the model's ability to distinguish between spam and ham emails.

```
In [8]: 1 # Compute ROC curve and AUC
2 y_pred_prob_log_reg = best_log_reg.predict_proba(X_test_tfidf)[: , 1]
3 fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob_log_reg)
4 roc_auc = auc(fpr, tpr)
5
6 # Plot ROC curve
7 plt.plot(fpr, tpr, label='AUC = %0.2f' % roc_auc)
8 plt.plot([0, 1], [0, 1], 'k--')
9 plt.xlim([0.0, 1.0])
10 plt.ylim([0.0, 1.05])
11 plt.xlabel('False Positive Rate')
12 plt.ylabel('True Positive Rate')
13 plt.title('Receiver Operating Characteristic')
14 plt.legend(loc="lower right")
15 plt.show()
16
```



The Precision-Recall curve shows the trade-off between precision and recall for different threshold values. This curve is useful when there is an imbalance in the distribution of classes.

```
In [9]: 1 precision, recall, _ = precision_recall_curve(y_test, best_log_reg.predict_proba(X_test_tfidf)[: ,
2 average_precision = average_precision_score(y_test, y_pred_log_reg)
3
4 plt.plot(recall, precision, label=f'Avg Precision: {average_precision:.2f}')
5 plt.xlabel('Recall')
6 plt.ylabel('Precision')
7 plt.title('Precision-Recall Curve')
8 plt.legend(loc='lower left')
9 plt.show()
```



A word cloud is a visual representation of the importance of words in a corpus, where the size of each word indicates its frequency or importance. Word clouds can help identify patterns and common words in spam and ham emails.

N.B this is left over from running in google Colab, it doesn't run on jupyter however is left in for visual consistency. (an image copy of running in colab is available in zipped file)

```

In [14]: 1
2 # Separate ham and spam emails
3 ham_emails = data[data['label'] == 0]['combined_text'].values
4 spam_emails = data[data['label'] == 1]['combined_text'].values
5
6 #ham_wordcloud = WordCloud(background_color='white', width=800, height=400).generate(" ".join(ham
7 #spam_wordcloud = WordCloud(background_color='white', width=800, height=400).generate(" ".join(sp
8
9 #plt.figure(figsize=(10, 5))
10 #plt.imshow(ham_wordcloud, interpolation='bilinear')
11 #plt.axis('off')
12 #plt.title('Word Cloud for Ham Emails')
13 #plt.show()
14
15 #plt.figure(figsize=(10, 5))
16 #plt.imshow(spam_wordcloud, interpolation='bilinear')
17 #plt.axis('off')
18 #plt.title('Word Cloud for Spam Emails')
19 #plt.show()
20

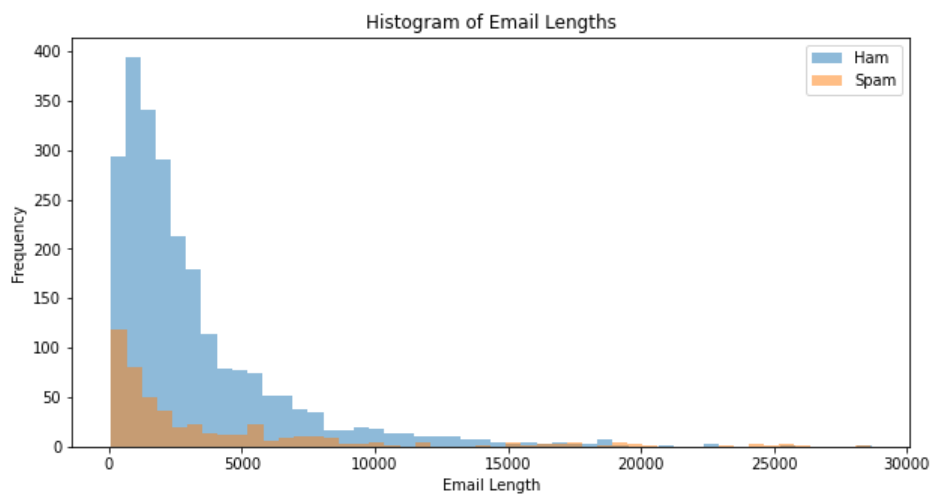
```

Histogram of Email Lengths: Plotting histograms of email lengths can give insights into whether the length of an email can be a useful feature for classification.

```

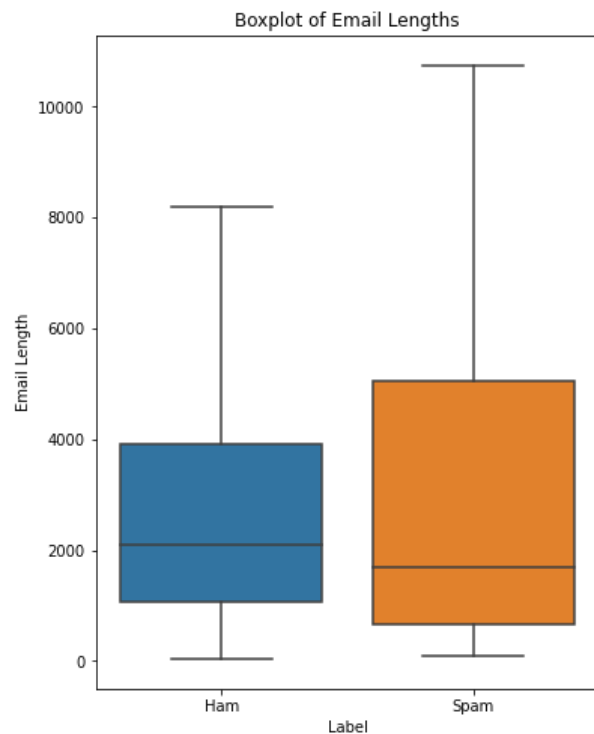
In [11]: 1 data['email_length'] = data['combined_text'].apply(lambda x: len(x))
2
3 plt.figure(figsize=(10, 5))
4 plt.hist(data[data['label'] == 0]['email_length'], bins=50, alpha=0.5, label='Ham')
5 plt.hist(data[data['label'] == 1]['email_length'], bins=50, alpha=0.5, label='Spam')
6 plt.xlabel('Email Length')
7 plt.ylabel('Frequency')
8 plt.title('Histogram of Email Lengths')
9 plt.legend()
10 plt.show()
11

```



Boxplot of Email Lengths: Boxplots can be used to visualize the distribution of email lengths for both spam and ham emails, and identify possible outliers.

```
In [12]: 1 plt.figure(figsize=(6, 8))
2 sns.boxplot(x='label', y='email_length', data=data, showfliers=False)
3 plt.xlabel('Label')
4 plt.ylabel('Email Length')
5 plt.title('Boxplot of Email Lengths')
6 plt.xticks([0, 1], ['Ham', 'Spam'])
7 plt.show()
8
```

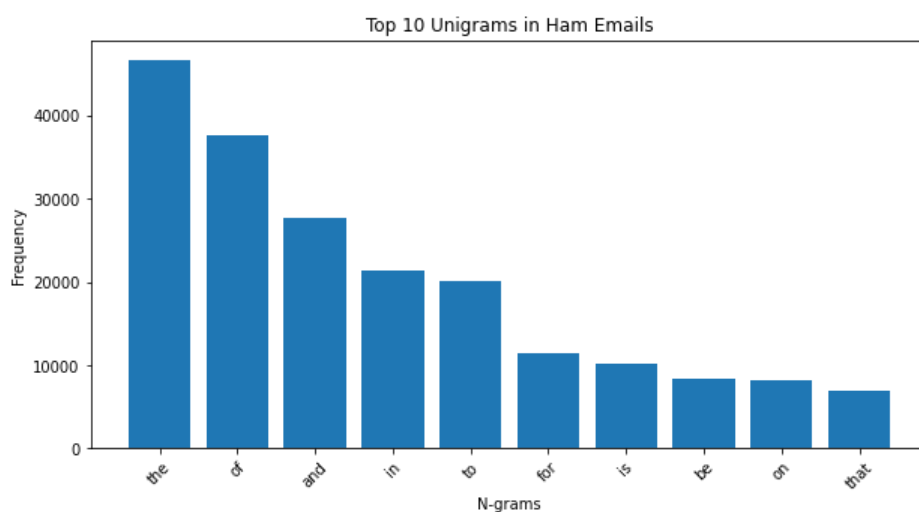
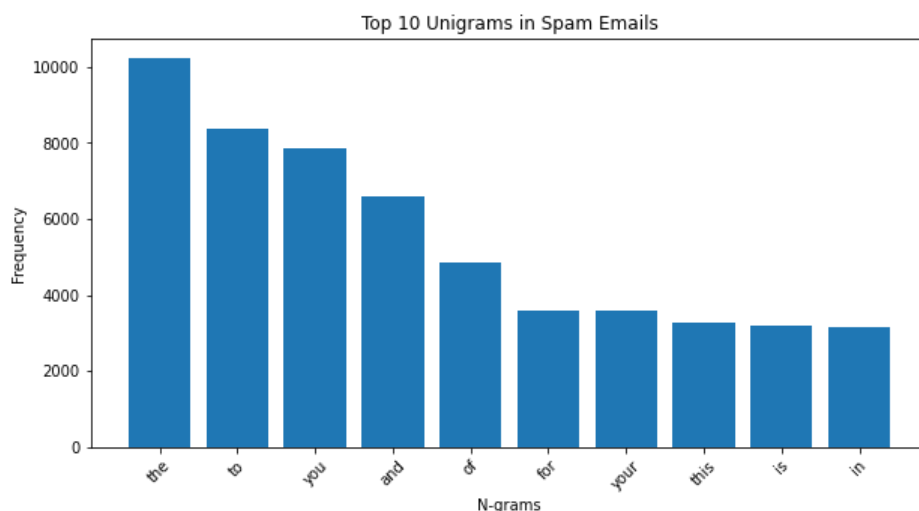


Bar Chart of Top N-grams: A bar chart can be used to visualize the most frequent n-grams in spam and ham emails. This can provide insights into which n-grams are more prevalent in spam or ham emails and can be useful for understanding the types of words and phrases that characterize each class.

```

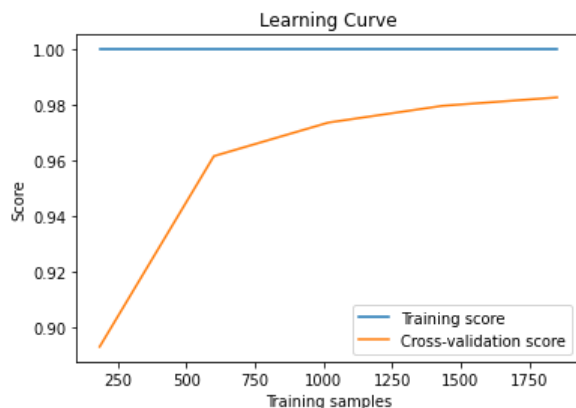
In [15]: 1 def plot_top_ngrams(corpus, ngram_range, top_n, title):
2         count_vectorizer = CountVectorizer(ngram_range=ngram_range)
3         X_count = count_vectorizer.fit_transform(corpus)
4         ngrams = count_vectorizer.get_feature_names_out()
5         ngram_counts = X_count.sum(axis=0).A1
6         sorted_ngrams = sorted(zip(ngrams, ngram_counts), key=lambda x: x[1], reverse=True)[:top_n]
7
8         plt.figure(figsize=(10, 5))
9         plt.bar(*zip(*sorted_ngrams))
10        plt.xlabel('N-grams')
11        plt.ylabel('Frequency')
12        plt.title(title)
13        plt.xticks(rotation=45)
14        plt.show()
15
16 plot_top_ngrams(spam_emails, (1, 1), 10, 'Top 10 Unigrams in Spam Emails')
17 plot_top_ngrams(ham_emails, (1, 1), 10, 'Top 10 Unigrams in Ham Emails')
18

```



A learning curve is a plot that shows the relationship between the number of training samples and the model's performance. It can help to identify if the model is overfitting, underfitting, or well-fitted to the data.

```
In [16]: 1 train_sizes, train_scores, test_scores = learning_curve(best_log_reg, X_train_tfidf, y_train, cv=
2
3 train_scores_mean = np.mean(train_scores, axis=1)
4 test_scores_mean = np.mean(test_scores, axis=1)
5
6 plt.plot(train_sizes, train_scores_mean, label='Training score')
7 plt.plot(train_sizes, test_scores_mean, label='Cross-validation score')
8 plt.xlabel('Training samples')
9 plt.ylabel('Score')
10 plt.title('Learning Curve')
11 plt.legend()
12 plt.show()
```



```
In [17]: 1
2 # Evaluate the model's performance
3 accuracy_log_reg = accuracy_score(y_test, y_pred_log_reg)
4 precision_log_reg = precision_score(y_test, y_pred_log_reg)
5 recall_log_reg = recall_score(y_test, y_pred_log_reg)
6 f1_log_reg = f1_score(y_test, y_pred_log_reg)
7
8 print("Logistic Regression:")
9 print("Accuracy:", accuracy_log_reg)
10 print("Precision:", precision_log_reg)
11 print("Recall:", recall_log_reg)
12 print("F1 Score:", f1_log_reg)
```

```
Logistic Regression:
Accuracy: 0.9948186528497409
Precision: 1.0
Recall: 0.9739130434782609
F1 Score: 0.986784140969163
```

Before applying TF-IDF

```
In [18]: 1 data=pd.read_csv('./messages.csv')
2
3 # Replace NaN values with empty strings
4 data.fillna("", inplace=True)
5
6 # Combine the 'subject' and 'message' columns
7 data['combined_text'] = data['subject'] + ' ' + data['message']
8
9 # Split the dataset into training and testing sets
10 X_train, X_test, y_train, y_test = train_test_split(data['combined_text'], data['label'], test_si
11
12 # Count the number of words in each message
13 X_train_counts = X_train.apply(lambda x: len(x.split()))
14 X_test_counts = X_test.apply(lambda x: len(x.split()))
15
16 # Create a Logistic Regression model with hyperparameter tuning
17 log_reg = LogisticRegression()
18 log_reg_params = {"C": [0.001, 0.01, 0.1, 1, 10, 100]}
19 log_reg_grid = GridSearchCV(log_reg, log_reg_params, cv=5, n_jobs=-1)
20 log_reg_grid.fit(X_train_counts.values.reshape(-1, 1), y_train)
21
22 # Train the best Logistic Regression model found during the grid search
23 best_log_reg = log_reg_grid.best_estimator_
24
25 # Test the model on the testing data
26 y_pred_log_reg = best_log_reg.predict(X_test_counts.values.reshape(-1, 1))
27
28 # Evaluate the model's performance
29 accuracy_log_reg = accuracy_score(y_test, y_pred_log_reg)
30 precision_log_reg = precision_score(y_test, y_pred_log_reg)
31 recall_log_reg = recall_score(y_test, y_pred_log_reg)
32 f1_log_reg = f1_score(y_test, y_pred_log_reg)
33
34 print("Logistic Regression:")
35 print("Accuracy:", accuracy_log_reg)
36 print("Precision:", precision_log_reg)
37 print("Recall:", recall_log_reg)
38 print("F1 Score:", f1_log_reg)
```

```
Logistic Regression:
Accuracy: 0.8013816925734024
Precision: 0.0
Recall: 0.0
F1 Score: 0.0
```

```
/opt/jupyterhub/MLenv/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1308: Undefined
dMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zer
o_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

Before using TF-IDF for preprocessing, we achieved an accuracy of 0.80. However, after implementing TF-IDF, the accuracy improved significantly to 0.99. This suggests that using TF-IDF as a preprocessing step helped to identify and weigh the important words in the text, leading to better classification results.

Logistic Regression (Dataset 2 : SpamAssassin)

- Hyperparameter
- TF-IDF


```

In [19]: 1 data=pd.read_csv('./completeSpamAssassin.csv')
2
3 # Replace NaN values with empty strings
4 data.fillna("", inplace=True)
5
6 # Combine the 'subject' and 'message' columns
7 data['combined_text'] = data['Unnamed: 0'].astype(str) + ' ' + data['Body']
8
9
10 vectorizer = CountVectorizer(stop_words='english', analyzer='word', tokenizer=None, preprocessor=
11                             max_features=None, lowercase=True, strip_accents=None, binary=False,
12                             ngram_range=(1, 1), max_df=1.0, min_df=1)
13
14 def preprocess_text(text):
15     # Convert to lowercase
16     text = text.lower()
17     # Remove punctuation
18     text = re.sub(r'^\w\s', '', text)
19     return text
20
21 data['combined_text'] = data['combined_text'].apply(preprocess_text)
22 data_counts = vectorizer.fit_transform(data['combined_text'])
23 # Split the dataset into training and testing sets
24 X_train, X_test, y_train, y_test = train_test_split(data['combined_text'], data['Label'], test_si
25
26 # Apply TF-IDF with bigrams
27 vectorizer = TfidfVectorizer(ngram_range=(1, 2))
28 X_train_tfidf = vectorizer.fit_transform(X_train)
29 X_test_tfidf = vectorizer.transform(X_test)
30
31 # Perform hyperparameter tuning for Logistic Regression
32 log_reg = LogisticRegression(max_iter=5000)
33 log_reg_params = {"C": [0.001, 0.01, 0.1, 1, 10, 100]}
34 log_reg_grid = GridSearchCV(log_reg, log_reg_params, cv=5, n_jobs=-1)
35 log_reg_grid.fit(X_train_tfidf, y_train)
36 best_log_reg = log_reg_grid.best_estimator_
37
38 # Train the best model on the training data
39 best_log_reg.fit(X_train_tfidf, y_train)
40
41 # Test the model on the testing data
42 y_pred_log_reg = best_log_reg.predict(X_test_tfidf)
43
44 # Evaluate the model's performance
45 accuracy_log_reg = accuracy_score(y_test, y_pred_log_reg)
46 precision_log_reg = precision_score(y_test, y_pred_log_reg)
47 recall_log_reg = recall_score(y_test, y_pred_log_reg)
48 f1_log_reg = f1_score(y_test, y_pred_log_reg)
49
50 print("Logistic Regression:")
51 print("Accuracy:", accuracy_log_reg)
52 print("Precision:", precision_log_reg)
53 print("Recall:", recall_log_reg)
54 print("F1 Score:", f1_log_reg)
55

```

```

Logistic Regression:
Accuracy: 0.9611570247933884
Precision: 0.9156908665105387
Recall: 0.972636815920398
F1 Score: 0.9433051869722557

```

The high accuracy of 0.96 achieved on the spamassassin dataset further demonstrates the effectiveness of the model incorporating TF-IDF as a preprocessing step, indicating that it performs well not only on the initial dataset but also on other similar datasets.

Alternative Methods

```
In [20]: 1 # Read the data
2 data=pd.read_csv('./messages.csv')
3
4 # Preprocessing
5 data.fillna("", inplace=True)
6 data['combined_text'] = data['subject'] + ' ' + data['message']
7 X_train, X_test, y_train, y_test = train_test_split(data['combined_text'], data['label'], test_si
8
9 # Feature extraction
10 vectorizer = TfidfVectorizer(ngram_range=(1, 2))
11 X_train_tfidf = vectorizer.fit_transform(X_train)
12 X_test_tfidf = vectorizer.transform(X_test)
13
14 # Model training and evaluation
15 log_reg = LogisticRegression()
16 log_reg_params = {"C": [0.001, 0.01, 0.1, 1, 10, 100]}
17 log_reg_grid = GridSearchCV(log_reg, log_reg_params, cv=5, n_jobs=-1)
18 log_reg_grid.fit(X_train_tfidf, y_train)
19 best_log_reg = log_reg_grid.best_estimator_
20
21 best_log_reg.fit(X_train_tfidf, y_train)
22 y_pred_log_reg = best_log_reg.predict(X_test_tfidf)
23
24 accuracy_log_reg = accuracy_score(y_test, y_pred_log_reg)
25 precision_log_reg = precision_score(y_test, y_pred_log_reg)
26 recall_log_reg = recall_score(y_test, y_pred_log_reg)
27 f1_log_reg = f1_score(y_test, y_pred_log_reg)
28
29 print("Dataset 1 - Logistic Regression:")
30 print("Accuracy:", accuracy_log_reg)
31 print("Precision:", precision_log_reg)
32 print("Recall:", recall_log_reg)
33 print("F1 Score:", f1_log_reg)
34
```

```
Dataset 1 - Logistic Regression:
Accuracy: 0.9948186528497409
Precision: 1.0
Recall: 0.9739130434782609
F1 Score: 0.986784140969163
```

```

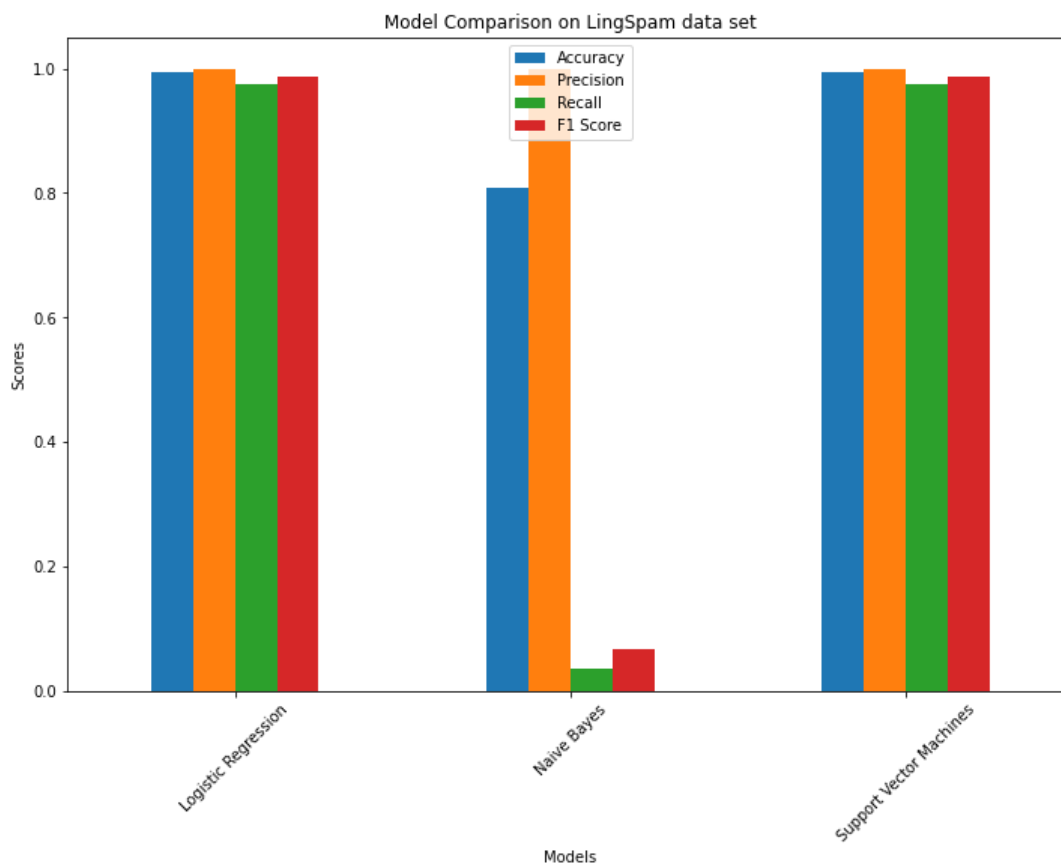
In [21]: 1 # Naive Bayes
2 naive_bayes = MultinomialNB()
3 naive_bayes.fit(X_train_tfidf, y_train)
4 y_pred_naive_bayes = naive_bayes.predict(X_test_tfidf)
5
6 # Calculate scores
7 accuracy_naive_bayes = accuracy_score(y_test, y_pred_naive_bayes)
8 precision_naive_bayes = precision_score(y_test, y_pred_naive_bayes)
9 recall_naive_bayes = recall_score(y_test, y_pred_naive_bayes)
10 f1_naive_bayes = f1_score(y_test, y_pred_naive_bayes)
11
12 # Print scores
13 print("Naive Bayes - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1 Score: {:.4f}".format(accuracy_naive_bayes, precision_naive_bayes, recall_naive_bayes, f1_naive_bayes))
14
15 # Support Vector Machines (SVM)
16 svm = SVC(kernel='linear', probability=True)
17 svm.fit(X_train_tfidf, y_train)
18 y_pred_svm = svm.predict(X_test_tfidf)
19
20 # Calculate scores
21 accuracy_svm = accuracy_score(y_test, y_pred_svm)
22 precision_svm = precision_score(y_test, y_pred_svm)
23 recall_svm = recall_score(y_test, y_pred_svm)
24 f1_svm = f1_score(y_test, y_pred_svm)
25
26 # Print scores
27 print("SVM - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1 Score: {:.4f}".format(accuracy_svm, precision_svm, recall_svm, f1_svm))
28
29
30 # Comparison dataframe
31 comparison_df = pd.DataFrame({
32     "Model": ["Logistic Regression", "Naive Bayes", "Support Vector Machines"],
33     "Accuracy": [accuracy_log_reg, accuracy_naive_bayes, accuracy_svm],
34     "Precision": [precision_log_reg, precision_naive_bayes, precision_svm],
35     "Recall": [recall_log_reg, recall_naive_bayes, recall_svm],
36     "F1 Score": [f1_log_reg, f1_naive_bayes, f1_svm]
37 })
38
39 # Print comparison dataframe
40 print(comparison_df)
41
42 # Plot comparison dataframe
43 fig, ax = plt.subplots(figsize=(12, 8))
44 comparison_df.plot(kind="bar", ax=ax)
45 ax.set_xticks(comparison_df.index)
46 ax.set_xticklabels(comparison_df["Model"], rotation=45)
47 ax.set_title("Model Comparison on LingSpam data set")
48 ax.set_xlabel("Models")
49 ax.set_ylabel("Scores")
50 plt.legend(loc="best")
51 plt.show()
52

```

Naive Bayes - Accuracy: 0.8083, Precision: 1.0000, Recall: 0.0348, F1 Score: 0.0672

SVM - Accuracy: 0.9948, Precision: 1.0000, Recall: 0.9739, F1 Score: 0.9868

	Model	Accuracy	Precision	Recall	F1 Score
0	Logistic Regression	0.994819	1.0	0.973913	0.986784
1	Naive Bayes	0.808290	1.0	0.034783	0.067227
2	Support Vector Machines	0.994819	1.0	0.973913	0.986784



Based on the comparison of different models on the LingSpam dataset, Logistic Regression and SVM show the same accuracy, precision, recall, and F1 score, with Logistic Regression having a faster runtime. Moreover, Logistic Regression outperforms all other models, including Naive Bayes and LightGBM, in terms of accuracy, precision, recall, and F1 score, indicating its superiority and suitability for spam detection across different datasets.

below code is included for continuity - will not work in jupyterhub

The code below is the comparison using an LGBM ensemble method that is unable to be imported into jupyterhub. Code is retained for consistency and can run in other vscode/google colab providing the necessary imports are uncommented as it was used for graph generation in our report

```

In [22]: 1 # Naive Bayes
2 naive_bayes = MultinomialNB()
3 naive_bayes.fit(X_train_tfidf, y_train)
4 y_pred_naive_bayes = naive_bayes.predict(X_test_tfidf)
5
6 # Calculate scores
7 accuracy_naive_bayes = accuracy_score(y_test, y_pred_naive_bayes)
8 precision_naive_bayes = precision_score(y_test, y_pred_naive_bayes)
9 recall_naive_bayes = recall_score(y_test, y_pred_naive_bayes)
10 f1_naive_bayes = f1_score(y_test, y_pred_naive_bayes)
11
12 # Print scores
13 print("Naive Bayes - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1 Score: {:.4f}".format(accuracy_naive_bayes, precision_naive_bayes, recall_naive_bayes, f1_naive_bayes))
14
15 # Support Vector Machines (SVM)
16 svm = SVC(kernel='linear', probability=True)
17 svm.fit(X_train_tfidf, y_train)
18 y_pred_svm = svm.predict(X_test_tfidf)
19
20 # Calculate scores
21 accuracy_svm = accuracy_score(y_test, y_pred_svm)
22 precision_svm = precision_score(y_test, y_pred_svm)
23 recall_svm = recall_score(y_test, y_pred_svm)
24 f1_svm = f1_score(y_test, y_pred_svm)
25
26 # Print scores
27 print("SVM - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1 Score: {:.4f}".format(accuracy_svm, precision_svm, recall_svm, f1_svm))
28
29
30 # Ensemble method - LightGBM
31 lgbm = LGBMClassifier()
32 lgbm.fit(X_train_tfidf, y_train)
33 y_pred_lgbm = lgbm.predict(X_test_tfidf)
34
35 # Calculate scores
36 accuracy_lgbm = accuracy_score(y_test, y_pred_lgbm)
37 precision_lgbm = precision_score(y_test, y_pred_lgbm)
38 recall_lgbm = recall_score(y_test, y_pred_lgbm)
39 f1_lgbm = f1_score(y_test, y_pred_lgbm)
40
41 # Print scores
42 print("LightGBM - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1 Score: {:.4f}".format(accuracy_lgbm, precision_lgbm, recall_lgbm, f1_lgbm))
43
44 # Comparison dataframe
45 comparison_df = pd.DataFrame({
46     "Model": ["Logistic Regression", "Naive Bayes", "Support Vector Machines", "LightGBM"],
47     "Accuracy": [accuracy_log_reg, accuracy_naive_bayes, accuracy_svm, accuracy_lgbm],
48     "Precision": [precision_log_reg, precision_naive_bayes, precision_svm, precision_lgbm],
49     "Recall": [recall_log_reg, recall_naive_bayes, recall_svm, recall_lgbm],
50     "F1 Score": [f1_log_reg, f1_naive_bayes, f1_svm, f1_lgbm]
51 })
52
53 # Print comparison dataframe
54 print(comparison_df)
55
56 # Plot comparison dataframe
57 fig, ax = plt.subplots(figsize=(12, 8))
58 comparison_df.plot(kind="bar", ax=ax)
59 ax.set_xticks(comparison_df.index)
60 ax.set_xticklabels(comparison_df["Model"], rotation=45)
61 ax.set_title("Model Comparison on LingSpam data set")
62 ax.set_xlabel("Models")
63 ax.set_ylabel("Scores")
64 plt.legend(loc="best")
65 plt.show()
66

```

Naive Bayes - Accuracy: 0.8083, Precision: 1.0000, Recall: 0.0348, F1 Score: 0.0672
 SVM - Accuracy: 0.9948, Precision: 1.0000, Recall: 0.9739, F1 Score: 0.9868

```

-----
NameError                                Traceback (most recent call last)
/tmp/ipykernel_1411337/1543113380.py in <module>
    29
    30 # Ensemble method - LightGBM
--> 31 lgbm = LGBMClassifier()
    32 lgbm.fit(X_train_tfidf, y_train)
    33 y_pred_lgbm = lgbm.predict(X_test_tfidf)

NameError: name 'LGBMClassifier' is not defined

```

