```
1 from google.colab import files
2 import pandas as pd
3
4 # Upload the CSV file
5 uploaded = files.upload()
```

Choose Files  spam.csv
  • spam.csv(text/csv) - 503663 bytes, last modified: 9/20/2019 - 100% done
  Saving spam.csv to spam.csv

```
1 # Load the dataset
2 df = pd.read_csv('spam.csv', encoding='latin1')
3
4 # Display the first few rows of the dataframe
5 df.head()
```

|   | v1   | v2                                      | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 |
|---|------|-----------------------------------------|------------|------------|------------|
| 0 | ham  | Go until jurong point, crazy.. Available only ... | NaN | NaN | NaN |
| 1 | ham  | Ok lar... Joking wif u oni...           | NaN | NaN | NaN |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... | NaN | NaN | NaN |
| 3 | ham  | U dun say so early hor... U c already then say... | NaN | NaN | NaN |
| 4 | ham  | Nah I don't think he goes to usf, he lives aro... | NaN | NaN | NaN |

Next steps:   Generate code with df      ○ View recommended plots      New interactive sheet

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.feature_extraction.text import TfidfVectorizer
```

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   v1          5572 non-null   object
 1   v2          5572 non-null   object
 2   Unnamed: 2  50 non-null     object
 3   Unnamed: 3  12 non-null     object
 4   Unnamed: 4  6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB
```

```
1 # Drop unnecessary columns
2 df = df[['v1', 'v2']]
3
4 # Rename columns for convenience
5 df.columns = ['label', 'message']
6
7 # Encode the labels: 'spam' -> 1, 'ham' -> 0
8 df['label'] = df['label'].map({'ham': 0, 'spam': 1})
```

```
<ipython-input-12-f40901fbba71>:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-v
  df['label'] = df['label'].map({'ham': 0, 'spam': 1})
```

```
1 # Check the cleaned data
2 print(df.head())
```

```
   label                                            message
0      0  Go until jurong point, crazy.. Available only ...
1      0                      Ok lar... Joking wif u oni...
2      1  Free entry in 2 a wkly comp to win FA Cup fina...
3      0  U dun say so early hor... U c already then say...
4      0  Nah I don't think he goes to usf, he lives aro...
```

```
1 # Split the data into features and target
2 X = df['message']
3 y = df['label']
4
5 # Split the data into training and test sets
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

```
1 # Initialize the TF-IDF Vectorizer
2 tfidf_vectorizer = TfidfVectorizer(stop_words='english')
3
4 # Fit and transform the training data
5 X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
6
7 # Transform the test data
8 X_test_tfidf = tfidf_vectorizer.transform(X_test)
```

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.svm import SVC
4 from sklearn.metrics import classification_report
```

```
1 # Naive Bayes
2 nb_classifier = MultinomialNB()
3 nb_classifier.fit(X_train_tfidf, y_train)
4 y_pred_nb = nb_classifier.predict(X_test_tfidf)
5 print("Naive Bayes Classification Report:")
6 print(classification_report(y_test, y_pred_nb))
```

```
Naive Bayes Classification Report:
              precision    recall  f1-score   support

           0       0.96      1.00      0.98      1453
           1       1.00      0.75      0.86       219

    accuracy                           0.97      1672
   macro avg       0.98      0.88      0.92      1672
weighted avg       0.97      0.97      0.97      1672
```

```
1 # Logistic Regression
2 lr_classifier = LogisticRegression(max_iter=1000)
3 lr_classifier.fit(X_train_tfidf, y_train)
4 y_pred_lr = lr_classifier.predict(X_test_tfidf)
5 print("Logistic Regression Classification Report:")
6 print(classification_report(y_test, y_pred_lr))
```

```
Logistic Regression Classification Report:
              precision    recall  f1-score   support

           0       0.95      1.00      0.97      1453
           1       0.98      0.65      0.78       219

    accuracy                           0.95      1672
   macro avg       0.96      0.83      0.88      1672
weighted avg       0.95      0.95      0.95      1672
```

```
1 # Support Vector Machines (SVM)
2 svm_classifier = SVC()
3 svm_classifier.fit(X_train_tfidf, y_train)
4 y_pred_svm = svm_classifier.predict(X_test_tfidf)
5 print("SVM Classification Report:")
6 print(classification_report(y_test, y_pred_svm))
```

```
SVM Classification Report:
              precision    recall  f1-score   support

           0       0.98      1.00      0.99      1453
           1       0.99      0.84      0.91       219

    accuracy                           0.98      1672
   macro avg       0.99      0.92      0.95      1672
weighted avg       0.98      0.98      0.98      1672
```
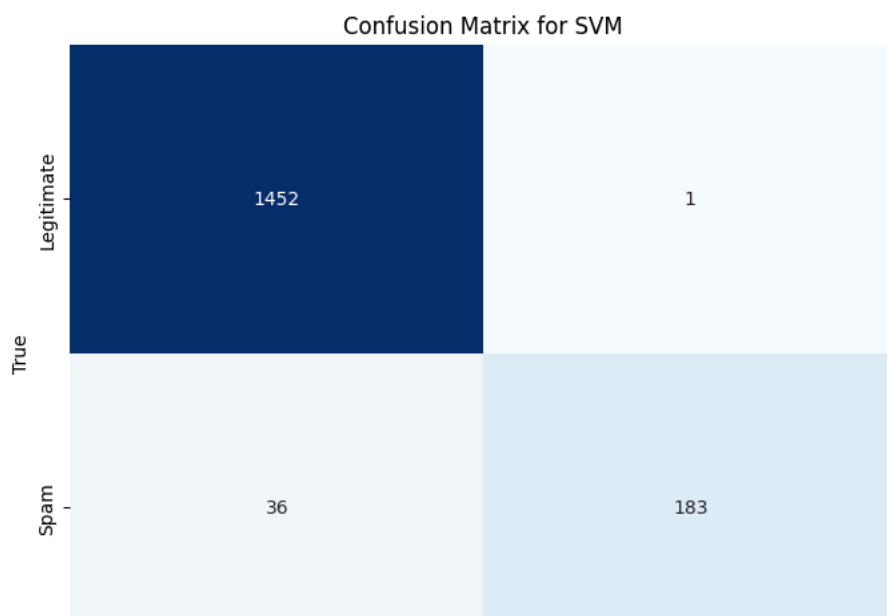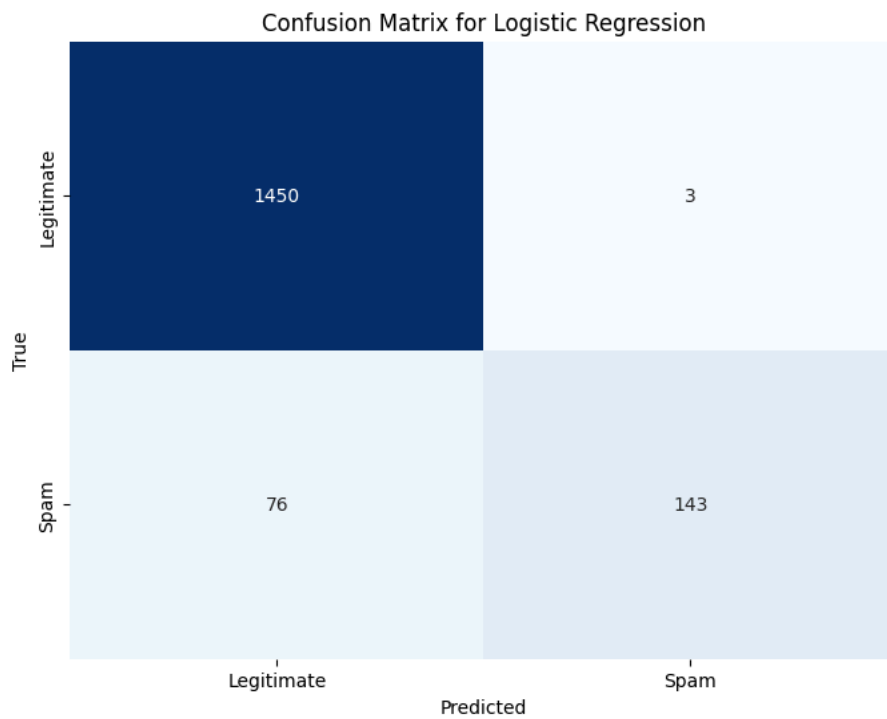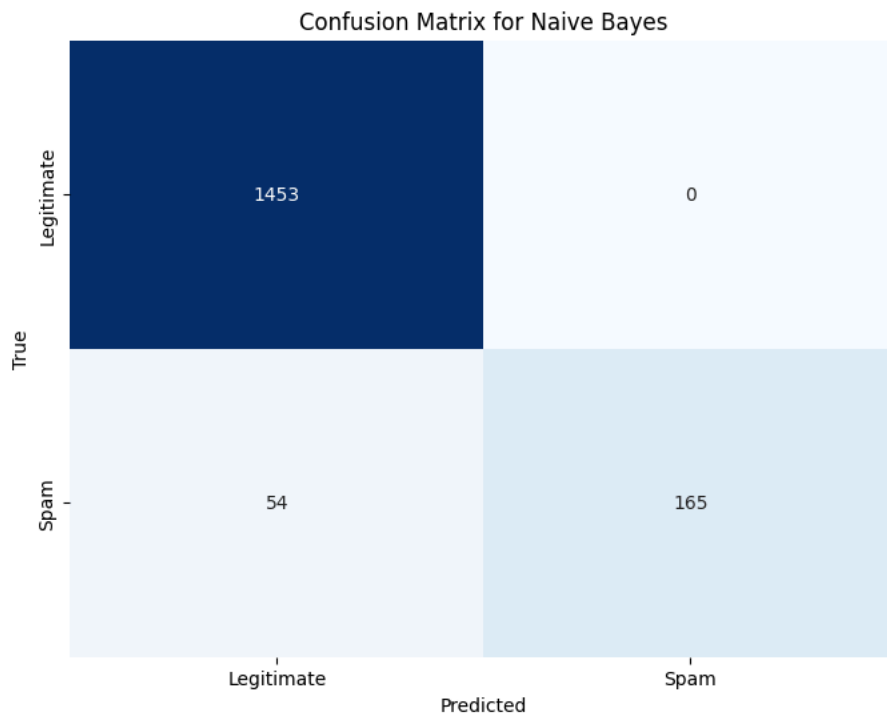
```
1  import matplotlib.pyplot as plt
2  import seaborn as sns
3  from sklearn.metrics import confusion_matrix
4
5  # Function to plot confusion matrix
6  def plot_confusion_matrix(y_true, y_pred, title='Confusion Matrix'):
7      cm = confusion_matrix(y_true, y_pred)
8      plt.figure(figsize=(8, 6))
9      sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False,
10                 xticklabels=['Legitimate', 'Spam'], yticklabels=['Legitimate', 'Spam'])
11     plt.xlabel('Predicted')
12     plt.ylabel('True')
13     plt.title(title)
14     plt.show()
15
16 # Plot confusion matrix for Naive Bayes
17 plot_confusion_matrix(y_test, y_pred_nb, title='Confusion Matrix for Naive Bayes')
18
19 # Plot confusion matrix for Logistic Regression
20 plot_confusion_matrix(y_test, y_pred_lr, title='Confusion Matrix for Logistic Regression')
21
22 # Plot confusion matrix for SVM
23 plot_confusion_matrix(y_test, y_pred_svm, title='Confusion Matrix for SVM')
```

```
1  import matplotlib.pyplot as plt
2  import seaborn as sns
3  from sklearn.metrics import confusion_matrix
4
5  # Function to plot confusion matrix
```

## Confusion Matrix for Naive Bayes



## Confusion Matrix for Logistic Regression



## Confusion Matrix for SVM

Legitimate                                                    Spam
                                        Predicted

```
1  from sklearn.metrics import roc_curve, auc
2
3  # Function to plot ROC curve
4  def plot_roc_curve(y_true, y_score, title='ROC Curve'):
5      fpr, tpr, _ = roc_curve(y_true, y_score)
6      roc_auc = auc(fpr, tpr)
7      plt.figure(figsize=(8, 6))
8      plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
9      plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
10     plt.xlim([0.0, 1.0])
11     plt.ylim([0.0, 1.05])
12     plt.xlabel('False Positive Rate')
13     plt.ylabel('True Positive Rate')
14     plt.title(title)
15     plt.legend(loc='lower right')
16     plt.show()
17
18 # Calculate probabilities for ROC curve
19 y_score_nb = nb_classifier.predict_proba(X_test_tfidf)[:, 1]
20 y_score_lr = lr_classifier.predict_proba(X_test_tfidf)[:, 1]
21 y_score_svm = svm_classifier.decision_function(X_test_tfidf)
22
23 # Plot ROC curves
24 plot_roc_curve(y_test, y_score_nb, title='ROC Curve for Naive Bayes')
25 plot_roc_curve(y_test, y_score_lr, title='ROC Curve for Logistic Regression')
26 plot_roc_curve(y_test, y_score_svm, title='ROC Curve for SVM')
```