

Start coding or [generate](#) with AI.

```
# Import Libraries
# Importing Numpy & Pandas for data processing & data wrangling
import numpy as np
import pandas as pd

# Importing tools for visualization
import matplotlib.pyplot as plt
import seaborn as sns

# Import evaluation metric libraries
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, roc_curve

# Word Cloud library
from wordcloud import WordCloud, STOPWORDS

# Library used for data preprocessing
from sklearn.feature_extraction.text import CountVectorizer

# Import model selection libraries
from sklearn.model_selection import train_test_split

# Library used for ML Model implementation
from sklearn.naive_bayes import MultinomialNB

# Importing the Pipeline class from scikit-learn
from sklearn.pipeline import Pipeline

# Library used for ignore warnings
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

```
# Load Dataset from github repository
df = pd.read_csv("https://raw.githubusercontent.com/Apaulgithub/oibsip_taskno4/main/spam.csv", encoding='ISO-8859-1')
```

```
# Dataset Columns
df.columns
```

```
Index(['v1', 'v2', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], dtype='object')
```

```
# Dataset Describe (all columns included)
df.describe(include= 'all').round(2)
```

| | v1 | v2 | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 |
|---------------|------|------------------------|--|----------------------|------------|
| count | 5572 | 5572 | 50 | 12 | 6 |
| unique | 2 | 5169 | 43 | 10 | 5 |
| top | ham | Sorry, I'll call later | bt not his girfrnd... G o o d n i g h t . . . @" | MK17 92H. 450Ppw 16" | GNT:-)" |
| freq | 4825 | 30 | 3 | 2 | 2 |

```
# Check Unique Values for each variable using a for loop.
for i in df.columns.tolist():
    print("No. of unique values in",i,"is",df[i].nunique())
```

```
No. of unique values in v1 is 2
No. of unique values in v2 is 5169
No. of unique values in Unnamed: 2 is 43
No. of unique values in Unnamed: 3 is 10
No. of unique values in Unnamed: 4 is 5
```

```
# Change the v1 & v2 columns as Category and Message
df.rename(columns={"v1": "Category", "v2": "Message"}, inplace=True)
```

```
# Removing the all unnamed columns (its include much number of missing values)
df.drop(columns=['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], inplace=True)
```

```
# Create a binary 'Spam' column: 1 for 'spam' and 0 for 'ham', based on the 'Category' column.
df['Spam'] = df['Category'].apply(lambda x: 1 if x == 'spam' else 0)
```

```
# Updated new dataset
df.head()
```

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

Next steps:

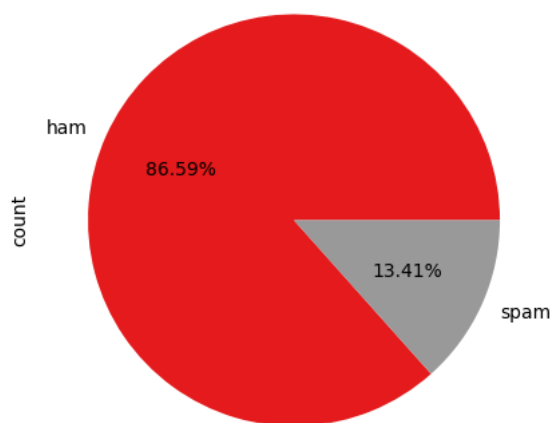
[Generate code with df](#)[New interactive sheet](#)

```
# Chart - 1 Pie Chart Visualization Code For Distribution of Spam vs Ham Messages
spread = df['Category'].value_counts()
plt.rcParams['figure.figsize'] = (5,5)

# Set Labels
spread.plot(kind = 'pie', autopct='%1.2f%%', cmap='Set1')
plt.title(f'Distribution of Spam vs Ham')

# Display the Chart
plt.show()
```

Distribution of Spam vs Ham



```
# Splitting Spam Messages
df_spam = df[df['Category']=='spam'].copy()
```

Start coding or [generate](#) with AI.

```
# Chart - 2 WordCloud Plot Visualization Code For Most Used Words in Spam Messages
# Create a String to Store All The Words
comment_words = ''

# Remove The Stopwords
stopwords = set(STOPWORDS)

# Iterate Through The Column
for val in df_spam.Message:

    # Typecaste Each Val to String
    val = str(val)

    # Split The Value
    tokens = val.split()

    # Converts Each Token into lowercase
    for i in range(len(tokens)):
        tokens[i] = tokens[i].lower()

    comment_words += " ".join(tokens)+" "

# Set Parameters
wordcloud = WordCloud(width = 1000, height = 500,
                      background_color = 'white',
                      stopwords = stopwords,
```

Most Used Words In Spam Messages



Start coding or generate with AI.

<https://colab.research.google.com/drive/1pLoaCF0ky6KDVRpNkoIAQMhsJNurmXxE?authuser=0#scrollTo=aqRDNDyr5F5a&printMode=true>

```

fig, ax = plt.subplots(1, 2, figsize=(11,4))

print("\nConfusion Matrix:")
sns.heatmap(cm_train, annot=True, xticklabels=['Negative', 'Positive'], yticklabels=['Negative', 'Positive'], cmap="Orange")
ax[0].set_xlabel("Predicted Label")
ax[0].set_ylabel("True Label")
ax[0].set_title("Train Confusion Matrix")

sns.heatmap(cm_test, annot=True, xticklabels=['Negative', 'Positive'], yticklabels=['Negative', 'Positive'], cmap="Orange")
ax[1].set_xlabel("Predicted Label")
ax[1].set_ylabel("True Label")
ax[1].set_title("Test Confusion Matrix")

plt.tight_layout()
plt.show()

# calculate classification report
cr_train = classification_report(y_train, y_pred_train, output_dict=True)
cr_test = classification_report(y_test, y_pred_test, output_dict=True)
print("\nTrain Classification Report:")
crt = pd.DataFrame(cr_train).T
print(crt.to_markdown())
# sns.heatmap(pd.DataFrame(cr_train).T.iloc[:, :-1], annot=True, cmap="Blues")
print("\nTest Classification Report:")
crt2 = pd.DataFrame(cr_test).T
print(crt2.to_markdown())
# sns.heatmap(pd.DataFrame(cr_test).T.iloc[:, :-1], annot=True, cmap="Blues")

precision_train = cr_train['weighted avg']['precision']
precision_test = cr_test['weighted avg']['precision']

recall_train = cr_train['weighted avg']['recall']
recall_test = cr_test['weighted avg']['recall']

acc_train = accuracy_score(y_true = y_train, y_pred = y_pred_train)
acc_test = accuracy_score(y_true = y_test, y_pred = y_pred_test)

F1_train = cr_train['weighted avg']['f1-score']
F1_test = cr_test['weighted avg']['f1-score']

model_score = [precision_train, precision_test, recall_train, recall_test, acc_train, acc_test, roc_auc_train, roc_auc_test]
return model_score

```

```

# ML Model - 1 Implementation
# Create a machine learning pipeline using scikit-learn, combining text vectorization (CountVectorizer)
# and a Multinomial Naive Bayes classifier for email spam detection.
clf = Pipeline([
    ('vectorizer', CountVectorizer()), # Step 1: Text data transformation
    ('nb', MultinomialNB()) # Step 2: Classification using Naive Bayes
])

# Model is trained (fit) and predicted in the evaluate model

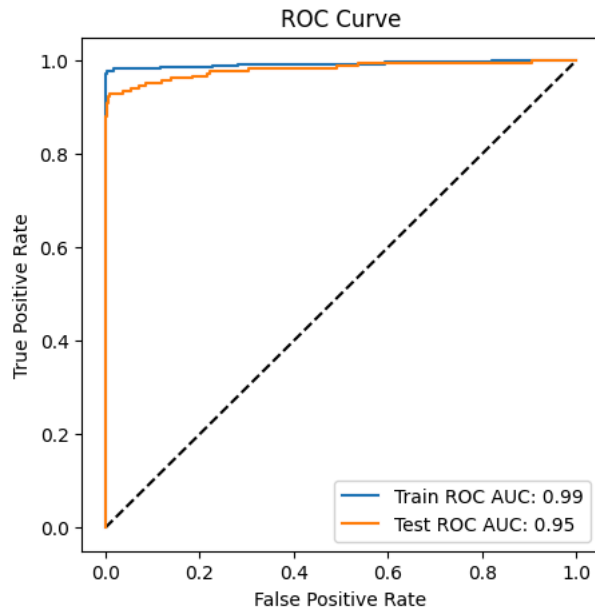
```

```

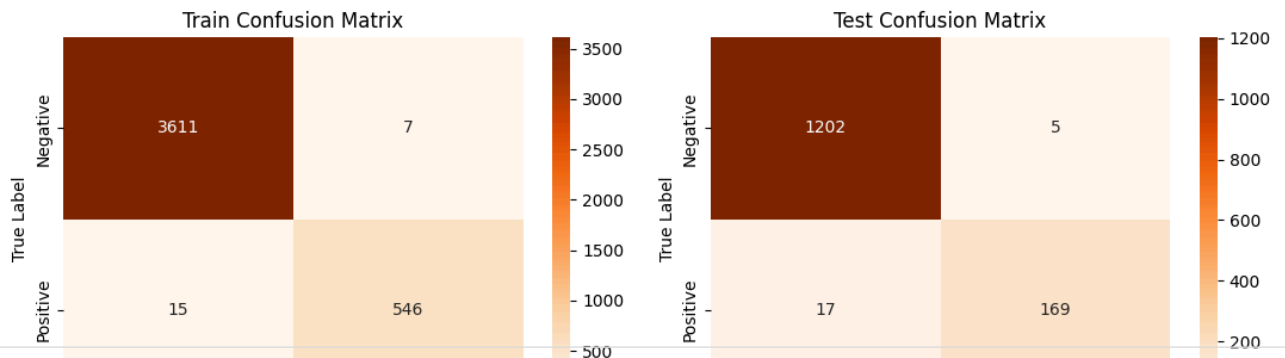
# Visualizing evaluation Metric Score chart
MultinomialNB_score = evaluate_model(clf, X_train, X_test, y_train, y_test)

```

Train ROC AUC: 0.9856636307470371
 Test ROC AUC: 0.9522298242331918



Confusion Matrix:



```
# Defining a function for the Email Spam Detection System
def detect_spam(email_text):
    # Load the trained classifier (clf) here
    # Replace the comment with your code to load the classifier model

    # Make a prediction using the loaded classifier
    prediction = clf.predict([email_text])

    if prediction == 0:
        return "This is a Ham Email!"
    else:
        return "This is a Spam Email!"
```

Test Classification Report:

```
sample_email = 'Bonus point'
result = detect_spam(sample_email)
print(result)
```

| | | | | |
|--------------|----------|----------|----------|----------|
| accuracy | 0.984207 | 0.984207 | 0.984207 | 0.984207 |
| macro avg | 0.978659 | 0.95223 | 0.96491 | 1393 |
| weighted avg | 0.984079 | 0.984207 | 0.983983 | 1393 |