

DATA SCIENCE PROJECT

Email/SMS Spam Classifier

Abstract:

The Email/SMS Spam Classifier project aims to develop a robust and accurate system for identifying spam messages in both email and SMS communication channels. With the exponential growth of digital communication, the problem of spam messages has become increasingly prevalent and poses a significant challenge for individuals and organizations. This project leverages machine learning techniques and natural language processing to create a classifier capable of distinguishing between legitimate and spam messages, enabling users to effectively filter unwanted content. To accomplish this goal, the project utilizes a diverse dataset comprising a wide range of legitimate and spam messages. The dataset is preprocessed to extract relevant features from the text, including word frequency, character patterns, and structural attributes. Various machine learning algorithms, such as Naive Bayes, Support Vector Machines, and ensemble methods, are trained and evaluated using the dataset to identify the most effective approach for spam detection. Additionally, advanced natural language processing techniques are employed to enhance the performance of the classifier. These techniques include text tokenization, stemming, and feature engineering, which provide a deeper understanding of the content and context of messages. The project also explores the utilization of domain-specific knowledge, such as known spam keywords and sender reputation, to further improve the classifier's accuracy. The evaluation of the developed classifier involves rigorous testing on a separate test dataset, measuring key performance metrics such as precision, recall, and F1-score. The results are compared against existing spam detection systems to determine the effectiveness and efficiency of the proposed approach. The Email/SMS Spam Classifier project has the potential to significantly alleviate the problem of unwanted messages by providing an intelligent and adaptable solution. By accurately identifying spam messages, users can save time, protect themselves from malicious content, and maintain a clean and secure communication environment. Furthermore, the project's findings and insights can contribute to the ongoing research and development of anti-spam technologies, benefiting both individuals and organizations in the digital age.

PROJECT PRESENTATION:

DATA SCIENCE PROJECT

Email/SMS Spam Classifier

TEAM MEMBERS

HEMA HARIHARAN S
DHATCHAYANI U

- ❖ INTRODUCTION
- ❖ DATA COLLECTION
- ❖ FEATURE EXTRACTION
- ❖ MODEL BUILDING
- ❖ EVALUATION
- ❖ CONCLUSION



Introduction:

Email and SMS spam has become a major issue in today's digital world. These unsolicited messages can be annoying, time-consuming, and even harmful. The good news is that data science can help us tackle this problem by building an effective spam classifier. A spam classifier is a machine learning model that can distinguish between spam and legitimate messages. By analyzing various features of the message such as sender, subject, content, etc., the classifier can predict whether a message is spam or not.



Data Collection

The first step in building a spam classifier is to collect data. This data should include both spam and legitimate messages. There are many publicly available datasets that can be used for this purpose. Some popular ones include the Enron email dataset and the SMS Spam Collection dataset. Once the data is collected, it needs to be preprocessed. This involves cleaning the data, removing any irrelevant information, and converting it into a format that can be used by the machine learning algorithm.





Feature Extraction

After preprocessing the data, the next step is to extract features from it. Features are the characteristics of the message that the classifier will use to make predictions. Some common features used in spam classification include the sender's email address, the subject line, the presence of certain keywords, and the length of the message. Feature extraction is an important step because it determines the quality of the classifier. The more relevant and informative features we extract, the better our classifier will perform.



Model Building

Once the features are extracted, we can start building our spam classifier model. There are many machine learning algorithms that can be used for this purpose, including Naive Bayes, Support Vector Machines (SVM), and Random Forests.

The model is trained on the preprocessed data with extracted features. During training, the algorithm learns to distinguish between spam and legitimate messages based on the provided features. Once the model is trained, it can be used to predict whether a new message is spam or not.



Evaluation

After building the model, it is important to evaluate its performance. This involves testing the model on a separate dataset that was not used during training. The performance of the model is measured using metrics such as accuracy, precision, recall, and F1 score. If the model's performance is not satisfactory, we can tweak the feature extraction process or try different machine learning algorithms until we get the desired results.



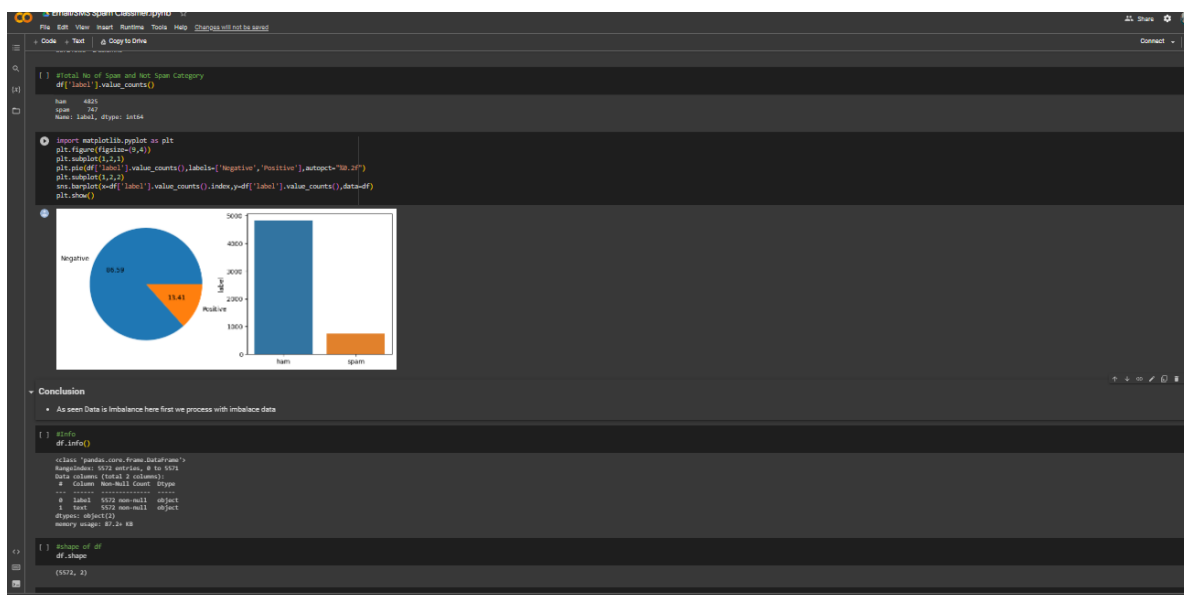
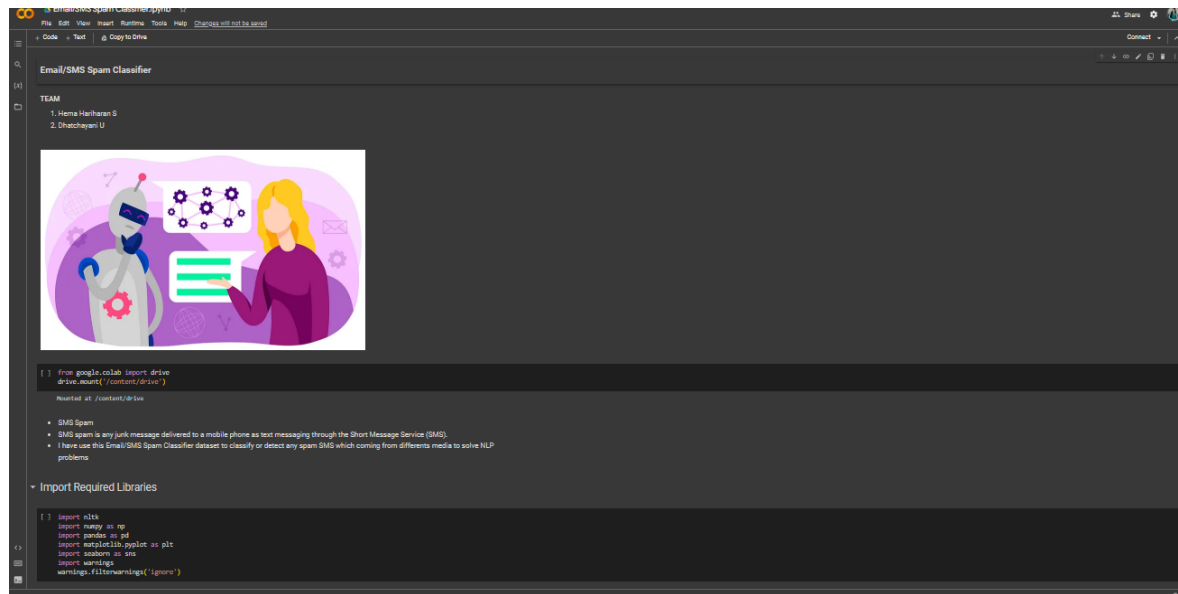
Conclusion

In conclusion, building an email and SMS spam classifier using data science is an effective way to tackle the problem of unsolicited messages. By collecting data, preprocessing it, extracting relevant features, building a model, and evaluating its performance, we can create a classifier that can accurately distinguish between spam and legitimate messages.

This technology can be used by individuals, businesses, and organizations to reduce the amount of time and resources spent on dealing with spam messages.



PROJECT SCREENSHOTS:



```

EmailSMS Spam Classifier.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
Code Text Copy to Drive Connect
(55/72, 2)

[] Check for Null Values
df.isnull().sum()

label: 0
text: 0
dtype: int64

[] df['label'].value_counts()/df.shape[0]*100

label
ham    86.55543
spam   13.44457
Name: label, dtype: float64

1.1 Feature Engineering

We Create Separate Features to extract Some Information

1. Total No of Characters
2. Total No of Words
3. Total No of Sentences

#1. Total No of Char
df['num_char']=df['text'].apply(len)

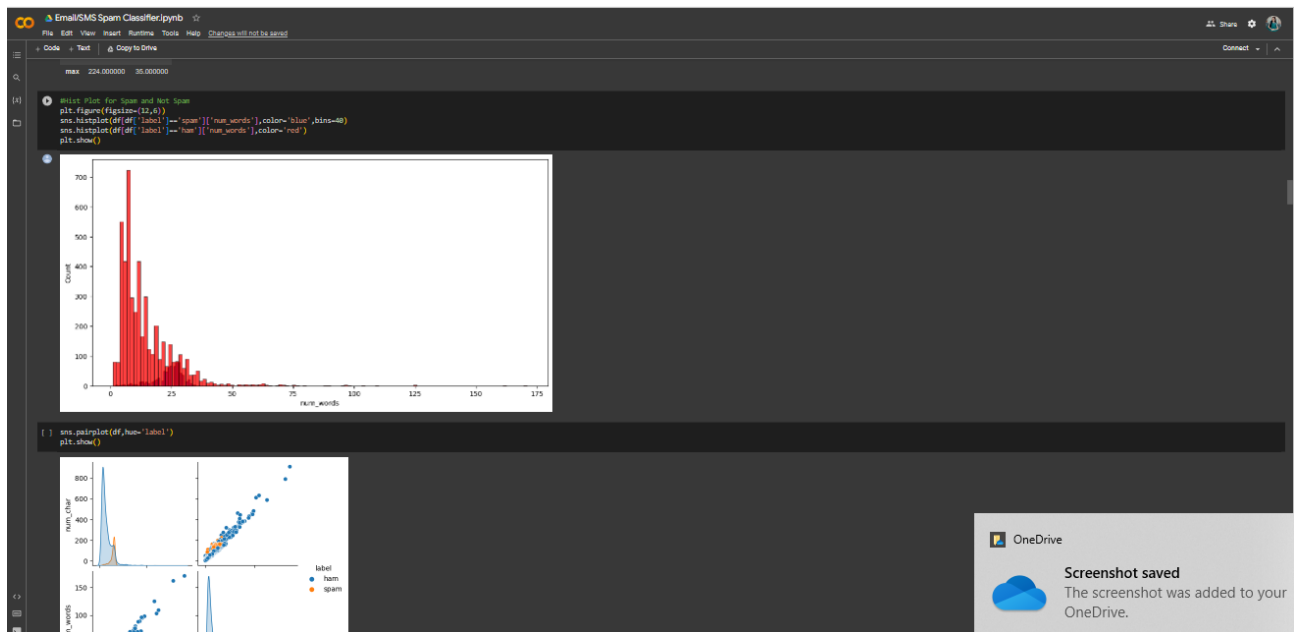
#2. Total No of Words
df['num_words']=df['text'].apply(lambda x: len(str(x).split()))

#3. Total No of Sentences
nltk.download('punkt')
df['num_sen']=df['text'].apply(lambda x: len(nltk.sent_tokenize(x)))

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

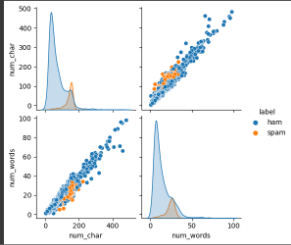
df
  label  text  num_char  num_words
0     ham  Go until juring point, crazy... Available only ...    111     20
1     ham  Ok lar... Joking wif u oni...         29      6
2    spam  Free entry in 2 a wky comp to win PA Cup final...    105     28
3     ham  U dun say so early hor... U s already then say...     49     11
4     ham  Nah i dont think he goes to ust, he lives and...     61     13
...    ...
6687  spam  This is the 2nd time we have tried 2 contact u...    161     30
6688  ham   Will u going to exstinate? home?         37      8
6689  ham  Pff, i was in mood for that. So... any other s...     57     10

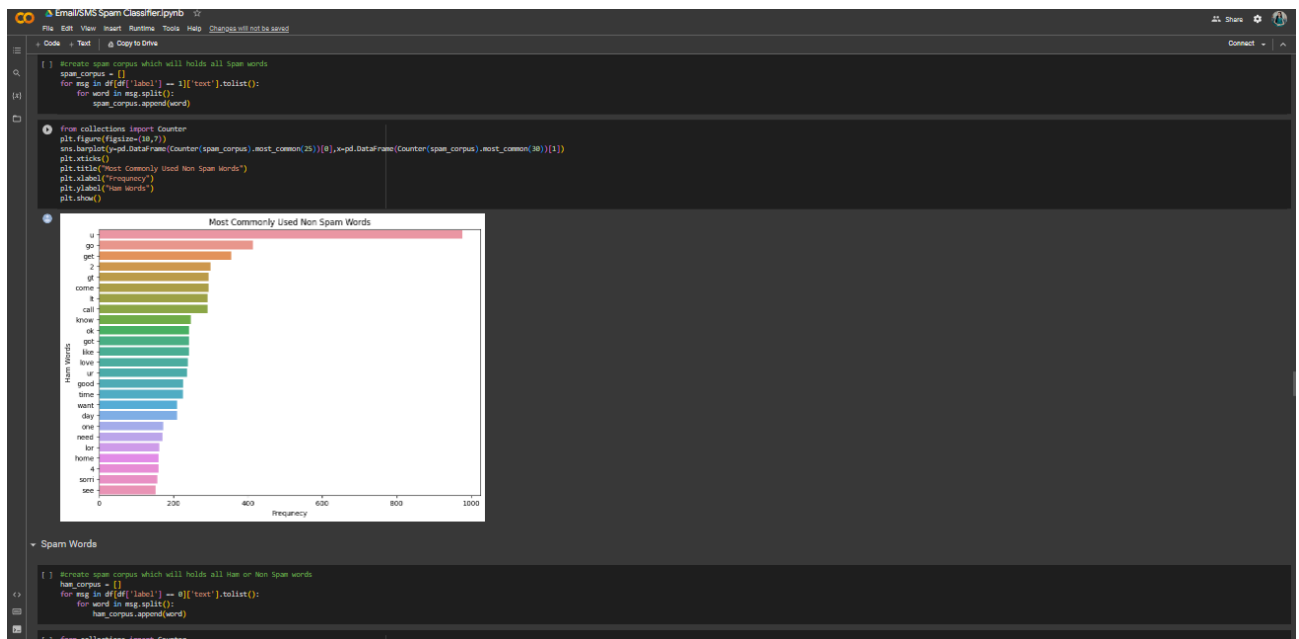
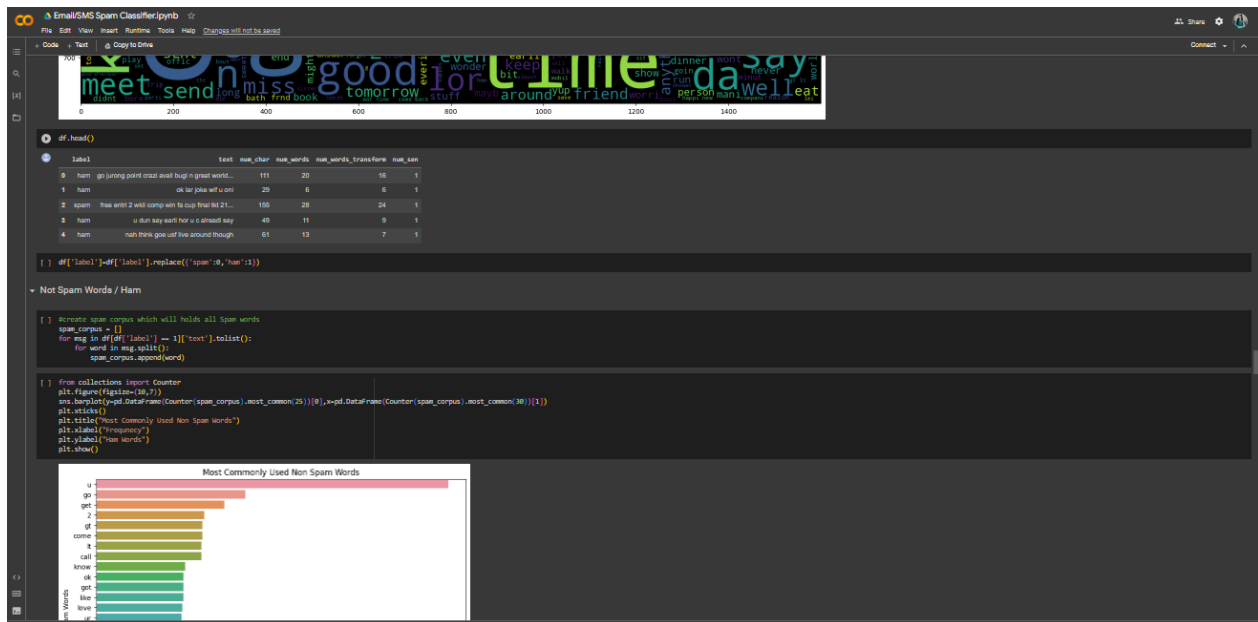
```

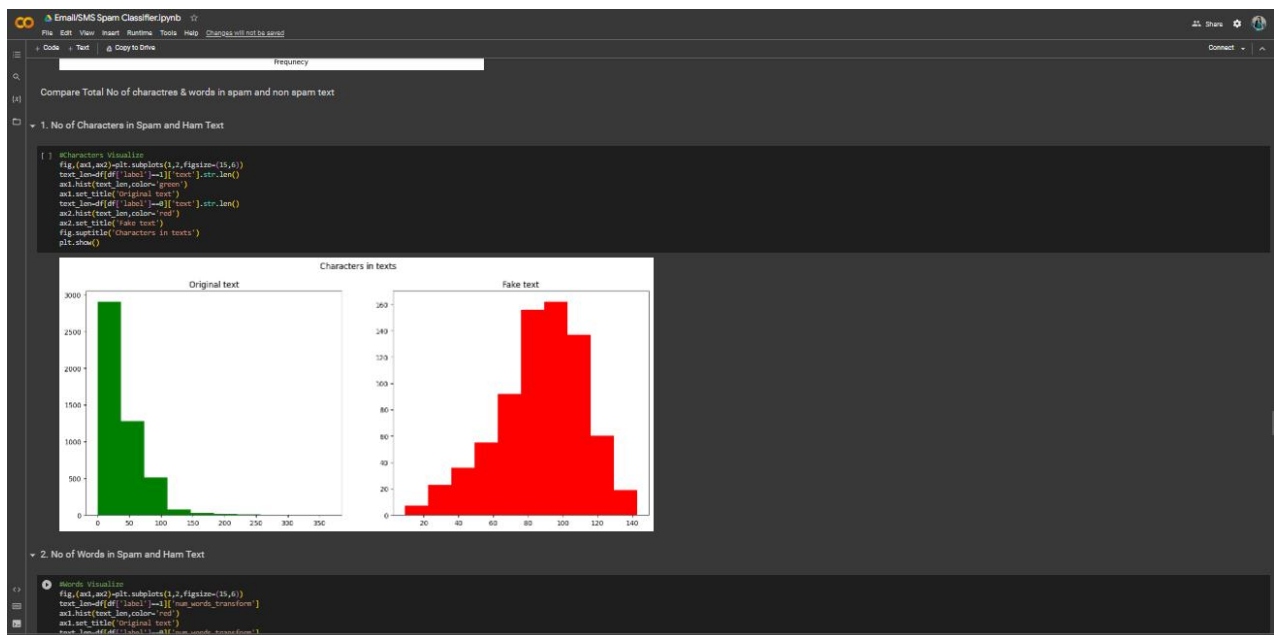
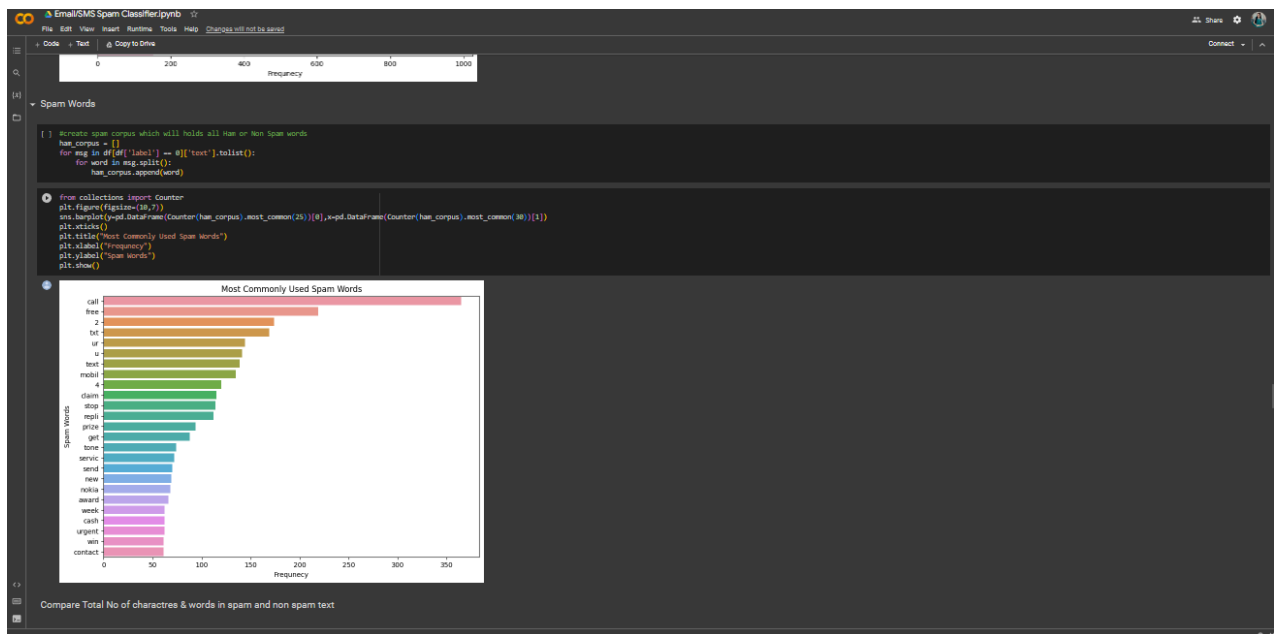


OneDrive

Screenshot saved
The screenshot was added to your OneDrive.

[illegible]







```

EmailSMS Spam Classifier.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
+ Code + Text @ Copy to Drive
Connect

Conclusion
After visualize we can conclude that spam text has more words and characters as compare to Ham text
Average characters includes in spam SMS is around 90 characters
Average words includes in spam SMS is around 15 words

Model Building

Text Vectorization

[ ] #Text Vectorization
from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer
cv = CountVectorizer()
tfidf = TfidfVectorizer(max_features=1000)

Input Output Selection

[ ] #Input and Output Features
X = tfidf.fit_transform(df['text']).toarray()
y = df['label'].values

[ ] X.shape
(5540, 2000)

[ ] y.shape
(5540,1)

Train Test Split

[ ] from sklearn.model_selection import train_test_split

[ ] X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=2)

Imports Libraries

[ ] #Model Training
from sklearn.naive_bayes import GaussianNB,MultinomialNB,BernoulliNB
from sklearn.metrics import accuracy_score,confusion_matrix,precision_score

[ ] !pip install --upgrade scikit-learn

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.2.2)
Requirement already satisfied: numpy>=1.7.3 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.22.4)
Requirement already satisfied: scipy>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.10.1)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.2.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.1.0)

[ ]
```

```

EmailSMS Spam Classifier.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
+ Code + Text @ Copy to Drive
Connect

[ ] X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,random_state=2)

Imports Libraries

[ ] #Model Training
from sklearn.naive_bayes import GaussianNB,MultinomialNB,BernoulliNB
from sklearn.metrics import accuracy_score,confusion_matrix,precision_score

[ ] !pip install --upgrade scikit-learn

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.2.2)
Requirement already satisfied: numpy>=1.7.3 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.22.4)
Requirement already satisfied: scipy>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.10.1)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.2.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.1.0)

[ ] #Importing lib
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.naive_bayes import MultinomialNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import ExtraTreeClassifier
from sklearn.ensemble import GradientBoostingClassifier
from xgboost import XGBClassifier
from sklearn.metrics import precision_score, recall_score, classification_report, accuracy_score, f1_score
from sklearn.model_selection import cross_val_score
from sklearn.metrics import plot_confusion_matrix
from sklearn.model_selection import cross_val_score

[ ] #GaussianNB
gnb = GaussianNB()
gnb.fit(X_train,y_train)
y_pred = gnb.predict(X_test)
print("Accuracy Score :",accuracy_score(y_test,y_pred))
print("Precision Score :",precision_score(y_test,y_pred))
print(confusion_matrix(y_test,y_pred))

Accuracy Score = 0.84601106110606
Precision Score = 0.8768181111111111
[[40 25]
 [48 81]]

[ ] #MultinomialNB
mnb = MultinomialNB()
mnb.fit(X_train,y_train)
y_pred = mnb.predict(X_test)
print("Accuracy Score :",accuracy_score(y_test,y_pred))
print("Precision Score :",precision_score(y_test,y_pred))
print(confusion_matrix(y_test,y_pred))
```

```
print(confusion_matrix(y_test, y_pred))

Accuracy Score = 0.840213449592823
Precision Score = 0.879318282512825
[[148  22]
 [148  81]]

[ ] #multinomial
mb = MultinomialNB()
mb.fit(X_train, y_train)
y_pred = mb.predict(X_test)
print("Accuracy Score : ", accuracy_score(y_test, y_pred))
print("Precision Score : ", precision_score(y_test, y_pred))
print(confusion_matrix(y_test, y_pred))

Accuracy Score = 0.879318282512825
Precision Score = 0.879318282512825
[[148  22]
 [ 148  81]]

[ ] classifiers={"svc": SVC(kernel="sigmoid", gamma=1.0),
               "knn": KNeighborsClassifier(),
               "nb": MultinomialNB(),
               "dtr": DecisionTreeClassifier(max_depth=5),
               "lr": LogisticRegression(solver="liblinear", penalty="l1"),
               "rf": RandomForestClassifier(n_estimators=50, random_state=2),
               "ada": AdaBoostClassifier(n_estimators=50, random_state=2),
               "gb": GradientBoostingClassifier(n_estimators=50, random_state=2)}

for (model in classifiers.items()): cv_score=cross_val_score(model, X_train, y_train, scoring="accuracy", cv=5) print("Mean Accuracy : ",
"cv_score.mean().round(3))

Model Evaluation

[ ] #common function for model train
def train_classifier(clf, X_train, y_train, X_test, y_test):
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)

    accuracy = accuracy_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred)
    train_accuracy = clf.score(X_train, y_train)

    return accuracy, precision, train_accuracy

[ ] svc=SVC(kernel="sigmoid", gamma=1.0)
train_classifier(svc, X_train, y_train, X_test, y_test)

(0.840213449592823, 0.8891282512825, 0.8491844871878)

[ ] accuracy_scores = []
precision_scores = []
train_accuracy_scores = []
```

```
[ ] classifiers={"svc": SVC(kernel="sigmoid", gamma=1.0),
               "knn": KNeighborsClassifier(),
               "nb": MultinomialNB(),
               "dtr": DecisionTreeClassifier(max_depth=5),
               "lr": LogisticRegression(solver="liblinear", penalty="l1"),
               "rf": RandomForestClassifier(n_estimators=50, random_state=2),
               "ada": AdaBoostClassifier(n_estimators=50, random_state=2),
               "gb": GradientBoostingClassifier(n_estimators=50, random_state=2)}

for (model in classifiers.items()): cv_score=cross_val_score(model, X_train, y_train, scoring="accuracy", cv=5) print("Mean Accuracy : ",
"cv_score.mean().round(3))

Model Evaluation

[ ] #common function for model train
def train_classifier(clf, X_train, y_train, X_test, y_test):
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X_test)

    accuracy = accuracy_score(y_test, y_pred)
    precision = precision_score(y_test, y_pred)
    train_accuracy = clf.score(X_train, y_train)

    return accuracy, precision, train_accuracy

[ ] svc=SVC(kernel="sigmoid", gamma=1.0)
train_classifier(svc, X_train, y_train, X_test, y_test)

(0.840213449592823, 0.8891282512825, 0.8491844871878)

[ ] accuracy_scores = []
precision_scores = []
train_accuracy_scores = []

for name, clf in classifiers.items():
    current_accuracy, current_precision, current_train_score = train_classifier(clf, X_train, y_train, X_test, y_test)

    print("For ", name)
    print("Accuracy : ", current_accuracy)
    print("Precision : ", current_precision)

    accuracy_scores.append(current_accuracy)
    precision_scores.append(current_precision)
    train_accuracy_scores.append(current_train_score)
    print()

For svc
Accuracy = 0.840213449592823
Precision = 0.8891282512825
For knn
```



```

EmailSMS Spam Classifier.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
Code Text Copy to Drive
(0.0882113464991821, 0.0891128281128281, 0.08495848871878)

accuracy_scores = []
precision_scores = []
train_accuracy_score = []

for name, clf in classifiers.items():
    current_accuracy, current_precision, current_train_score = train_classifier(clf, X_train, y_train, X_test, y_test)
    print(f'For {name}')
    print(f'Accuracy - {current_accuracy}')
    print(f'Precision - {current_precision}')
    accuracy_scores.append(current_accuracy)
    precision_scores.append(current_precision)
    train_accuracy_score.append(current_train_score)
    print()

For svm
Accuracy - 0.9882113464991821
Precision - 0.9891128281128281

For mnb
Accuracy - 0.98754835487387
Precision - 0.983811828741956

For mbd
Accuracy - 0.9793538884388797
Precision - 0.977512828481112

For dtc
Accuracy - 0.932675844811854
Precision - 0.935788512861124

For lr
Accuracy - 0.96228821346499
Precision - 0.963861877828181

For rfc
Accuracy - 0.978878817811214
Precision - 0.971338628549118

For adt
Accuracy - 0.981386118117171
Precision - 0.971378541492841

For xgb
Accuracy - 0.978851088872848
Precision - 0.979481188881111

For gbc
Accuracy - 0.951111111111111
Precision - 0.958822858822858

[] df1=pd.DataFrame({'Algorithm':classifiers.keys(),'Precision':precision_scores,
                    'Test Accuracy':accuracy_scores,'Train Accuracy':train_accuracy_score}).round(3)

df2=df1.sort_values(['Precision','Test Accuracy'],ascending=False)
df2

```

```

EmailSMS Spam Classifier.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
Code Text Copy to Drive
"Test Accuracy":accuracy_scores,"Train Accuracy":train_accuracy_score}).round(3)

[] df2=df1.sort_values(['Precision','Test Accuracy'],ascending=False)
df2

Algorithm Precision Test Accuracy Train Accuracy
0 svm 0.981 0.980 0.985
1 mnb 0.978 0.979 0.983
2 mbd 0.975 0.975 0.985
3 xgb 0.973 0.973 1.000
4 dtc 0.971 0.967 0.975
5 lr 0.963 0.962 0.965
6 gbc 0.954 0.955 0.966
7 dtc 0.937 0.933 0.947
8 knc 0.903 0.908 0.927

[] df3 = pd.melt(df2, id_vars = "Algorithm")
df3.head()

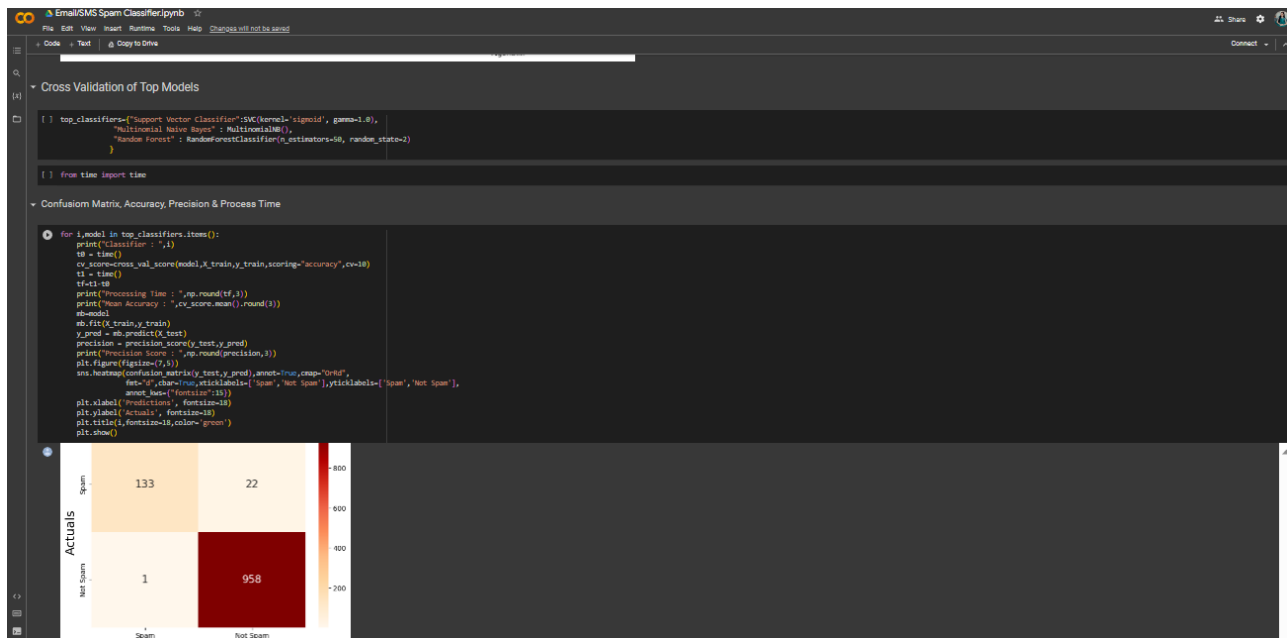
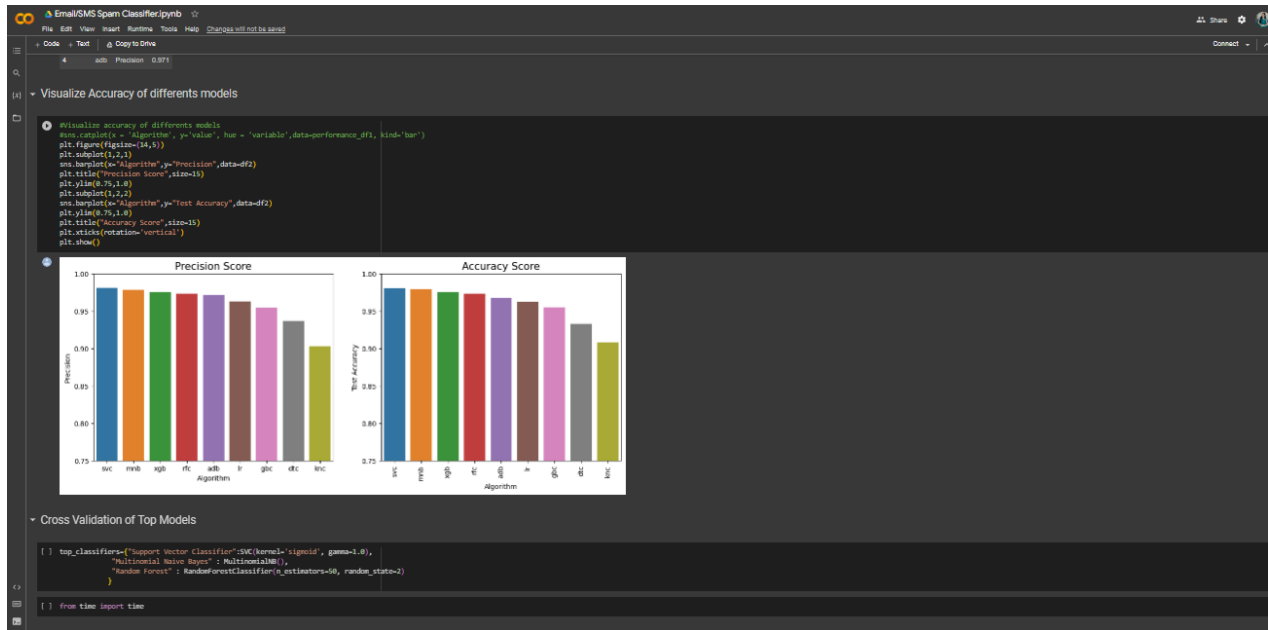
Algorithm variable value
0 svm Precision 0.981
1 mnb Precision 0.978
2 mbd Precision 0.975
3 xgb Precision 0.973
4 dtc Precision 0.971

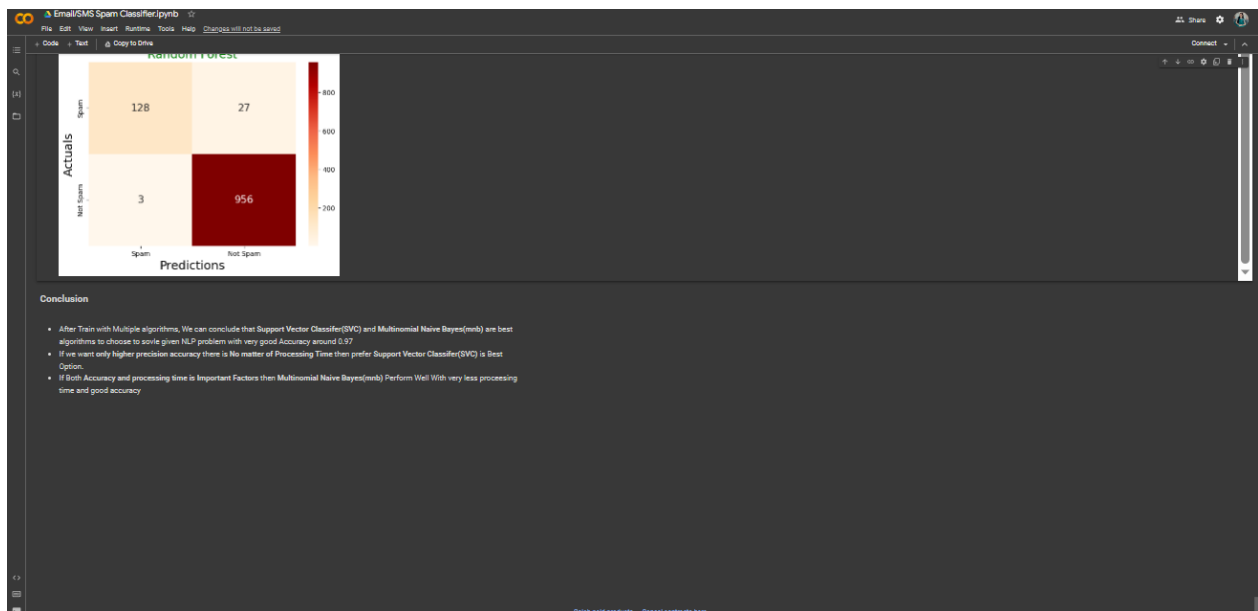
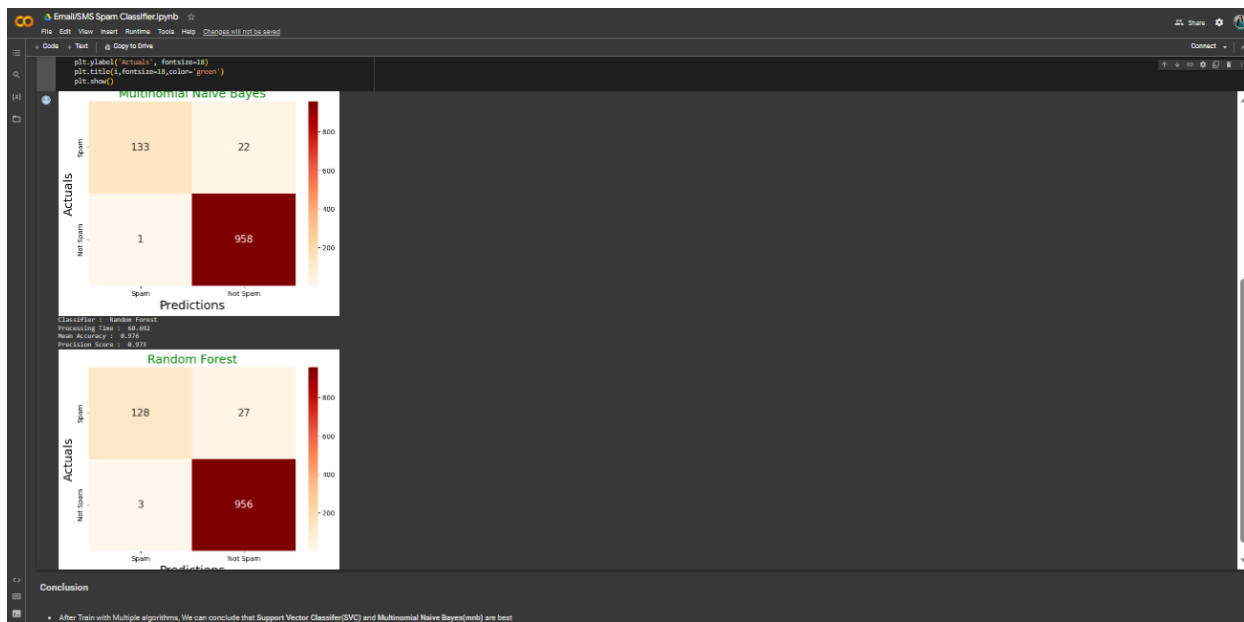
Visualize Accuracy of different models

Visualize accuracy of different models
sns.pairplot(x='Algorithm', y='value', hue = 'variable', data=performance_df3, kind='bar')
plt.figure(figsize=(14,5))
plt.subplot(1,2,1)
sns.barplot(x='Algorithm',y='Precision',data=df2)
plt.title('Precision Score',size=15)
plt.ylim(0.75,1.0)
plt.subplot(1,2,2)
sns.barplot(x='Algorithm',y='Test Accuracy',data=df2)
plt.ylim(0.75,1.0)
plt.title('Accuracy Score',size=15)
plt.xticks(rotation='vertical')
plt.show()

1.00 Precision Score 1.00 Accuracy Score

```





PROJECT CODE:

```
from google.colab import drive  
drive.mount('/content/drive')
```

```
import nltk  
  
import numpy as np  
  
import pandas as pd  
  
import matplotlib.pyplot as plt  
  
import seaborn as sns  
  
import warnings  
warnings.filterwarnings('ignore')
```

#import DataFrame

```
df = pd.read_csv('/content/drive/MyDrive/Datasets/spam.csv',encoding='latin-1')  
df
```

#Remove Unwanted Columns

```
df=df.iloc[:,2]
```

#Rename Columns

```
df=df.rename(columns={"v1":"label","v2":"text"})  
df
```

#Total No of Spam and Not Spam Category

```
df['label'].value_counts()
```

```
import matplotlib.pyplot as plt
```

```
plt.figure(figsize=(9,4))
```

```
plt.subplot(1,2,1)
```

```
plt.pie(df['label'].value_counts(),labels=['Negative','Positive'],autopct="%0.2f")
```

```
plt.subplot(1,2,2)
```

```
sns.barplot(x=df['label'].value_counts().index,y=df['label'].value_counts(),data=df)
```

```
plt.show()
```

#Info

```
df.info()
```

#shape of df

```
df.shape
```

#Check for Null Values

```
df.isnull().sum()
```

```
df['label'].value_counts()/df.shape[0]*100
```

#1. Total No of Char

```
df['num_char']=df['text'].apply(len)
```

#2. Total No of Words

```
df['num_words']=df['text'].apply(lambda x: len(str(x).split()))
```

3. Total No of Sentences

```
nlTK.download('punkt')
```

```
df['num_sen']=df['text'].apply(lambda x: len(nltk.sent_tokenize(x)))
```

#Statistical Info about dataset

```
df.describe()
```

```
df[df['label']=='ham'].describe()
```

```
df[df['label']=='spam'].describe()
```

#Hist Plot for Spam and Not Spam

```
plt.figure(figsize=(12,6))
```

```
sns.histplot(df[df['label']=='spam']['num_words'],color='blue',bins=40)
```

```
sns.histplot(df[df['label']=='ham']['num_words'],color='red')
```

```
plt.show()
```

```
sns.pairplot(df,hue='label')
```

```
plt.show()
```

#Remove few Outliers present in dataset

```
i=df[df['num_char']>500].index
```

```
df.drop(i,axis=0,inplace=True)
```

```
df=df.reset_index()
```

```
df.drop("index",inplace=True,axis=1)
```

```
sns.pairplot(df,hue='label')
```

```
plt.show()
```

#HeatMap

```
sns.heatmap(df.corr(),annot=True)
```

```
plt.show()
```

#Import lib required for text processing

```
from nltk.corpus import stopwords
```

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

```
stopwords.words('english')
```

```
from nltk.stem import PorterStemmer
```

```
from wordcloud import WordCloud
```

```
import string,time
```

```
string.punctuation
```

#punctuation

```
string.punctuation
```

#Stopwards

```
stopwords.words('english')
```

#User Define Funtion for Text processing

```
def remove_website_links(text):
```

```
    no_website_links = text.replace(r"http\S+", "")
```

```
    return no_website_links
```

```
def remove_numbers(text):
```

```
    removed_numbers = text.replace(r'\d+', '')
```

```
    return removed_numbers
```

```
def remove_emails(text):
```

```
    no_emails = text.replace(r"\S*@\S*\s?", "")
```

```
    return no_emails
```

#Call Function

```
df['text'] = df['text'].apply(remove_website_links)
```

```
df['text'] = df['text'].apply(remove_numbers)
```

```
df['text'] = df['text'].apply(remove_emails)
```

#Create Common Function


```
def transform_text(text):  
    #1.lower casing  
    text=text.lower()  
  
    #2.tokenization  
    lst=nlk.word_tokenize(text)  
  
    #3.remove spcl characters stopwords and punctuation  
    l1=[]  
    useless_words=stopwords.words('english')+list(string.punctuation)  
    for word in lst:  
        if word.isalnum()==True and word not in useless_words:  
            l1.append(word)  
  
    #4.stemming  
    l2=[]  
    for word in l1:  
        ps=PorterStemmer()  
        l2.append(ps.stem(word))  
  
    return " ".join(l2).strip()  
    l1.clear()  
    l2.clear()
```

#call function

```
nlTK.download('punkt')

df['text'] = df['text'].apply(transform_text)

df['num_words_transform']=df['text'].apply(lambda x: len(str(x).split()))

#conda install -c conda-forge wordcloud

from wordcloud import WordCloud

plt.figure(figsize = (20,20)) # Text that is not fraudulent(0)

wc = WordCloud(width = 1600 , height = 800 , max_words = 500,background_color='black').generate(" ".join(df[df.label == 'spam'].text))

plt.imshow(wc , interpolation = 'bilinear')

plt.figure(figsize = (20,20)) # Text that is not fraudulent(0)

wc = WordCloud(width = 1600 , height = 800 , max_words = 500).generate(" ".join(df[df.label == 'ham'].text))

plt.imshow(wc , interpolation = 'bilinear')


df.head()

df['label']=df['label'].replace({'spam':0,'ham':1})

#create spam corpus which will holds all Spam words

spam_corpus = []

for msg in df[df['label'] == 1]['text'].tolist():

    for word in msg.split():

        spam_corpus.append(word)

from collections import Counter
```

```
plt.figure(figsize=(10,7))

sns.barplot(y=pd.DataFrame(Counter(spam_corpus).most_common(25))[0],x=pd.
DataFrame(Counter(spam_corpus).most_common(30))[1])

plt.xticks()

plt.title("Most Commonly Used Non Spam Words")

plt.xlabel("Frequency")

plt.ylabel("Ham Words")

plt.show()
```

#create spam corpus which will holds all Ham or Non Spam words

```
ham_corpus = []
```

```
for msg in df[df['label'] == 0]['text'].tolist():
```

```
    for word in msg.split():
```

```
        ham_corpus.append(word)
```

```
from collections import Counter
```

```
plt.figure(figsize=(10,7))
```

```
sns.barplot(y=pd.DataFrame(Counter(ham_corpus).most_common(25))[0],x=pd.D
ataFrame(Counter(ham_corpus).most_common(30))[1])
```

```
plt.xticks()
```

```
plt.title("Most Commonly Used Spam Words")
```

```
plt.xlabel("Frequency")
```

```
plt.ylabel("Spam Words")
```

```
plt.show()
```

#Characters Visualize

```
fig,(ax1,ax2)=plt.subplots(1,2,figsize=(15,6))
```

```
text_len=df[df['label']==1]['text'].str.len()
```

```
ax1.hist(text_len,color='green')
```

```
ax1.set_title('Original text')
```

```
text_len=df[df['label']==0]['text'].str.len()
```

```
ax2.hist(text_len,color='red')
```

```
ax2.set_title('Fake text')
```

```
fig.suptitle('Characters in texts')
```

```
plt.show()
```

```
#Words Visualize
```

```
fig,(ax1,ax2)=plt.subplots(1,2,figsize=(15,6))
```

```
text_len=df[df['label']==1]['num_words_transform']
```

```
ax1.hist(text_len,color='red')
```

```
ax1.set_title('Original text')
```

```
text_len=df[df['label']==0]['num_words_transform']
```

```
ax2.hist(text_len,color='green')
```

```
ax2.set_title('Fake text')
```

```
fig.suptitle('Words in texts')
```

```
plt.show()
```

```
fig,(ax1,ax2)=plt.subplots(1,2,figsize=(15,6))
```

```
word=df[df['label']==1]['num_words_transform']
```

```
sns.distplot(word.map(lambda x: np.mean(x)),ax=ax1,color='red')
```

```
ax1.set_title('Original text')
```

```

word=df[df['label']==0]['num_words_transform']

sns.distplot(word.map(lambda x: np.mean(x)),ax=ax2,color='green')

ax2.set_title('Fake text')

fig.suptitle('Average word length in each text')

#Text Vecorization

from sklearn.feature_extraction.text import CountVectorizer,TfidfVectorizer

cv = CountVectorizer()

tfidf = TfidfVectorizer(max_features=3000)

#Input and Output Features

X = tfidf.fit_transform(df['text']).toarray()

y = df['label'].values

X.shape

y.shape

from sklearn.model_selection import train_test_split

#Model Training

from sklearn.naive_bayes import GaussianNB,MultinomialNB,BernoulliNB

from sklearn.metrics import accuracy_score,confusion_matrix,precision_score

!pip install --upgrade scikit-learn

#Imports Lib

from sklearn.linear_model import LogisticRegression

from sklearn.svm import SVC

from sklearn.naive_bayes import MultinomialNB

from sklearn.tree import DecisionTreeClassifier

```

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.ensemble import ExtraTreesClassifier
from sklearn.ensemble import GradientBoostingClassifier
from xgboost import XGBClassifier

from sklearn.metrics import precision_score, recall_score, classification_report,
accuracy_score, f1_score

#from sklearn.metrics import plot_confusion_matrix

from sklearn.model_selection import cross_val_score

#from sklearn.metrics import precision_score, recall_score,
plot_confusion_matrix, classification_report, accuracy_score, f1_score

from sklearn.model_selection import cross_val_score

#GaussianNB

gnb = GaussianNB()

gnb.fit(X_train,y_train)

y_pred1 = gnb.predict(X_test)

print("Accuracy Score -",accuracy_score(y_test,y_pred1))

print("Precision Score -",precision_score(y_test,y_pred1))

print(confusion_matrix(y_test,y_pred1))

#MultinomialNB

mnb = MultinomialNB()
```

```
mnb.fit(X_train,y_train)

y_pred2 = mnb.predict(X_test)

print("Accuracy Score -",accuracy_score(y_test,y_pred2))

print("Precision Score -",precision_score(y_test,y_pred2))

print(confusion_matrix(y_test,y_pred2))
```

```
classifiers={"svc":SVC(kernel='sigmoid', gamma=1.0),

            "knc": KNeighborsClassifier(),

            "mnb" : MultinomialNB(),

            "dtc" : DecisionTreeClassifier(max_depth=5),

            "lr" : LogisticRegression(solver='liblinear', penalty='l1'),

            "rfc" : RandomForestClassifier(n_estimators=50, random_state=2),

            "adb" : AdaBoostClassifier(n_estimators=50, random_state=2),

            "xgb" : XGBClassifier(n_estimators=50,random_state=2),

            "gbc" : GradientBoostingClassifier(n_estimators=50,random_state=2)

            }
```

#Common Function for model train

```
def train_classifier(clf,X_train,y_train,X_test,y_test):

    clf.fit(X_train,y_train)

    y_pred = clf.predict(X_test)

    accuracy = accuracy_score(y_test,y_pred)

    precision = precision_score(y_test,y_pred)
```

```

train_accuracy = clf.score(X_train,y_train)

return accuracy,precision,train_accuracy
svc=SVC(kernel='sigmoid', gamma=1.0)
train_classifier(svc,X_train,y_train,X_test,y_test)
accuracy_scores = []
precision_scores = []
train_accuracy_score=[]

for name,clf in classifiers.items():

    current_accuracy,current_precision,current_train_score = train_classifier(clf,
X_train,y_train,X_test,y_test)

    print("For ",name)
    print("Accuracy - ",current_accuracy)
    print("Precision - ",current_precision)

    accuracy_scores.append(current_accuracy)
    precision_scores.append(current_precision)
    train_accuracy_score.append(current_train_score)
    print()

df1=pd.DataFrame({'Algorithm':classifiers.keys(),'Precision':precision_scores,

```



```
        'Test Accuracy':accuracy_scores,'Train Accuracy':train_accuracy_score}).round(3)

df2=df1.sort_values(['Precision','Test Accuracy'],ascending=False)

df2

df3 = pd.melt(df2, id_vars = "Algorithm")

df3.head()
```

#Visualize accuracy of different models

```
#sns.catplot(x = 'Algorithm', y='value', hue = 'variable',data=performance_df1,
kind='bar')

plt.figure(figsize=(14,5))

plt.subplot(1,2,1)

sns.barplot(x="Algorithm",y="Precision",data=df2)

plt.title("Precision Score",size=15)

plt.ylim(0.75,1.0)

plt.subplot(1,2,2)

sns.barplot(x="Algorithm",y="Test Accuracy",data=df2)

plt.ylim(0.75,1.0)

plt.title("Accuracy Score",size=15)

plt.xticks(rotation='vertical')

plt.show()

top_classifiers={"Support Vector Classifier":SVC(kernel='sigmoid', gamma=1.0),
```

```

        "Multinomial Naive Bayes" : MultinomialNB(),

        "Random Forest" : RandomForestClassifier(n_estimators=50,
random_state=2)

    }

from time import time

for i,model in top_classifiers.items():

    print("Classifier : ",i)

    t0 = time()

    cv_score=cross_val_score(model,X_train,y_train,scoring="accuracy",cv=10)

    t1 = time()

    tf=t1-t0

    print("Processing Time : ",np.round(tf,3))

    print("Mean Accuracy : ",cv_score.mean().round(3))

    mb=model

    mb.fit(X_train,y_train)

    y_pred = mb.predict(X_test)

    precision = precision_score(y_test,y_pred)

    print("Precision Score : ",np.round(precision,3))

    plt.figure(figsize=(7,5))

    sns.heatmap(confusion_matrix(y_test,y_pred),annot=True,cmap="OrRd",

                fmt="d",cbar=True,xticklabels=['Spam','Not
Spam'],yticklabels=['Spam','Not Spam'],

```

```
        annot_kws={"fontsize":15})  
plt.xlabel('Predictions', fontsize=18)  
plt.ylabel('Actuals', fontsize=18)  
plt.title(i,fontsize=18,color='green')  
plt.show()
```

TEAM MEMBERS DETAILS:

TEAM NAME: Shadow knights

MEMBER 1:

NAME : HEMA HARIHARAN S
DEPT : CSE
Email : hemaharharansamson@gmail.com
Mobile Number : 9080602796

MEMBER 2:

NAME : DHATCHAYANI U
DEPT : CSE
Email : dhatchayani8102@gmail.com
Mobile Number : 9385448919