

```

1 ! pip install nltk
2 ! pip install keras
3 ! pip install lifelines
4

```

```

Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: nltk in /usr/local/lib/python3.9/dist-packages (3.8.1)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.9/dist-packages (from nltk) (2022.10.31)
Requirement already satisfied: joblib in /usr/local/lib/python3.9/dist-packages (from nltk) (1.1.1)
Requirement already satisfied: tqdm in /usr/local/lib/python3.9/dist-packages (from nltk) (4.65.0)
Requirement already satisfied: click in /usr/local/lib/python3.9/dist-packages (from nltk) (8.1.3)
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: keras in /usr/local/lib/python3.9/dist-packages (2.12.0)
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: lifelines in /usr/local/lib/python3.9/dist-packages (0.27.4)
Requirement already satisfied: pandas>=1.0.0 in /usr/local/lib/python3.9/dist-packages (from lifelines) (1.4.4)
Requirement already satisfied: autograd>=1.5 in /usr/local/lib/python3.9/dist-packages (from lifelines) (1.5)
Requirement already satisfied: numpy>=1.14.0 in /usr/local/lib/python3.9/dist-packages (from lifelines) (1.22.4)
Requirement already satisfied: scipy>=1.2.0 in /usr/local/lib/python3.9/dist-packages (from lifelines) (1.10.1)
Requirement already satisfied: matplotlib>=3.0 in /usr/local/lib/python3.9/dist-packages (from lifelines) (3.7.1)
Requirement already satisfied: autograd-gamma>=0.3 in /usr/local/lib/python3.9/dist-packages (from lifelines) (0.5.0)
Requirement already satisfied: formulaic>=0.2.2 in /usr/local/lib/python3.9/dist-packages (from lifelines) (0.5.2)
Requirement already satisfied: future>=0.15.2 in /usr/local/lib/python3.9/dist-packages (from autograd>=1.5->lifelines) (0.18.2)
Requirement already satisfied: typing-extensions>=4.2.0 in /usr/local/lib/python3.9/dist-packages (from formulaic>=0.2.2->lifelines) (4.5.0)
Requirement already satisfied: astor>=0.8 in /usr/local/lib/python3.9/dist-packages (from formulaic>=0.2.2->lifelines) (0.9.1)
Requirement already satisfied: interface-meta>=1.2.0 in /usr/local/lib/python3.9/dist-packages (from formulaic>=0.2.2->lifelines) (1.3.0)
Requirement already satisfied: wrapt>=1.0 in /usr/local/lib/python3.9/dist-packages (from formulaic>=0.2.2->lifelines) (1.14.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (1.0.7)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (4.22.0)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (21.3)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (1.4.4)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (3.0.9)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (9.0.1)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (2.8.2)
Requirement already satisfied: importlib-resources>=3.2.0 in /usr/local/lib/python3.9/dist-packages (from matplotlib>=3.0->lifelines) (5.10.0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.9/dist-packages (from pandas>=1.0.0->lifelines) (2021.3)
Requirement already satisfied: zipp>=3.1.0 in /usr/local/lib/python3.9/dist-packages (from importlib-resources>=3.2.0->lifelines) (3.10.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.9/dist-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)

```

```

1 !pip install -q wordcloud
2 import wordcloud
3 import nltk
4 nltk.download('stopwords')
5 nltk.download('wordnet')
6 nltk.download('punkt')
7 nltk.download('averaged_perceptron_tagger')

```

```

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

```

```

1 #Importing the libraries
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
7 from sklearn.naive_bayes import MultinomialNB
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
10 from sklearn.svm import SVC
11 from sklearn import tree
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.ensemble import RandomForestClassifier
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.neighbors import KNeighborsClassifier
16 from sklearn.neural_network import MLPClassifier
17 from sklearn.metrics import accuracy_score
18 from sklearn.ensemble import GradientBoostingClassifier
19 from sklearn.naive_bayes import GaussianNB
20 from scipy.sparse import csr_matrix
21 from sklearn.metrics import roc_curve, auc
22 from keras.models import Sequential
23 from keras.layers import Dense, Conv1D, MaxPooling1D, Flatten
24 import numpy as np

```

```

25 from sklearn.metrics import precision_recall_curve
26 import matplotlib.pyplot as plt
27 # Importing the warnings
28 import warnings
29 warnings.filterwarnings('ignore')
30 from sklearn.datasets import make_classification
31 from sklearn.model_selection import train_test_split
32 from sklearn.linear_model import LogisticRegression
33 import numpy as np
34 import matplotlib.pyplot as plt
35 from sklearn.metrics import precision_recall_curve, average_precision_score
36

```

```

1 #Loading the dataset
2 df = pd.read_csv("messages.csv",encoding='latin-1')

```

```
1 df.head()
```

	subject	message	label
0	job posting - apple-iss research center	content - length : 3386 apple-iss research cen...	0
1	NaN	lang classification grimes , joseph e . and ba...	0
2	query : letter frequencies for text identifica...	i am posting this inquiry for sergei atamas (...	0
3	risk	a colleague and i are researching the differin...	0
4	request book information	earlier this morning i was on the phone with a...	0

```

1 #Checking information of dataset
2 df.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2893 entries, 0 to 2892
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0    subject    2831 non-null   object
1    message    2893 non-null   object
2    label      2893 non-null   int64
dtypes: int64(1), object(2)
memory usage: 67.9+ KB

```

```

1 #Checking the shape of the dataset
2 print("Shape of the dataset:", df.shape)

```

```
Shape of the dataset: (2893, 3)
```

```

1 #Checking for the null values
2 df.isnull().values.any()

```

```
True
```

```

1 #Checkin for the null values in columns
2 df.isnull().sum()

```

```

subject      62
message       0
label        0
dtype: int64

```

```

1 # replace null values with empty strings
2 df["message"] = df["message"].replace(np.nan, "", regex=True)
3
4 # remove messages with empty strings
5 df = df[df["message"].str.strip().astype(bool)]
6
7 # update csv file
8 df.to_csv("messages_cleaned.csv", index=False)
9

```

```

1 # 62 row are missing in the subject columns that means 62 emails are without subject heading.
2 # Here, not dropping Nan rows for subject column as it of no use in building model.

```

```

1 #Checking total number of mails
2 print("Count of label:\n",df['label'].value_counts())

```

```
Count of label:
0      2412
1       481
Name: label, dtype: int64
```

```
1 # Note:- Here in our dataset 1 stands for Spam mail and 0 stands for not a spam mail.
2 #Checking the Ratio of labels
3 print("Not a Spam Email Ratio i.e. 0 label:",round(len(df[df['label']==0])/len(df['label']),2)*100,"%")
4 print("Spam Email Ratio that is 1 label:",round(len(df[df['label']==1])/len(df['label']),2)*100,"%")
```

```
Not a Spam Email Ratio i.e. 0 label: 83.0 %
Spam Email Ratio that is 1 label: 17.0 %
```

```
1 # so here 17 % of the data is a spam email
```

```
1 #Creating the new column for length of message column
2 df['length'] = df.message.str.len()
3 df.head()
```

	subject	message	label	length
0	job posting - apple-iss research center	content - length : 3386 apple-iss research cen...	0	2856
1	NaN	lang classification grimes , joseph e . and ba...	0	1800
2	query : letter frequencies for text identifica...	i am posting this inquiry for sergei atamas (...	0	1435
3	risk	a colleague and i are researching the differin...	0	324
4	request book information	earlier this morning i was on the phone with a...	0	1046

```
1 #Converting all messages to lower case
2 df['message'] = df['message'].str.lower()
3 df.head()
```

	subject	message	label	length
0	job posting - apple-iss research center	content - length : 3386 apple-iss research cen...	0	2856
1	NaN	lang classification grimes , joseph e . and ba...	0	1800
2	query : letter frequencies for text identifica...	i am posting this inquiry for sergei atamas (...	0	1435
3	risk	a colleague and i are researching the differin...	0	324
4	request book information	earlier this morning i was on the phone with a...	0	1046

```
1 # regular expressions
2 # Replace email addresses with 'email'
3 df['message'] = df['message'].str.replace(r'^.+@[^\.\.]*\.[a-z]{2,}$','emailaddress')
4
5 # Replace URLs with 'webaddress'
6 df['message'] = df['message'].str.replace(r'^http://[a-zA-Z0-9\-\.\.]+\.[a-zA-Z]{2,3}(/\S*)?$','webaddress')
7
8 # Replace currency symbols with 'moneysymb' (£ can by typed with ALT key + 156)
9 df['message'] = df['message'].str.replace(r'£|$', 'dollars')
10
11 # Replace 10 digit phone numbers (formats include paranthesis, spaces, no spaces, dashes) with 'phonenumber'
12 df['message'] = df['message'].str.replace(r'^((?[\d]{3})?[\s-]?[\d]{3}[\s-]?[\d]{4})$', 'phonenumber')
13
14 # Replace numeric characters with 'numbr'
15 df['message'] = df['message'].str.replace(r'\d+(\.\d+)?', 'numbr')
```

```
1 # Remove punctuation
2 df['message'] = df['message'].str.replace(r'[^\w\d\s]', ' ')
3
4 # Replace whitespace between terms with a single space
5 df['message'] = df['message'].str.replace(r'\s+', ' ')
6
7 # Remove leading and trailing whitespace
8 df['message'] = df['message'].str.replace(r'^\s+|\s+$', '')
```

```
1 # now re-checking the data
2 df.head()
```

	subject	message	label	length
0	job posting - apple-iss research center	content length numbr apple iss research center...	0	2856
1	NaN	lang classification grimes joseph e and barbar...	0	1800

```

1 #Removing the stopwords
2 import string
3 import nltk
4 from nltk.corpus import stopwords
5
6 stop_words = set(stopwords.words('english') + ['u', 'ü', 'ur', '4', '2', 'im', 'dont', 'doin', 'ure'])
7
8 df['message'] = df['message'].apply(lambda x: " ".join(term for term in x.split() if term not in stop_words))

```

```

1 # New column (clean_length) after punctuations,stopwords removal
2 df['clean_length'] = df.message.str.len()
3 df.head()

```

	subject	message	label	length	clean_length
0	job posting - apple-iss research center	content length numbr apple iss research center...	0	2856	2179
1	NaN	lang classification grimes joseph e barbara f...	0	1800	1454
2	query : letter frequencies for text identifica...	posting inquiry sergei atamas satamas umabnet ...	0	1435	1064
3	risk	colleague researching differing degrees risk p...	0	324	210
4	request book information	earlier morning phone friend mine living south...	0	1046	629

```

1 #Total length removal
2 print("Original Length:",df.length.sum())
3 print("Cleaned Length:",df.clean_length.sum())
4 print("Total Words Removed:",(df.length.sum()) - (df.clean_length.sum()))

```

```

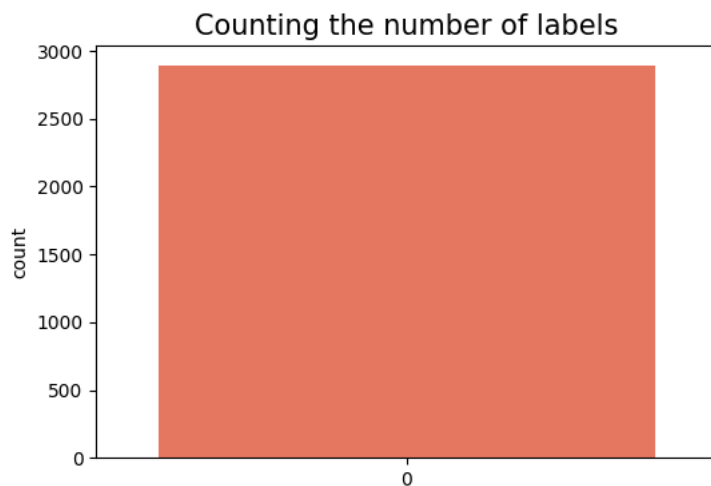
Original Length: 9344743
Cleaned Length: 6767857
Total Words Removed: 2576886

```

```

1 #Graphical Visualisation for counting number of labels.
2 plt.figure(figsize=(6,4))
3 sns.countplot(df['label'],palette= 'Reds')
4 plt.title("Counting the number of labels",fontsize=15)
5 plt.xticks(rotation='horizontal')
6 plt.show()
7
8 print(df.label.value_counts())

```



```

0    2412
1     481
Name: label, dtype: int64

```

```

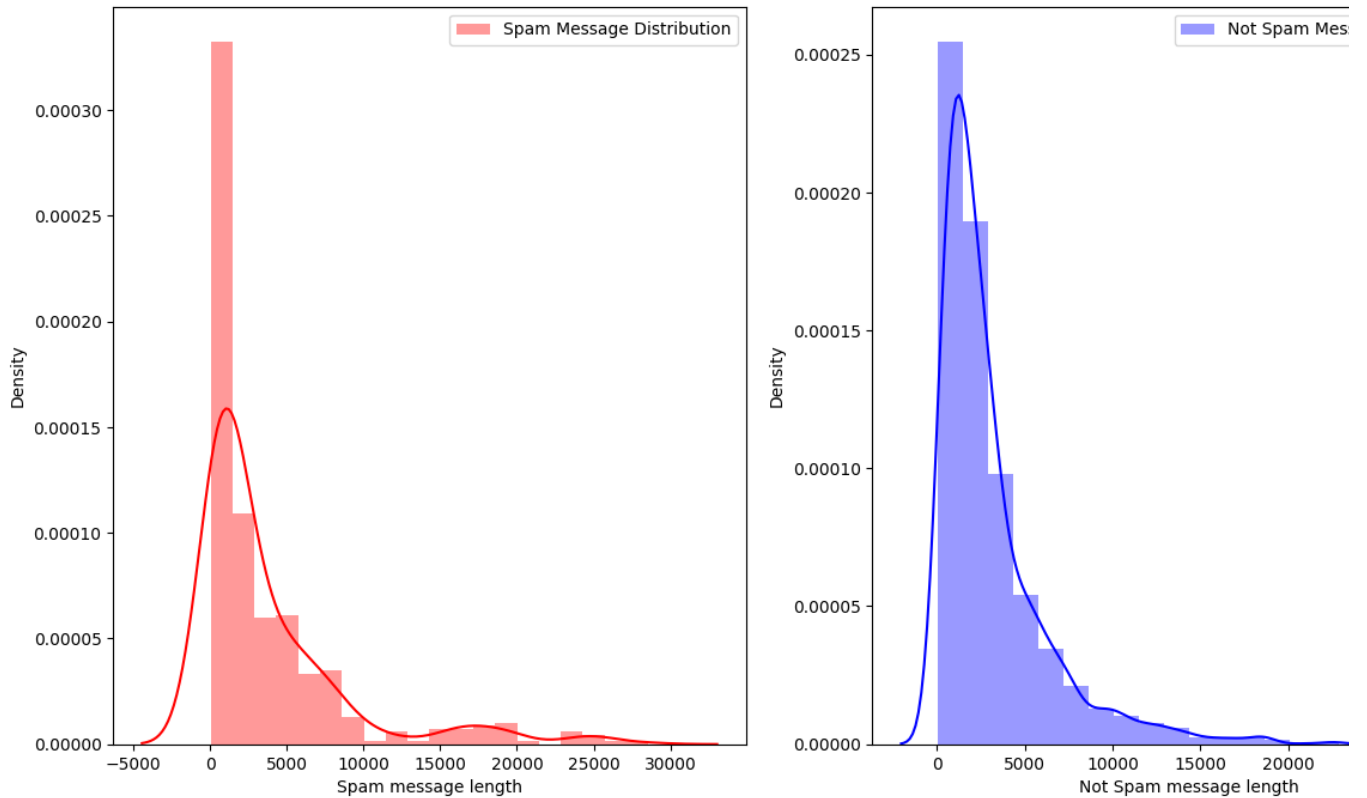
1 #Message distribution before cleaning
2 f,ax = plt.subplots(1,2,figsize=(15,8))
3
4 sns.distplot(df[df['label']==1]['length'],bins=20, ax=ax[0],label='Spam Message Distribution',color='r')
5 ax[0].set_xlabel('Spam message length')
6 ax[0].legend()
7
8 sns.distplot(df[df['label']==0]['length'],bins=20, ax=ax[1],label='Not Spam Message Distribution',color='b')

```

```

9 ax[1].set_xlabel('Not Spam message length')
10 ax[1].legend()
11
12 plt.show()

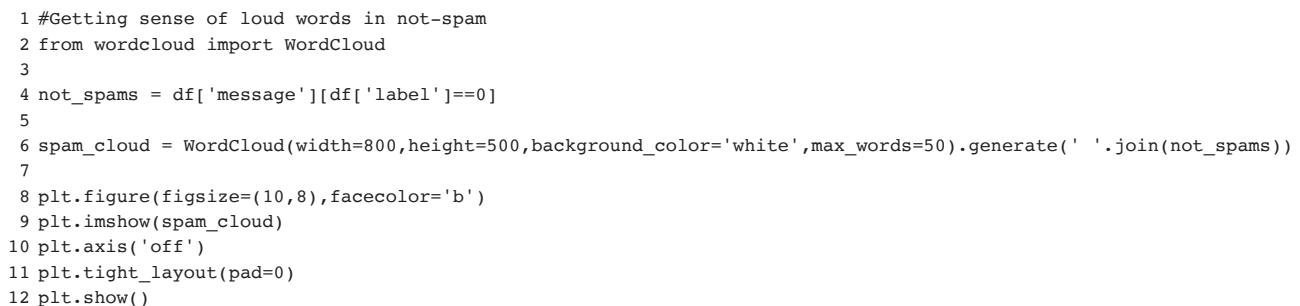
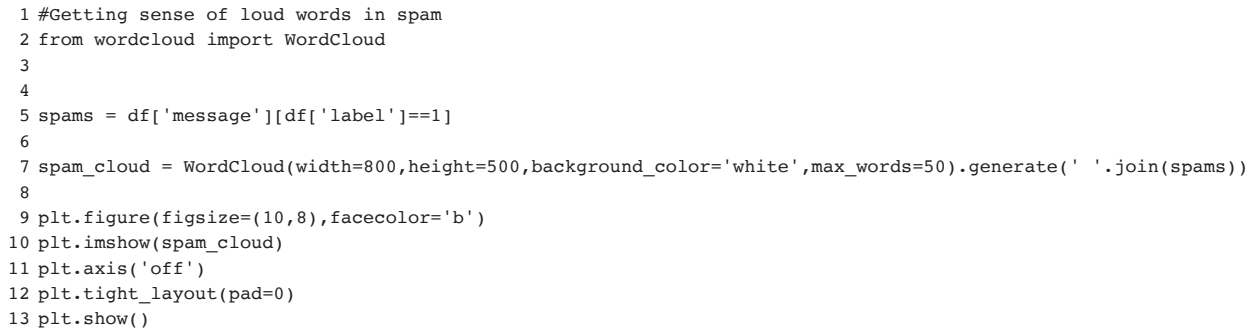
```

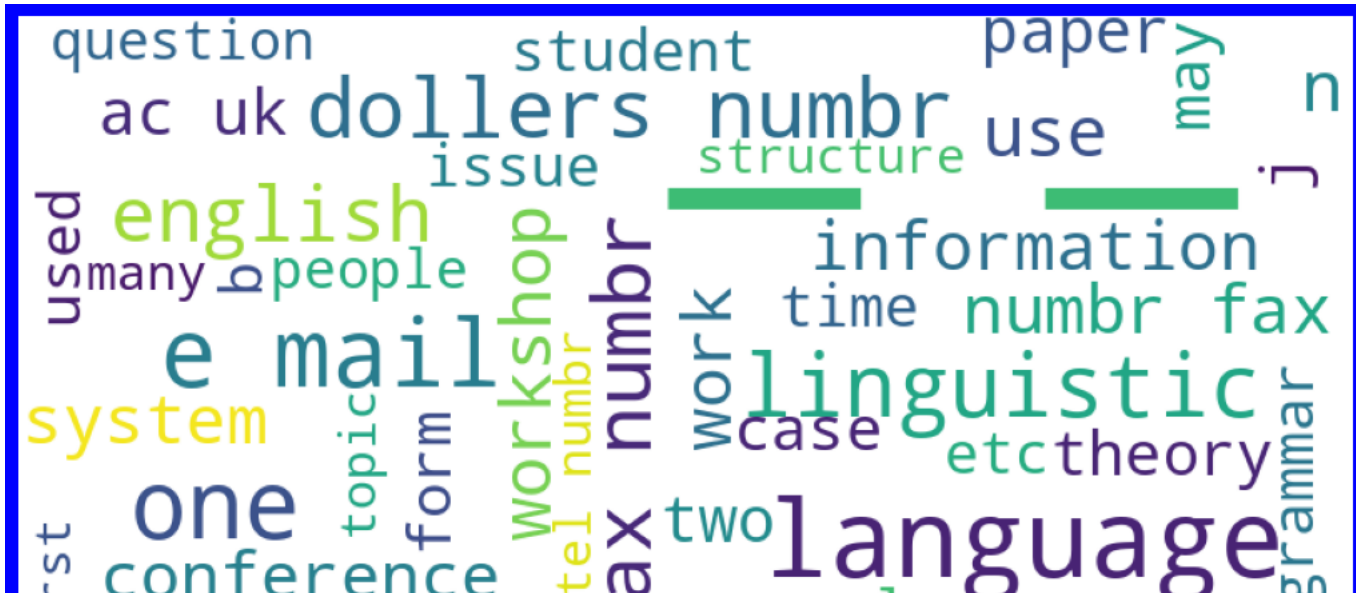


```

1 #Message distribution after cleaning
2 f,ax = plt.subplots(1,2,figsize=(15,8))
3
4 sns.distplot(df[df['label']==1]['clean_length'],bins=20, ax=ax[0],label='Spam Message Distribution',color='r')
5 ax[0].set_xlabel('Spam message length')
6 ax[0].legend()
7
8 sns.distplot(df[df['label']==0]['clean_length'],bins=20, ax=ax[1],label='Not Spam Message Distribution',color='g')
9 ax[1].set_xlabel('Not a Spam message length')
10 ax[1].legend()
11
12 plt.show()

```





```

1 # Converting the text into vectors using TF-IDF, as text cannot be the input in the model
2 # 1. Convert text into vectors using TF-IDF
3 # 2. Instantiate MultinomialNB classifier
4 # 3. Split feature and label
5
6
7 tf_vec = TfidfVectorizer()
8
9 naive = MultinomialNB()
10
11 SVM = SVC(C=1.0, kernel='linear', degree=3 , gamma='auto')
12
13 decision = DecisionTreeClassifier()
14 classifier= RandomForestClassifier(n_estimators= 10, criterion="entropy")
15
16 clf = LogisticRegression()
17
18
19 features = tf_vec.fit_transform(df['message'])
20
21 X = features
22 y = df['label']

```

```

1 # Train and predict for naive bayes model
2 X_train,x_test,Y_train,y_test = train_test_split(X,y,random_state=42)
3
4 #test_size=0.20 random_state=42 test_size=0.15
5
6 naive.fit(X_train,Y_train)
7 y_pred= naive.predict(x_test)
8
9
10
11 print ('Final score = > ', accuracy_score(y_test,y_pred))
12
13
14

```

Final score = > 0.8342541436464088

```

1 # Train and predict for SVM model
2 X_train,x_test,Y_train,y_test = train_test_split(X,y,random_state=42)
3
4 #test_size=0.20 random_state=42 test_size=0.15
5
6
7
8 SVM.fit(X_train,Y_train)
9 y_pred = SVM.predict(x_test)
10
11 print ('Final score = > ', accuracy_score(y_test,y_pred))
12

```

Final score = > 0.9875690607734806

```

1 # train and predict for the Decision tree model
2 X_train,x_test,Y_train,y_test = train_test_split(X,y,random_state=42)
3 decision.fit(X_train,Y_train)
4 #test_size=0.20 random_state=42 test_size=0.15
5
6 y_pred = decision.predict(x_test)
7 print ('Final score = > ', accuracy_score(y_test,y_pred))
8
9

```

Final score = > 0.9571823204419889

```

1 # train and predict using random forest classifier
2 X_train,x_test,Y_train,y_test = train_test_split(X,y,random_state=42)
3 classifier.fit(X_train, Y_train)
4 #test_size=0.20 random_state=42 test_size=0.15
5 y_pred= classifier.predict(x_test)
6 print ('Final score = > ', accuracy_score(y_test,y_pred))
7

```

Final score = > 0.9585635359116023

```

1 # train and predict using logistic regression
2 X_train,x_test,Y_train,y_test = train_test_split(X,y,random_state=42)
3 clf.fit(X_train, Y_train)
4 #test_size=0.20 random_state=42 test_size=0.15
5 y_pred= clf.predict(x_test)
6 print ('Final score = > ', accuracy_score(y_test,y_pred))
7

```

Final score = > 0.9475138121546961

```

1 # train and predict using KNN
2 # Split data into training and testing sets
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
4
5 # Create KNN model and fit it to the training data
6 knn = KNeighborsClassifier(n_neighbors=5)
7 knn.fit(X_train, y_train)
8
9 # Make predictions on the testing data
10 y_pred = knn.predict(X_test)
11
12 # Measure accuracy of the model
13 accuracy = accuracy_score(y_test, y_pred)
14 print("Accuracy:", accuracy)
15

```

Accuracy: 0.9758203799654577

```

1 # train and predict using Neural Network
2 # Split data into training and testing sets
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
4
5 # Create a neural network model and fit it to the training data
6 nn = MLPClassifier(hidden_layer_sizes=(10,), max_iter=1000)
7 nn.fit(X_train, y_train)
8
9 # Make predictions on the testing data
10 y_pred = nn.predict(X_test)
11
12 # Measure accuracy of the model
13 accuracy = accuracy_score(y_test, y_pred)
14 print("Accuracy:", accuracy)
15

```

Accuracy: 0.9930915371329879

```

1 # train and predict using Gradient Boosting Algorithm
2 # Split data into training and testing sets
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
4
5 # Create gradient boosting model and fit it to the training data
6 gb = GradientBoostingClassifier()
7 gb.fit(X_train, y_train)
8
9 # Make predictions on the testing data
10 y_pred = gb.predict(X_test)
11
12 # Measure accuracy of the model

```



```

13 accuracy = accuracy_score(y_test, y_pred)
14 print("Accuracy:", accuracy)
15

```

Accuracy: 0.9740932642487047

```

1 # train and predict using Gaussian Naive Bayes model
2 # Convert sparse matrix X to dense numpy array
3 X_dense = X.toarray()
4
5 # Split data into training and testing sets
6 X_train, X_test, y_train, y_test = train_test_split(X_dense, y, test_size=0.2, random_state=42)
7
8 # Create Gaussian Naive Bayes model and fit it to the training data
9 nb = GaussianNB()
10 nb.fit(X_train, y_train)
11
12 # Make predictions on the testing data
13 y_pred = nb.predict(X_test)
14
15 # Measure accuracy of the model
16 accuracy = accuracy_score(y_test, y_pred)
17 print("Accuracy:", accuracy)
18

```

Accuracy: 0.9430051813471503

```

1 # train and predict using CNN
2
3 # Convert sparse matrix X to dense numpy array
4 X_dense = X.toarray()
5
6 # Reshape X_dense for use in 1D Convolutional Neural Network
7 X_resaped = np.expand_dims(X_dense, axis=2)
8
9 # Split data into training and testing sets
10 X_train, X_test, y_train, y_test = train_test_split(X_resaped, y, test_size=0.2, random_state=42)
11
12 # Define the CNN model
13 model = Sequential()
14 model.add(Conv1D(filters=32, kernel_size=3, activation='relu', input_shape=(X_train.shape[1], X_train.shape[2])))
15 model.add(MaxPooling1D(pool_size=2))
16 model.add(Flatten())
17 model.add(Dense(1, activation='sigmoid'))
18
19 # Compile the model
20 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
21
22 # Fit the model to the training data
23 model.fit(X_train, y_train, epochs=10, batch_size=32)
24
25 # Make predictions on the testing data
26 y_pred_prob = model.predict(X_test)
27 y_pred = np.argmax(y_pred_prob, axis=1)
28
29 # Measure accuracy of the model
30 accuracy = accuracy_score(y_test, y_pred)
31 print("Accuracy:", accuracy)
32

```

```

Epoch 1/10
73/73 [=====] - 84s 1s/step - loss: 0.4445 - accuracy: 0.8392
Epoch 2/10
73/73 [=====] - 95s 1s/step - loss: 0.2247 - accuracy: 0.9097
Epoch 3/10
73/73 [=====] - 82s 1s/step - loss: 0.0910 - accuracy: 0.9831
Epoch 4/10
73/73 [=====] - 84s 1s/step - loss: 0.0458 - accuracy: 0.9948
Epoch 5/10
73/73 [=====] - 85s 1s/step - loss: 0.0259 - accuracy: 0.9991
Epoch 6/10
73/73 [=====] - 82s 1s/step - loss: 0.0153 - accuracy: 0.9991
Epoch 7/10
73/73 [=====] - 84s 1s/step - loss: 0.0100 - accuracy: 0.9996
Epoch 8/10
73/73 [=====] - 92s 1s/step - loss: 0.0065 - accuracy: 1.0000
Epoch 9/10
73/73 [=====] - 102s 1s/step - loss: 0.0045 - accuracy: 1.0000
Epoch 10/10
73/73 [=====] - 87s 1s/step - loss: 0.0033 - accuracy: 1.0000
19/19 [=====] - 7s 341ms/step
Accuracy: 0.8013816925734024

```

[illegible]

```
1 tree.plot_tree(decision)
```

```

[Text(0.8589449541284404, 0.9821428571428571, 'x[18089] <= 0.015\ngini = 0.266\nsamples = 2169\nvalue = [1827, 342]'),
Text(0.7766055045871559, 0.9464285714285714, 'x[41154] <= 0.009\ngini = 0.157\nsamples = 1901\nvalue = [1738, 163]'),
Text(0.7146788990825688, 0.9107142857142857, 'x[32052] <= 0.076\ngini = 0.12\nsamples = 1853\nvalue = [1734, 119]'),
Text(0.6642201834862386, 0.875, 'x[8614] <= 0.053\ngini = 0.1\nsamples = 1828\nvalue = [1732, 96]'),
Text(0.6073394495412844, 0.8392857142857143, 'x[19527] <= 0.068\ngini = 0.08\nsamples = 1804\nvalue = [1729, 75]'),
Text(0.5669724770642202, 0.8035714285714286, 'x[23195] <= 0.019\ngini = 0.065\nsamples = 1780\nvalue = [1720, 60]'),
Text(0.5302752293577981, 0.7678571428571429, 'x[49044] <= 0.029\ngini = 0.057\nsamples = 1771\nvalue = [1719, 52]'),
Text(0.5009174311926605, 0.7321428571428571, 'x[13671] <= 0.05\ngini = 0.049\nsamples = 1760\nvalue = [1716, 44]'),
Text(0.48623853211009177, 0.6964285714285714, 'x[7268] <= 0.044\ngini = 0.043\nsamples = 1755\nvalue = [1716, 39]'),
Text(0.45688073394495415, 0.6607142857142857, 'x[11366] <= 0.082\ngini = 0.039\nsamples = 1750\nvalue = [1715, 35]'),
Text(0.44220183486238535, 0.625, 'x[39836] <= 0.082\ngini = 0.036\nsamples = 1747\nvalue = [1715, 32]'),
Text(0.3963302752293578, 0.5892857142857143, 'x[30114] <= 0.082\ngini = 0.032\nsamples = 1741\nvalue = [1713, 28]'),
Text(0.3486238532110092, 0.5535714285714286, 'x[11408] <= 0.076\ngini = 0.027\nsamples = 1735\nvalue = [1711, 24]'),
Text(0.29724770642201837, 0.5178571428571429, 'x[35084] <= 0.123\ngini = 0.023\nsamples = 1728\nvalue = [1708, 20]'),
Text(0.26788990825688075, 0.48214285714285715, 'x[49688] <= 0.023\ngini = 0.02\nsamples = 1724\nvalue = [1707, 17]'),
Text(0.25321100917431194, 0.44642857142857145, 'x[28346] <= 0.135\ngini = 0.017\nsamples = 1722\nvalue = [1707, 15]'),
Text(0.21284403669724772, 0.4107142857142857, 'x[39102] <= 0.076\ngini = 0.015\nsamples = 1719\nvalue = [1706, 13]'),
Text(0.1761467889908257, 0.375, 'x[6931] <= 0.052\ngini = 0.013\nsamples = 1716\nvalue = [1705, 11]'),
Text(0.14678899082568808, 0.3392857142857143, 'x[40841] <= 0.015\ngini = 0.01\nsamples = 1713\nvalue = [1704, 9]'),
Text(0.13211009174311927, 0.30357142857142855, 'x[30559] <= 0.122\ngini = 0.009\nsamples = 1712\nvalue = [1704, 8]'),
Text(0.11743119266055047, 0.26785714285714285, 'x[48312] <= 0.173\ngini = 0.008\nsamples = 1711\nvalue = [1704, 7]'),
Text(0.10275229357798166, 0.23214285714285715, 'x[53919] <= 0.087\ngini = 0.007\nsamples = 1710\nvalue = [1704, 6]'),
Text(0.08807339449541285, 0.19642857142857142, 'x[29608] <= 0.159\ngini = 0.006\nsamples = 1709\nvalue = [1704, 5]'),
Text(0.07339449541284404, 0.16071428571428573, 'x[8947] <= 0.102\ngini = 0.005\nsamples = 1708\nvalue = [1704, 4]'),
Text(0.05871559633027523, 0.125, 'x[50332] <= 0.287\ngini = 0.004\nsamples = 1707\nvalue = [1704, 3]'),
Text(0.044036697247706424, 0.08928571428571429, 'x[24241] <= 0.129\ngini = 0.002\nsamples = 1706\nvalue = [1704, 2]'),
Text(0.029357798165137616, 0.05357142857142857, 'x[36100] <= 0.069\ngini = 0.001\nsamples = 1705\nvalue = [1704, 1]'),
Text(0.014678899082568808, 0.017857142857142856, 'gini = 0.0\nsamples = 1704\nvalue = [1704, 0]'),
Text(0.044036697247706424, 0.017857142857142856, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.05871559633027523, 0.05357142857142857, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.07339449541284404, 0.08928571428571429, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.08807339449541285, 0.125, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.10275229357798166, 0.16071428571428573, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.11743119266055047, 0.19642857142857142, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.13211009174311927, 0.23214285714285715, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.14678899082568808, 0.26785714285714285, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.1614678899082568, 0.30357142857142855, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.20550458715596331, 0.3392857142857143, 'x[8377] <= 0.094\ngini = 0.444\nsamples = 3\nvalue = [1, 2]'),
Text(0.1908256880733945, 0.30357142857142855, 'gini = 0.0\nsamples = 2\nvalue = [0, 2]'),
Text(0.22018348623853212, 0.30357142857142855, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.24954128440366974, 0.375, 'x[6002] <= 0.03\ngini = 0.444\nsamples = 3\nvalue = [1, 2]'),
Text(0.23486238532110093, 0.3392857142857143, 'gini = 0.0\nsamples = 2\nvalue = [0, 2]'),
Text(0.26422018348623855, 0.3392857142857143, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.29357798165137616, 0.4107142857142857, 'x[53983] <= 0.049\ngini = 0.444\nsamples = 3\nvalue = [1, 2]'),
Text(0.27889908256880735, 0.375, 'gini = 0.0\nsamples = 2\nvalue = [0, 2]'),
Text(0.30825688073394497, 0.375, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.28256880733944956, 0.44642857142857145, 'gini = 0.0\nsamples = 2\nvalue = [0, 2]'),
Text(0.326605504587156, 0.48214285714285715, 'x[1745] <= 0.092\ngini = 0.375\nsamples = 4\nvalue = [1, 3]'),
Text(0.3119266055045872, 0.44642857142857145, 'gini = 0.0\nsamples = 3\nvalue = [0, 3]'),
Text(0.3412844036697248, 0.44642857142857145, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.4, 0.5178571428571429, 'x[35511] <= 0.059\ngini = 0.49\nsamples = 7\nvalue = [3, 4]'),
Text(0.3853211009174312, 0.48214285714285715, 'x[33644] <= 0.062\ngini = 0.375\nsamples = 4\nvalue = [3, 1]'),
Text(0.3706422018348624, 0.44642857142857145, 'gini = 0.0\nsamples = 3\nvalue = [3, 0]'),
Text(0.4, 0.44642857142857145, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),
Text(0.41467889908256883, 0.48214285714285715, 'gini = 0.0\nsamples = 3\nvalue = [0, 3]'),
Text(0.44403669724770645, 0.5535714285714286, 'x[54053] <= 0.018\ngini = 0.444\nsamples = 6\nvalue = [2, 4]'),
Text(0.42935779816513764, 0.5178571428571429, 'gini = 0.0\nsamples = 4\nvalue = [0, 4]'),
Text(0.45871559633027525, 0.5178571428571429, 'gini = 0.0\nsamples = 2\nvalue = [2, 0]'),
Text(0.48807339449541287, 0.5892857142857143, 'x[12037] <= 0.017\ngini = 0.444\nsamples = 6\nvalue = [2, 4]'),
Text(0.47339449541284406, 0.5535714285714286, 'gini = 0.0\nsamples = 4\nvalue = [0, 4]'),
Text(0.5027522935779817, 0.5535714285714286, 'gini = 0.0\nsamples = 2\nvalue = [2, 0]'),
Text(0.47155963302752296, 0.625, 'gini = 0.0\nsamples = 3\nvalue = [0, 3]'),
Text(0.5155963302752293, 0.6607142857142857, 'x[9817] <= 0.012\ngini = 0.32\nsamples = 5\nvalue = [1, 4]'),
Text(0.5009174311926605, 0.625, 'gini = 0.0\nsamples = 4\nvalue = [0, 4]'),
Text(0.5302752293577981, 0.625, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.5155963302752293, 0.6964285714285714, 'gini = 0.0\nsamples = 5\nvalue = [0, 5]'),
Text(0.5596330275229358, 0.7321428571428571, 'x[10456] <= 0.017\ngini = 0.397\nsamples = 11\nvalue = [3, 8]'),
Text(0.544954128440367, 0.6964285714285714, 'gini = 0.0\nsamples = 8\nvalue = [0, 8]'),
Text(0.5743119266055046, 0.6964285714285714, 'gini = 0.0\nsamples = 3\nvalue = [3, 0]'),
Text(0.6036697247706422, 0.7678571428571429, 'x[6950] <= 0.023\ngini = 0.198\nsamples = 9\nvalue = [1, 8]'),
Text(0.5889908256880734, 0.7321428571428571, 'gini = 0.0\nsamples = 8\nvalue = [0, 8]'),
Text(0.618348623853211, 0.7321428571428571, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]'),
Text(0.6477064220183486, 0.8035714285714286, 'x[22172] <= 0.021\ngini = 0.469\nsamples = 24\nvalue = [9, 15]'),
Text(0.6330275229357798, 0.7678571428571429, 'gini = 0.0\nsamples = 7\nvalue = [7, 0]'),
Text(0.6623853211009174, 0.7678571428571429, 'x[51183] <= 0.018\ngini = 0.208\nsamples = 17\nvalue = [2, 15]')
]

```

```
1 # Checking Classification report
```

```
2 print(classification_report(y_test, y_pred))
```

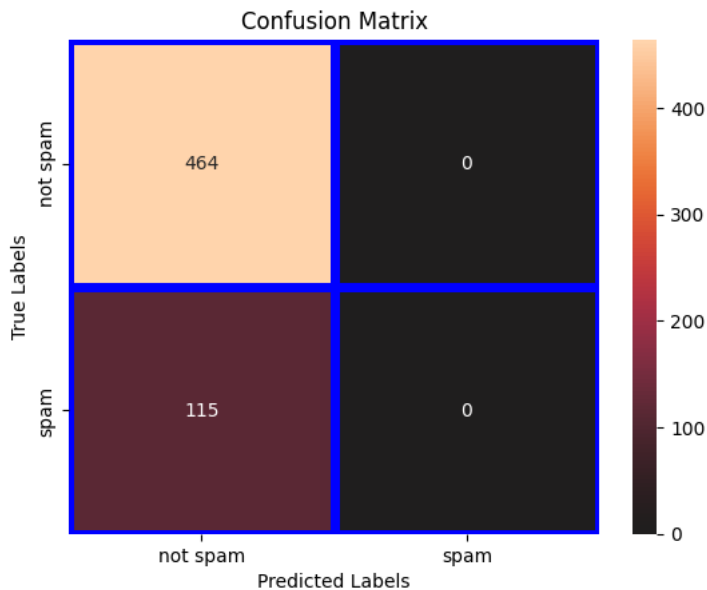
	precision	recall	f1-score	support
0	0.80	1.00	0.89	464
1	0.00	0.00	0.00	115
accuracy			0.80	579
macro avg	0.40	0.50	0.44	579
weighted avg	0.64	0.80	0.71	579

```
Text(0.8385321100917431, 0.8392857142857143, 'gini = 0.0\nsamples = 1\nvalue = [1, 0]')
```

```

1 from sklearn.metrics import confusion_matrix
2
3 Text/(0.8678899082568807, 0.8392857142857143, 'xrf478271 <= 0.016\ndini = 0.083\nsamples = 185\nvalue = {8. 1771}')
4
5 conf_mat = confusion_matrix(y_test,y_pred)
6
7 ax = plt.subplot()
8
9 sns.heatmap(conf_mat, annot=True, ax=ax, linewidths=5, linecolor='b', center=0, fmt='g')
10
11 ax.set_xlabel('Predicted Labels')
12 ax.set_ylabel('True Labels')
13 ax.set_title('Confusion Matrix')
14 ax.xaxis.set_ticklabels(['not spam','spam'])
15 ax.yaxis.set_ticklabels(['not spam','spam'])
16
17 plt.show()
18

```



```

1
2
3 # Generate a random classification dataset
4 X, y = make_classification(n_samples=1000, n_classes=2, random_state=42)
5
6 # Split the dataset into training and testing sets
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
8
9 # Train a logistic regression model on the training set
10 model = LogisticRegression(random_state=42)
11 model.fit(X_train, y_train)
12
13 # Use the model to make predictions on the testing set
14 y_pred = model.predict_proba(X_test)[:, 1]
15 y_true = y_test
16
17
18 #pr curve
19
20 # Assume y_true and y_pred are the true labels and predicted probabilities, respectively
21 precision, recall, thresholds = precision_recall_curve(y_true, y_pred)
22 average_precision = average_precision_score(y_true, y_pred)
23
24 # Plot the PR curve
25 plt.step(recall, precision, color='b', alpha=0.2, where='post')
26 plt.fill_between(recall, precision, step='post', alpha=0.2, color='b')
27
28 # Add labels and a legend to the plot
29 plt.xlabel('Recall')
30 plt.ylabel('Precision')
31 plt.ylim([0.0, 1.05])
32 plt.xlim([0.0, 1.0])
33 plt.title('Precision-Recall curve: AP={0:0.2f}'.format(average_precision))
34 plt.show()
35

```



