

```
In [1]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from matplotlib import pyplot as plt
from sklearn.feature_extraction.text import TfidfVectorizer

In [2]: mail=pd.read_csv("C:\\Users\\Pranav\\Desktop\\DATA SCIENCE DATA\\CVC file\\mail_data.csv")
mail.head()
```

Out[2]:

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

```
In [3]: #number of rows and columns
mail.shape
```

Out[3]: (5572, 2)

```
In [4]: #information about data
mail.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  -
0    Category    5572 non-null   object
1    Message     5572 non-null   object
dtypes: object(2)
memory usage: 87.2+ KB
```

```
In [5]: mail['Category'].value_counts()
```

Out[5]:

ham	4825
spam	747

Name: Category, dtype: int64

```
In [6]: #replacing the Ham as 0 and spam as 1
mail.replace({'Category':{'spam':0,'ham':1}},inplace=True)
mail.head()
```

Out[6]:

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

```
In [7]: #separting data
X=mail['Message']
y=mail['Category']
print(X)
print(y)

0      Go until jurong point, crazy.. Available only ...
1      Ok lar... Joking wif u oni...
2      Free entry in 2 a wkly comp to win FA Cup fina...
3      U dun say so early hor... U c already then say...
4      Nah I don't think he goes to usf, he lives aro...
...
5567    This is the 2nd time we have tried 2 contact u...
5568           Will ù b going to esplanade fr home?
5569    Pity, * was in mood for that. So...any other s...
5570    The guy did some bitching but I acted like i'd...
5571           Rofl. Its true to its name
Name: Message, Length: 5572, dtype: object
0      1
1      1
2      0
3      1
4      1
...
5567    0
5568    1
5569    1
5570    1
5571    1
Name: Category, Length: 5572, dtype: int64
```

```
In [8]: #training and test data
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=2)
print("shape of X_train= ",X_train.shape)
print("shape of X_test= ",X_test.shape)
print("shape of y_train= ",y_train.shape)
print("shape of y_test= ",y_test.shape)

shape of X_train=  (4457,)
shape of X_test=   (1115,)
shape of y_train=  (4457,)
shape of y_test=   (1115,)
```

```
In [9]: #Feature Extraction
#Transform the text data to feature vectors that can be used as input to the Logistics Regression
feature_extraction=TfidfVectorizer(min_df=1,stop_words='english',lowercase='True')
```

```
In [10]: feature_extraction
```

Out[10]: TfidfVectorizer(lowercase='True', stop_words='english')

```
In [18]: X_train_features=feature_extraction.fit_transform(X_train)
X_test_features=feature_extraction.transform(X_test)
```

```
In [19]: #conver y train and y test integer
y_train=y_train.astype('int')
y_test=y_test.astype('int')
```

```
In [20]: print(X_test_features)

(0, 6619)    0.33077540807715927
(0, 4752)    0.44421921026428457
(0, 2494)    0.359541012283057
(0, 2313)    0.37081499071603014
(0, 2110)    0.2538341210056606
(0, 1623)    0.47755798461662824
(0, 1153)    0.3660464944955722
(1, 4140)    0.7724156535136
(1, 3802)    0.40629294786687964
(1, 3352)    0.4881599110135932
(2, 3179)    0.3405136304031059
(2, 3169)    0.9402395798463798
(3, 6670)    0.4948874540031021
(3, 6543)    0.5505088255084791
(3, 2900)    0.6723291165103608
(4, 7417)    0.4582086641273852
(4, 6613)    0.6612385994559425
(4, 5583)    0.3946308162640678
(4, 1764)    0.443931136059295
(5, 7144)    0.2525030795568811
(5, 6017)    0.3435042181615311
(5, 5522)    0.37192637792006283
(5, 4761)    0.3253891605505013
(5, 4161)    0.44233446097815598
(5, 4048)    0.23654956954038084
:
:
(1111, 5132) 0.4888630580390552
(1111, 5071) 0.3867437918860694
(1111, 4094) 0.24494882973980492
(1111, 3138) 0.24402169398619392
(1111, 3084) 0.24749503861730665
(1111, 1031) 0.4888630580390552
(1112, 7203) 0.6546374185867087
(1112, 4471) 0.7559430204626075
(1113, 7417) 0.5146241230268624
(1113, 6304) 0.6835461063738834
(1113, 861)  0.5176163950841749
(1114, 6855) 0.15064835569263915
(1114, 5214) 0.21778432884602225
(1114, 4790) 0.20816334585240823
(1114, 4718) 0.23425427376646862
(1114, 4382) 0.231344342775171
(1114, 4379) 0.231344342775171
(1114, 4330) 0.19554545364082745
(1114, 3964) 0.2606227394501477
(1114, 3928) 0.2912663505498453
(1114, 3871) 0.3992082760935345
(1114, 2348) 0.44776345719647237
(1114, 1556) 0.24096532576878502
(1114, 1355) 0.24937006166328782
(1114, 50)  0.231344342775171
```

```
In [21]: # training the Logistics Regression model with the training data
model=LogisticRegression()
model.fit(X_train_features,y_train)
```

Out[21]: LogisticRegression()

```
In [22]: #prediction on training data
predction_on_training_data=model.predict(X_train_features)
training_data_accuracy=accuracy_score(predction_on_training_data,y_train)
print('Accuracy on training data:',training_data_accuracy*100)

Accuracy on training data: 96.83643706529055
```

```
In [25]: predction_on_testing_data=model.predict(X_test_features)
testing_data_accuracy=accuracy_score(predction_on_testing_data,y_test)
print('Accuracy on testing data:',testing_data_accuracy*100)

Accuracy on testing data: 95.24663677130046
```

```
In [26]: input_mail1='Thanks for your subscription to Ringtone UK your mobile will be charged £5/month Please confirm by replying YES or NO. If you reply NO you will not be charged']
input_data_features=feature_extraction.transform(input_mail1)
```

```
In [40]: predcection=model.predict(input_data_features)
predcection
```

Out[40]: array([0])

```
In [45]: if(predcection[0]==1):
print('Hub mail')
else:
print('Spam mail')

Spam mail
```

```
In [47]: input_mail2=['Ok lar... Joking wif u oni...']
input_data_features2=feature_extraction.transform(input_mail2)
```

```
In [48]: predcection=model.predict(input_data_features2)
predcection
```

Out[48]: array([1])

```
In [49]: if(predcection[0]==1):
print('Hub mail')
else:
print('Spam mail')

Hub mail
```

In []: