**FLIP ROBO**

# Email Spam Detection

Submitted by:

Lahari S.K

# INTRODUCTION

Spam emails are messages randomly sent to multiple addresses of all types of group but mostly by lazy advertisers and criminals leading to phishing sites.

Spam emails can not only be annoying but also dangerous. It is very important to detect spam emails and segregate them because they may cause following problems

- Critical mails are missed or delayed
- Millions of computers compromised
- Billions of dollars lost worldwide
- Identity theft
- Spam can crash mail servers and fill up hard drives.
-

Objective of this project is to give customers the knowledge of relevant emails and the fake emails.

# Data Description

The dataset contains 2893 rows and 3 columns. The columns are namely message, subject and the label. The label column is the target column. The dataset is in the form of CSV.

Out of 2893 instances 2412 is spam and 481 is ham mail. That is 83% of the given instances belongs to spam mails and the remaining are ham mails.

```
1  #Reading data
2  data
```

|  | 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 |
| ... | ... | ... | ... |
| 2888 | love your profile - ysuolvpv | hello thanks for stopping by ! ! we have taken... | 1 |
| 2889 | you have been asked to join kiddin | the list owner of : " kiddin " has invited you... | 1 |
| 2890 | anglicization of composers ' names | judging from the return post , i must have sou... | 0 |
| 2891 | re : 6 . 797 , comparative method : n - ary co... | gotcha ! there are two separate fallacies in t... | 0 |
| 2892 | re : american - english in australia | hello ! i ' m working on a thesis concerning a... | 0 |

2893 rows × 3 columns

# Data Pre-processing Done

## ◆ Creating additional columns to know length

Let's create two more columns known namely "sublength" and "messagelength" in order to calculate the length of the subject and message.

```
1 data["sublength"]=data.subject.str.len()
2 data["messagelength"]=data.message.str.len()
```

```
1 data
```

| | subject | message | label | sublength | messagelength |
|---|---|---|---|---|---|
| 0 | job posting - apple-iss research center | content - length : 3386 apple-iss research cen... | 0 | 39.0 | 2856 |
| 1 | NaN | lang classification grimes , joseph e . and ba... | 0 | NaN | 1800 |
| 2 | query : letter frequencies for text identifica... | i am posting this inquiry for sergei atamas ( ... | 0 | 50.0 | 1435 |
| 3 | risk | a colleague and i are researching the differin... | 0 | 4.0 | 324 |
| 4 | request book information | earlier this morning i was on the phone with a... | 0 | 24.0 | 1046 |
| ... | ... | ... | ... | ... | ... |
| 2888 | love your profile - ysuolvpv | hello thanks for stopping by ! ! we have taken... | 1 | 28.0 | 262 |
| 2889 | you have been asked to join kiddin | the list owner of : " kiddin " has invited you... | 1 | 34.0 | 2163 |
| 2890 | anglicization of composers ' names | judging from the return post , i must have sou... | 0 | 34.0 | 1039 |
| 2891 | re : 6 . 797 , comparative method : n - ary co... | gotcha ! there are two separate fallacies in t... | 0 | 54.0 | 2949 |
| 2892 | re : american - english in australia | hello ! i ' m working on a thesis concerning a... | 0 | 36.0 | 700 |

2893 rows × 5 columns

## ◆ Converting to lower case

Let's convert the subject and message column into lower case using the lower() function.

```
1 #Converting subject and messageto lower
2 data['subject']=data['subject'].str.lower()
3 data['message']=data['message'].str.lower()
```

```
1 data
```

| | subject | message | label | sublength | messagelength |
|---|---|---|---|---|---|
| 0 | job posting - apple-iss research center | content - length : 3386 apple-iss research cen... | 0 | 39.0 | 2856 |
| 1 | NaN | lang classification grimes , joseph e . and ba... | 0 | NaN | 1800 |
| 2 | query : letter frequencies for text identifica... | i am posting this inquiry for sergei atamas ( ... | 0 | 50.0 | 1435 |
| 3 | risk | a colleague and i are researching the differin... | 0 | 4.0 | 324 |
| 4 | request book information | earlier this morning i was on the phone with a... | 0 | 24.0 | 1046 |
| ... | ... | ... | ... | ... | ... |
| 2888 | love your profile - ysuolvpv | hello thanks for stopping by ! ! we have taken... | 1 | 28.0 | 262 |
| 2889 | you have been asked to join kiddin | the list owner of : " kiddin " has invited you... | 1 | 34.0 | 2163 |
| 2890 | anglicization of composers ' names | judging from the return post , i must have sou... | 0 | 34.0 | 1039 |
| 2891 | re : 6 . 797 , comparative method : n - ary co... | gotcha ! there are two separate fallacies in t... | 0 | 54.0 | 2949 |
| 2892 | re : american - english in australia | hello ! i ' m working on a thesis concerning a... | 0 | 36.0 | 700 |

## Replacing regular expression

Lets make use of regular expressions to replace some strings.

```
1  #Replacing emailaddresses with email
2  data["message"]=data["message"].str.replace(r'^.+@[^\.].*\.[a-z]{2,}$','email_address')
3
4  #Replacing webaddress
5  data["message"]=data["message"].str.replace(r'^http\://[a-zA-Z0-9\-\.]+\.[a-zA-Z]{2,3}(/\S*)?$','webaddress')
6
7  #Replacing moneysymbols
8  data["message"]=data["message"].str.replace(r'£|/$','dollars')
9
10 #Replacing phonenumber,paranthesis,spaces,dashes,nospaces)
11 data['message']=data['message'].str.replace(r'^\(?[\d]{3}\)?[\s-]?[\d]{3}[\s-]?[\d]{4}$','phone number')
12
13 #Replacing phone number with phone
14 data['message']=data['message'].replace(r'\d+(\.\d+)?','number')
```

## Removing punctuation

Lets remove all the punctuation marks .

```
1  #Removing punctuation marks
2  data['message']=data['message'].str.replace(r'[^\w\d\s]',' ')
3
4  #Removing whitespaes between terms with a single space
5  data['message']=data['message'].str.replace(r'\s+',' ')
6
7  #Removing tailing and heading whitespaces
8  data['message']=data['message'].str.replace(r'^\s+|\s+?$',' ')
9
```

# Removing stop words

Lets remove all the stop words and see the length after pre-processing steps.

```python
1  from nltk.corpus import stopwords
2  import string
3  import nltk
4
5  stop_words=stopwords.words("english")
6
7  data['message']=data['message'].apply(lambda x:' '.join(term for term in x.split() if term not in stop_words))
8
```

```python
1  #New column after cleasing
2  data["clean_length_message"]=data.message.str.len()
3  data["clean_length_subject"]=data.subject.str.len()
```

```python
1  #Total length removal
2  print("Original length of message",data.messagelength.sum())
3  print("After cleansing",data.clean_length_message.sum())
4
5  print("Original legth ofsubject",data.sublength.sum())
6  print("Aftre cleansing ",data.clean_length_subject.sum())
```

```
Original length of message 9070005
After cleansing 6308430
Original legth ofsubject 91647.0
Aftre cleansing  84741.0
```

# Model/s Development and Evaluation

## Algorithms Used

## 1.Naive Bayes

```
1  naive.fit(X_train,Y_train)
2  y_pred=naive.predict(x_test)
3  print("Finle score -->",accuracy_score(y_test,y_pred))
```

Finle score --> 0.835635359116022

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

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.83      | 1.00   | 0.91     | 585     |
| 1            | 1.00      | 0.14   | 0.25     | 139     |
|              |           |        |          |         |
| accuracy     |           |        | 0.84     | 724     |
| macro avg    | 0.92      | 0.57   | 0.58     | 724     |
| weighted avg | 0.86      | 0.84   | 0.78     | 724     |

## 2.SVM

```
1  from sklearn import svm
```

```
1  svm=svm.SVC()
2  svm.fit(X_train,Y_train)
3  y_pred=svm.predict(x_test)
4  print("Finle score -->",accuracy_score(y_test,y_pred))
```

Finle score --> 0.9737569060773481

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

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.97      | 1.00   | 0.98     | 585     |
| 1            | 1.00      | 0.86   | 0.93     | 139     |
|              |           |        |          |         |
| accuracy     |           |        | 0.97     | 724     |
| macro avg    | 0.98      | 0.93   | 0.96     | 724     |
| weighted avg | 0.97      | 0.97   | 0.97     | 724     |

The best model is SVM.