# Spam_SMS_email_Classifier

*-Mayank Srivastava*

- (https://linkedin.com/in/mayank-srivastava-6a8421105/)



Dataset: https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset (https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset)

```
# 1. Data cleaning
# 2. EDA
# 3. Text Preprocessing
# 4. Model building
# 5. Evaluation
# 6. Improvement
# 7. Website
# 8. Deploy
```

In [1]:
```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

In [2]:
```python
df= pd.read_csv('spam.csv', encoding='latin-1')
df.head()
```

Out[2]:

|   | v1 | v2 | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 |
|---|----|----|-----------|-----------|-----------|
| 0 | ham | Go until jurong point, crazy.. Available only ... | NaN | NaN | NaN |
| 1 | ham | Ok lar... Joking wif u oni... | NaN | NaN | NaN |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... | NaN | NaN | NaN |
| 3 | ham | U dun say so early hor... U c already then say... | NaN | NaN | NaN |
| 4 | ham | Nah I don't think he goes to usf, he lives aro... | NaN | NaN | NaN |

```
In [3]:  ▶   1  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   v1          5572 non-null   object
 1   v2          5572 non-null   object
 2   Unnamed: 2  50 non-null     object
 3   Unnamed: 3  12 non-null     object
 4   Unnamed: 4  6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB
```

```
Observations:

- No Null values
- size  = 5572 x 4
```

## 1. Data Cleaning

```
In [4]:  ▶   1  # Dropping the extra columns
              2  df.drop(columns=['Unnamed: 2','Unnamed: 3','Unnamed: 4'], inplace = True)
              3
              4  #renaming the columns
              5  df.rename(columns ={'v1': 'label', 'v2': 'text'}, inplace = True)
              6
              7  # replacing ham ->0 , spam -> 1
              8  df.label.replace(['ham','spam'],[0,1], inplace = True)
              9
             10  df.head()
```

Out[4]:

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

```
In [5]:  ▶   1  # checking for dupplicates
              2
              3  df.duplicated().sum()
```

Out[5]: 403

```
In [6]:  ▶   1  # removing duplicates
              2  df =df[~df.duplicated()]
              3  df.shape
```
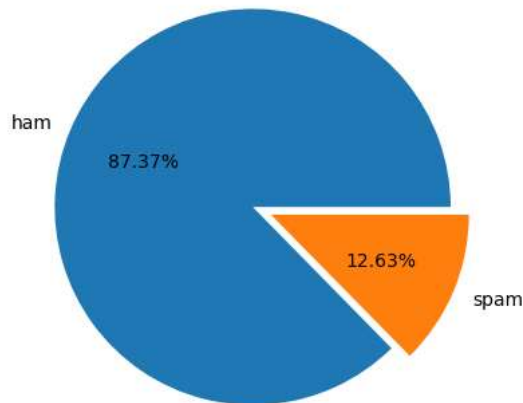
Out[6]: (5169, 2)

```
In [7]:  ▶   1  5572-403
```

Out[7]: 5169

## 2. EDA

```
In [8]:   1  data= df.label.value_counts()
          2  plt.pie(data.values, labels=['ham', 'spam'], autopct ="%1.2f%%", explode = [0.0,0.1])
          3  plt.title('Data Distribution Spam (1) and Not-Spam (0)')
          4  plt.show()
          5  data
```

Data Distribution Spam (1) and Not-Spam (0)



```
Out[8]:  label
         0    4516
         1     653
         Name: count, dtype: int64
```

**Feature Engineering**

```
In [9]:   1  # lets find no. of characters, words and sentences in spam and ham
          2  import nltk
          3  from nltk import word_tokenize,sent_tokenize
          4
          5  df['n_chars'] = df.text.str.len()
          6  df['n_words'] =df.text.apply(lambda x: len(word_tokenize(x)))
          7  df['n_sent'] =df.text.apply(lambda x: len(sent_tokenize(x)))
          8  df.head()
```

Out[9]:

|   | label | text | n_chars | n_words | n_sent |
|---|-------|------|---------|---------|--------|
| 0 | 0 | Go until jurong point, crazy.. Available only ... | 111 | 24 | 2 |
| 1 | 0 | Ok lar... Joking wif u oni... | 29 | 8 | 2 |
| 2 | 1 | Free entry in 2 a wkly comp to win FA Cup fina... | 155 | 37 | 2 |
| 3 | 0 | U dun say so early hor... U c already then say... | 49 | 13 | 1 |
| 4 | 0 | Nah I don't think he goes to usf, he lives aro... | 61 | 15 | 1 |

```
In [10]:  1  # plotting spam vs ham based on no. of characters
          2  plt.figure(figsize =(10,4))
          3  sns.histplot(df[df.label == 0]['n_chars'])
          4  sns.histplot(df[df.label == 1]['n_chars'], color= 'red')
```
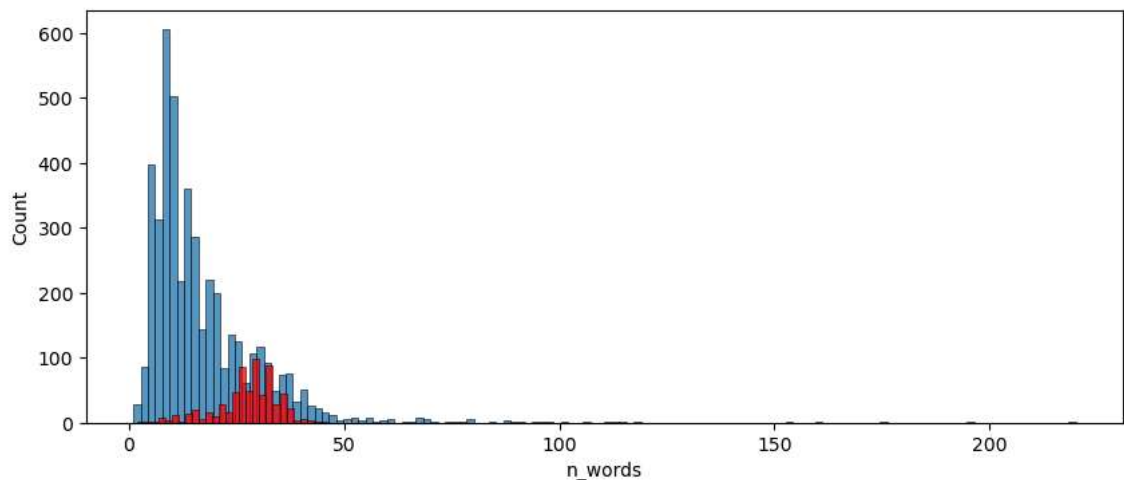
Out[10]: <Axes: xlabel='n_chars', ylabel='Count'>

Observation: Spam text has generally higher character count than ham

In [11]: ▶|
```python
1  # plotting spam vs ham based on no. of words
2  plt.figure(figsize =(10,4))
3  sns.histplot(df[df.label == 0]['n_words'])
4  sns.histplot(df[df.label == 1]['n_words'], color= 'red')
```

Out[11]: <Axes: xlabel='n_words', ylabel='Count'>

Observation: Spam text has generally higher word count than ham

In [12]: ▶|
```python
1  df.groupby('label')[['n_chars','n_words','n_sent']].mean().transpose().rename(columns ={0:'ham', 1:'spam'})
```

Out[12]:

| label | ham | spam |
|---|---|---|
| n_chars | 70.459256 | 137.891271 |
| n_words | 17.123782 | 27.667688 |
| n_sent | 1.820195 | 2.970904 |

Observation- On an average:
- ham text has avg 70 characters while spam has almost double characters.
- ham text has 17 words, spam has 27 words
- ham has ~ 2 sentences, spam has ~3 sentences.

## 3. Text preprocessing

In [13]: ▶|
```python
1  import nltk
2
3  nltk.download('stopwords')
4  nltk.download('punkt')
5  from nltk.corpus import stopwords
6
7  from nltk.stem import PorterStemmer
8  from nltk.stem import WordNetLemmatizer
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\hp\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\hp\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

In [14]: ▶|
```python
1  # lets define a helper fuction for standardization of text
```

In [15]: ▶|
```python
1  # remove punctuations
2  import string
3  punctuation= string.punctuation
4  punctuation
```

Out[15]: '!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~'

```
In [16]:  ▶|    1  # stop words
               2  from nltk.corpus import stopwords
               3  stops= set(stopwords.words('english'))
               4  stops
```

```
Out[16]: {'a',
          'about',
          'above',
          'after',
          'again',
          'against',
          'ain',
          'all',
          'am',
          'an',
          'and',
          'any',
          'are',
          'aren',
          "aren't",
          'as',
          'at',
          'be',
          'because',
```

```
In [17]:  ▶|    1  import re
               2  stem =PorterStemmer()
               3  lemm = WordNetLemmatizer()
               4  def standardize(text):
               5
               6      # change to lower case
               7      text = text.lower()
               8
               9      # keep only alpha numerics
              10      # assuming _ sign is used for space in text, replacing it with space
              11      text =re.sub('_','',text)
              12      text =re.findall(r"\w+", text)
              13
              14      # now text has been converted into list of words , after re.findall
              15
              16      # remove punctuations and stopwords
              17      text =[i for i in text if i not in stops and i not in punctuation]
              18
              19      # Lemmatization
              20      text= [lemm.lemmatize(i) for i in text]
              21
              22      # stemming
              23      text= [stem.stem(i) for i in text]
              24
              25      return(' '.join(text))
```

```
In [18]:  ▶|    1  standardize("Hello' there I am peter, from22nd_street the mayor's office, who__ are. you? Dear %. Dancing ate remote
```

```
Out[18]: 'hello peter from22ndstreet mayor offic dear danc ate remot'
```

```
In [19]:  ▶|    1  df['std_text'] = df.text.apply(standardize)
               2  df.std_text
```

```
Out[19]: 0           go jurong point crazi avail bugi n great world...
         1                                 ok lar joke wif u oni
         2       free entri 2 wkli comp win fa cup final tkt 21...
         3                     u dun say earli hor u c alreadi say
         4                     nah think go usf life around though
                                       ...
         5567    2nd time tri 2 contact u u å 750 pound prize 2...
         5568                         ì b go esplanad fr home
         5569                                piti mood suggest
         5570    guy bitch act like interest buy someth els nex...
         5571                                   rofl true name
         Name: std_text, Length: 5169, dtype: object
```
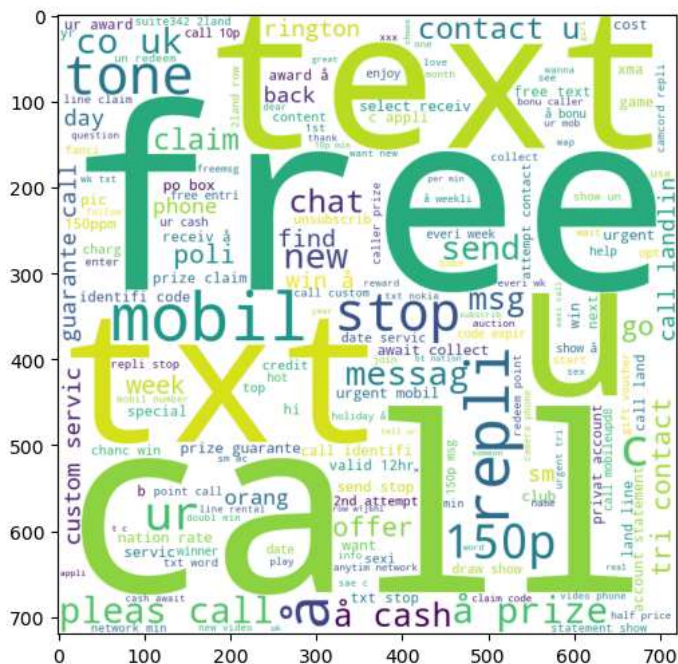
```
In [20]:  ▶|    1  # pip install wordcloud
```

```
In [21]:  ▶|    1  # word_cloud of spam messages
               2  from wordcloud import WordCloud
               3  wc= WordCloud(width =720, height =720, min_font_size =10, background_color ='white')
```

```
In [22]:  ▶|    1  spam_wc = wc.generate(df[df.label == 1]['std_text'].astype('str').str.cat(sep= " "))
               2  # concatentaing texts with seperator space
```

```
In [23]:  ▶|   1  plt.figure(figsize =(8,6))
             2  plt.imshow(spam_wc)
```

Out[23]: <matplotlib.image.AxesImage at 0x232ef89a810>



```
In [24]:  ▶|   1  ham_wc = wc.generate(df[df.label == 0]['std_text'].astype('str').str.cat(sep= " "))
             2  # concatentaing texts with seperator space
```

```
In [25]:  ▶|   1  plt.figure(figsize =(8,6))
             2  plt.imshow(ham_wc)
```

Out[25]: <matplotlib.image.AxesImage at 0x232ef9ba490>



```
In [26]:  ▶|   1  df[df.label == 1]['std_text']
```

Out[26]: 2        free entri 2 wkli comp win fa cup final tkt 21...
         5        freemsg hey darl 3 week word back like fun sti...
         8        winner valu network custom select receivea å 9...
         9        mobil 11 month u r entitl updat latest colour ...
         11       six chanc win cash 100 20 000 pound txt csh11 ...
                                        ...
         5537     want explicit sex 30 sec ring 02073162414 cost...
         5540     ask 3mobil 0870 chatlin inclu free min india c...
         5547     contract mobil 11 mnth latest motorola nokia e...
         5566     remind o2 get 2 50 pound free call credit deta...
         5567     2nd time tri 2 contact u u å 750 pound prize 2...
         Name: std_text, Length: 653, dtype: object
```

```
In [27]:  ▶  1  # Top 30 words in spam,
             2  # made a single list, used value_counts
             3  spam_words =[]
             4  for i in df[df.label == 1]['std_text'].str.split(' '):
             5      spam_words+=i
             6  print(spam_words)
```
'1', 'minmobsmorelkpobox177hp51fl', 'urgent', 'tri', 'contact', 'u', 'today', 'draw', 'show', 'å', '800', 'prize', 'g
uarante', 'call', '09050001295', 'land', 'line', 'claim', 'a21', 'valid', '12hr', 'monthli', 'password', 'wap', 'mobs
i', 'com', '391784', 'use', 'wap', 'phone', 'pc', 'today', 'vodafon', 'number', 'end', '0089', 'last', 'four', 'digi
t', 'select', 'receiv', 'å', '350', 'award', 'number', 'match', 'pleas', 'call', '09063442151', 'claim', 'å', '350',
'award', 'free', 'top', 'rington', 'sub', 'weekli', 'rington', 'get', '1st', 'week', 'free', 'send', 'subpoli', '8161
8', '3', 'per', 'week', 'stop', 'sm', '08718727870', 'free', 'msg', 'sorri', 'servic', 'order', '81303', 'could', 'de
liv', 'suffici', 'credit', 'pleas', 'top', 'receiv', 'servic', 'hard', 'live', '121', 'chat', '60p', 'min', 'choos',
'girl', 'connect', 'live', 'call', '09094646899', 'cheap', 'chat', 'uk', 'biggest', 'live', 'servic', 'vu', 'bcm1896w
c1n3xx', 'wow', 'boy', 'r', 'back', 'take', '2007', 'uk', 'tour', 'win', 'vip', 'ticket', 'pre', 'book', 'vip', 'clu
b', 'txt', 'club', '81303', 'trackmarqu', 'ltd', 'info', 'vipclub4u', 'hi', 'mandi', 'sullivan', 'call', 'hotmix', 'f
m', 'chosen', 'receiv', 'å', '5000', '00', 'easter', 'prize', 'draw', 'pleas', 'telephon', '09041940223', 'claim', '2
9', '03', '05', 'prize', 'transfer', 'someon', 'els', 'ur', 'go', '2', 'bahama', 'callfreefon', '08081560665', 'spea
k', 'live', 'oper', 'claim', 'either', 'bahama', 'cruis', 'ofå', '2000', 'cash', '18', 'opt', 'txt', 'x', '0778620011
7', 'someon', 'conact', 'date', 'servic', 'enter', 'phone', 'fanci', 'find', 'call', 'landlin', '09111030116', 'pobox
12n146tf15', 'hi', '07734396839', 'ibh', 'custom', 'loyalti', 'offer', 'new', 'nokia6600', 'mobil', 'å', '10', 'txtau
ction', 'txt', 'word', 'start', '81151', 'get', '4t', 'sm', 'auction', 'nokia', '7250i', 'get', 'win', 'free', 'aucti
on', 'take', 'part', 'send', 'nokia', '86021', 'hg', 'suite342', '2land', 'row', 'w1jhl', '16', 'call', 'freephon',
'0800', '542', '0578', 'buy', 'space', 'invad', '4', 'chanc', '2', 'win', 'orig', 'arcad', 'game', 'consol', 'press',
'0', 'game', 'arcad', 'std', 'wap', 'charg', 'see', 'o2', 'co', 'uk', 'game', '4', 'term', 'set', 'purchas', 'big',
'brother', 'alert', 'comput', 'select', 'u', '10k', 'cash', '150', 'voucher', 'call', '09064018838', 'ntt', 'po', 'bo

```
In [28]:  ▶  1  len(spam_words)
```
Out[28]: 11996

```
In [29]:  ▶  1  from collections import Counter
             2  spam_count =Counter(spam_words)
             3  spam_count.most_common(30)
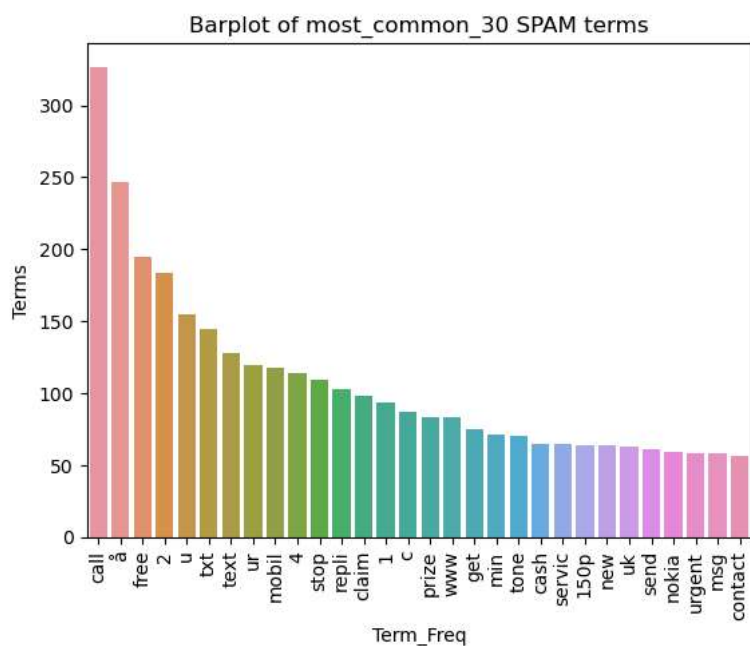```
Out[29]: [('call', 327),
         ('å', 247),
         ('free', 195),
         ('2', 184),
         ('u', 155),
         ('txt', 145),
         ('text', 128),
         ('ur', 119),
         ('mobil', 118),
         ('4', 114),
         ('stop', 109),
         ('repli', 103),
         ('claim', 98),
         ('1', 93),
         ('c', 87),
         ('prize', 83),
         ('www', 83),
         ('get', 75),
         ('min', 71),
         ('tone', 70),
         ('cash', 65),
         ('servic', 65),
         ('150p', 64),
         ('new', 64),
         ('uk', 63),
         ('send', 61),
         ('nokia', 59),
         ('urgent', 58),
         ('msg', 58),
         ('contact', 56)]

```
1  pd.DataFrame(spam_count.most_common(30))
```

Out[30]:

|    | 0       | 1   |
|----|---------|-----|
| 0  | call    | 327 |
| 1  | å       | 247 |
| 2  | free    | 195 |
| 3  | 2       | 184 |
| 4  | u       | 155 |
| 5  | txt     | 145 |
| 6  | text    | 128 |
| 7  | ur      | 119 |
| 8  | mobil   | 118 |
| 9  | 4       | 114 |
| 10 | stop    | 109 |
| 11 | repli   | 103 |
| 12 | claim   | 98  |
| 13 | 1       | 93  |
| 14 | c       | 87  |
| 15 | prize   | 83  |
| 16 | www     | 83  |
| 17 | get     | 75  |
| 18 | min     | 71  |
| 19 | tone    | 70  |
| 20 | cash    | 65  |
| 21 | servic  | 65  |
| 22 | 150p    | 64  |
| 23 | new     | 64  |
| 24 | uk      | 63  |
| 25 | send    | 61  |
| 26 | nokia   | 59  |
| 27 | urgent  | 58  |
| 28 | msg     | 58  |
| 29 | contact | 56  |

In [31]:

```
1  sns.barplot(x= pd.DataFrame(spam_count.most_common(30))[0],y= pd.DataFrame(spam_count.most_common(30))[1])
2  plt.xticks(rotation ='vertical')
3  plt.xlabel('Term_Freq')
4  plt.ylabel('Terms')
5  plt.title('Barplot of most_common_30 SPAM terms')
6  plt.show()
```

```python
# similarly for ham_words
ham_words = []
for i in df[df.label == 0]['std_text'].str.split():
    ham_words+=i
print(ham_words)
```

```
'soon', 'come', 'otherwis', 'tomorrow', 'msg', 'r', 'time', 'pa', 'silent', 'say', 'think', 'u', 'right', 'also', 'ma
ke', 'u', 'think', 'least', '4', 'moment', 'gd', 'nt', 'swt', 'drm', 'shesil', 'yeah', 'probabl', 'swing', 'roommat',
'finish', 'girl', 'happi', 'new', 'year', 'melodi', 'ìï', 'dun', 'need', 'pick', 'ur', 'gf', 'yay', 'better', 'told',
'5', 'girl', 'either', 'horribl', 'u', 'eat', 'mac', 'eat', 'u', 'forgot', 'abt', 'alreadi', 'rite', 'u', 'take', 'lo
ng', '2', 'repli', 'thk', 'toot', 'b4', 'b', 'prepar', 'wat', 'shall', 'eat', 'say', 'fantast', 'chanc', 'anyth', 'ne
ed', 'bigger', 'life', 'lift', 'lose', '2', 'live', 'think', 'would', 'first', 'person', '2', 'die', 'n', 'v', 'q',
'nw', 'came', 'hme', 'da', 'outsid', 'island', 'head', 'toward', 'hard', 'rock', 'run', 'day', 'class', 'class', 'che
nnai', 'velacheri', 'flippin', 'shit', 'yet', 'k', 'give', 'sec', 'break', 'lt', 'gt', 'cstore', 'much', 'bad', 'avoi
d', 'like', 'yo', 'around', 'got', 'car', 'back', 'annoy', 'goodmorn', 'today', 'late', 'lt', 'gt', 'min', 'point',
'hangin', 'mr', 'right', 'makin', 'u', 'happi', 'come', 'aliv', 'better', 'correct', 'good', 'look', 'fi
gur', 'case', 'guess', 'see', 'campu', 'lodg', 'done', 'come', 'home', 'one', 'last', 'time', 'wont', 'anyth', 'trus
t', 'night', 'worri', 'appt', 'shame', 'miss', 'girl', 'night', 'quiz', 'popcorn', 'hair', 'ok', 'c', 'ya', 'said',
'matter', 'mind', 'say', 'matter', 'al', 'moan', 'n', 'e', 'thin', 'go', 'wrong', 'fault', 'al', 'de', 'argument',
'r', 'fault', 'fed', 'himso', 'bother', 'hav', '2go', 'thanx', 'xx', 'neft', 'transact', 'refer', 'number', 'lt', 'g
t', 'r', 'lt', 'decim', 'gt', 'credit', 'beneficiari', 'account', 'lt', 'gt', 'lt', 'time', 'gt', 'lt', 'gt', 'otherw
is', 'part', 'time', 'job', 'na', 'tuition', 'know', 'call', 'also', 'da', 'feel', 'yesterday', 'night', 'wait', 'ti
l', '2day', 'night', 'dear', 'thank', 'understand', 'tri', 'tell', 'sura', 'whole', 'car', 'appreci', 'last', 'two',
'dad', 'map', 'read', 'semi', 'argument', 'apart', 'thing', 'go', 'ok', 'p', 'need', 'strong', 'arm', 'also', 'maaaa
n', 'miss', 'bday', 'real', 'april', 'guessin', 'gonna', '9', 'ok', 'come', 'ur', 'home', 'half', 'hour', 'yo', 'gam
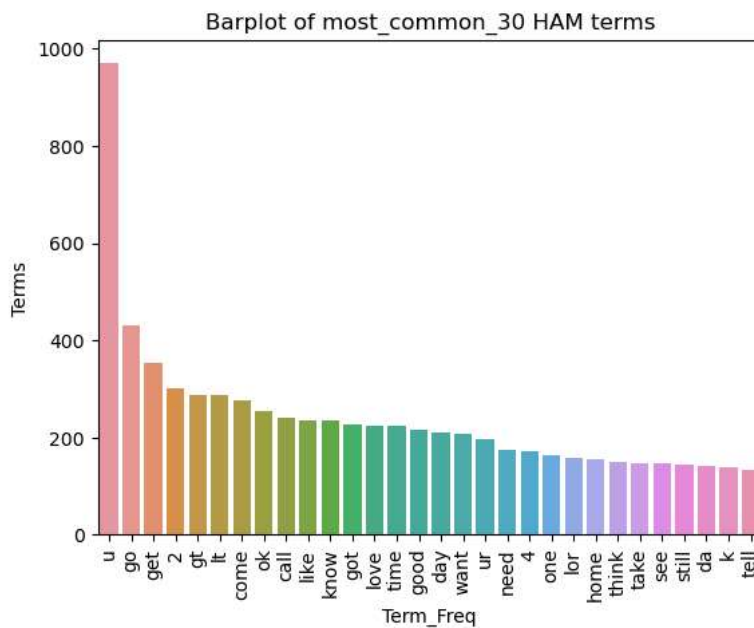e', 'almost', 'want', 'go', 'walmart', 'soon', 'yeah', 'probabl', 'sure', 'ilol', 'let', 'u', 'know', 'person', 'wuld
```

```python
len(ham_words)
```

```
36505
```

```python
from collections import Counter
ham_count =Counter(ham_words)
ham_count.most_common(30)
```

```
[('u', 969),
 ('go', 432),
 ('get', 354),
 ('2', 302),
 ('gt', 288),
 ('lt', 287),
 ('come', 276),
 ('ok', 255),
 ('call', 240),
 ('like', 236),
 ('know', 236),
 ('got', 226),
 ('love', 225),
 ('time', 223),
 ('good', 216),
 ('day', 210),
 ('want', 209),
 ('ur', 198),
 ('need', 174),
 ('4', 171),
 ('one', 165),
 ('lor', 159),
 ('home', 156),
 ('think', 150),
 ('take', 148),
 ('see', 147),
 ('still', 145),
 ('da', 143),
 ('k', 138),
 ('tell', 134)]
```

```
In [35]:  ▶  1  sns.barplot(x= pd.DataFrame(ham_count.most_common(30))[0],y= pd.DataFrame(ham_count.most_common(30))[1])
             2  plt.xticks(rotation ='vertical')
             3  plt.xlabel('Term_Freq')
             4  plt.ylabel('Terms')
             5  plt.title('Barplot of most_common_30 HAM terms')
             6  plt.show()
```



## 4. Text Vectorization

```
In [36]:  ▶  1  from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
             2  cv= CountVectorizer(stop_words = 'english')
             3  tv= TfidfVectorizer(stop_words = 'english')
             4
             5  X_cv= cv.fit_transform(df.std_text).toarray()
             6  X_tv= tv.fit_transform(df.std_text).toarray()
```

```
In [37]:  ▶  1  X_cv.shape, X_tv.shape
```

```
Out[37]:  ((5169, 7059), (5169, 7059))
```

- Both the vectorization methods give text vetors of same length

```
In [38]:  ▶  1  y=df.label
```

```
In [39]:  ▶  1  #Tf-idf  vectorization
             2  from sklearn.model_selection import cross_val_score, train_test_split
             3  xtrain, xtest, ytrain, ytest =train_test_split(X_tv, y, test_size =0.2, random_state =42)
```

## 5. Model Builidng

- Which is Better for Spam Detection?
    - If minimizing false positives is critical (i.e., you don't want to risk important emails being marked as spam), precision is more important.
    - This is often the case in professional or business environments where missing an important email can have serious consequences.
    - If minimizing false negatives is critical (i.e., you want to ensure that most spam emails are correctly identified), recall is more important.
    - This is often the case in personal email use where receiving spam emails is more of an annoyance than missing an important email.

```
In [40]:  ▶  1  # to check names of scorers for cross-val
             2  # import sklearn
             3  # sklearn.metrics.get_scorer_names()
```

```
In [41]:  ▶  1  from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
             2  from sklearn.model_selection import cross_val_score, train_test_split
             3  gnb= GaussianNB()
             4  mnb=MultinomialNB()
             5  bnb =BernoulliNB()
```

```
In [42]:  ▶  1  X={'BoW':X_cv, 'Tf-Idf':X_tv}
             2  model={'GaussianNB':gnb,'MultinomialNB':mnb,'BernoulliNB':bnb}
```

```
In [43]:    1  scores={}
            2  for i in X:
            3      for j in model:
            4          scores[i+'_'+j] = np.max(cross_val_score(model[j], X[i], y, cv=5, scoring ='precision'))
            5
            6      # our metric is recall score because, we dont want any false positives (Type 2 error), ie a Spam predicted as Ham
```

```
In [44]:    1  scores
```

```
Out[44]: {'BoW_GaussianNB': 0.5041322314049587,
          'BoW_MultinomialNB': 0.9333333333333333,
          'BoW_BernoulliNB': 1.0,
          'Tf-Idf_GaussianNB': 0.497907949790795,
          'Tf-Idf_MultinomialNB': 1.0,
          'Tf-Idf_BernoulliNB': 1.0}
```

```
In [45]:    1  sorted(scores.items(), key=lambda x: x[1], reverse = True)
```

```
Out[45]: [('BoW_BernoulliNB', 1.0),
          ('Tf-Idf_MultinomialNB', 1.0),
          ('Tf-Idf_BernoulliNB', 1.0),
          ('BoW_MultinomialNB', 0.9333333333333333),
          ('BoW_GaussianNB', 0.5041322314049587),
          ('Tf-Idf_GaussianNB', 0.497907949790795)]
```

- Clearly the winner is BoW with Multinomial Naive Bayes with the best f1_score

```
In [46]:    1  # training the model on 'BoW_BernoulliNB' or 'Tf-Idf_MultinomialNB' or 'Tf-Idf_BernoulliNB' is giving ideal Precisio
            2
            3  # randomly proceeding with Tf-Idf_MultinomialNB
            4  xtrain, xtest, ytrain, ytest =train_test_split(X_tv, y, test_size =0.2, random_state =42)
            5  mnb.fit(xtrain, ytrain)
```

```
Out[46]:  ▾ MultinomialNB
          MultinomialNB()
```

```
In [47]:    1  from sklearn.metrics import accuracy_score, precision_score, classification_report, confusion_matrix
            2  pred_train =mnb.predict(xtrain)
            3  pred_test =mnb.predict(xtest)
            4  print('training accuracy_score: ', accuracy_score(ytrain, pred_train))
            5  print('training precision_score: ', precision_score(ytrain, pred_train))
            6  print('training confusion_matrix: \n', confusion_matrix(ytrain, pred_train))
            7
            8  #predicting testing data
            9  print('testing accuracy_score: ', accuracy_score(ytest, pred_test))
           10  print('testing precision_score: ', precision_score(ytest, pred_test))
           11  print('testing confusion_matrix: \n', confusion_matrix(ytest, pred_test))
```

```
training accuracy_score:  0.9743651753325272
training precision_score:  1.0
training confusion_matrix:
 [[3627    0]
 [ 106  402]]
testing accuracy_score:  0.9661508704061895
testing precision_score:  1.0
testing confusion_matrix:
 [[889    0]
 [ 35 110]]
```

```
In [48]:    1  # 'Tf-Idf_BernoulliNB'
            2  xtrain, xtest, ytrain, ytest =train_test_split(X_tv, y, test_size =0.2, random_state =42)
            3  bnb.fit(xtrain, ytrain)
            4  from sklearn.metrics import accuracy_score, precision_score, classification_report, confusion_matrix
            5  pred_train =bnb.predict(xtrain)
            6  pred_test =bnb.predict(xtest)
            7  print('training accuracy_score: ', accuracy_score(ytrain, pred_train))
            8  print('training precision_score: ', precision_score(ytrain, pred_train))
            9  print('training confusion_matrix: \n', confusion_matrix(ytrain, pred_train))
           10
           11  #predicting testing data
           12  print('testing accuracy_score: ', accuracy_score(ytest, pred_test))
           13  print('testing precision_score: ', precision_score(ytest, pred_test))
           14  print('testing confusion_matrix: \n', confusion_matrix(ytest, pred_test))
```

```
training accuracy_score:  0.9835550181378476
training precision_score:  0.9932735426008968
training confusion_matrix:
 [[3624    3]
 [ 65  443]]
testing accuracy_score:  0.971953578336557
testing precision_score:  0.967741935483871
testing confusion_matrix:
 [[885    4]
 [ 25 120]]
```

```
In [49]:   1  # 'BoW_BernoulliNB'
           2  xtrain, xtest, ytrain, ytest =train_test_split(X_cv, y, test_size =0.2, random_state =42)
           3  bnb.fit(xtrain, ytrain)
           4  from sklearn.metrics import accuracy_score, precision_score, classification_report, confusion_matrix
           5  pred_train =bnb.predict(xtrain)
           6  pred_test =bnb.predict(xtest)
           7  print('training accuracy_score: ', accuracy_score(ytrain, pred_train))
           8  print('training precision_score: ', precision_score(ytrain, pred_train))
           9  print('training confusion_matrix: \n', confusion_matrix(ytrain, pred_train))
          10
          11  #predicting testing data
          12  print('testing accuracy_score: ', accuracy_score(ytest, pred_test))
          13  print('testing precision_score: ', precision_score(ytest, pred_test))
          14  print('testing confusion_matrix: \n', confusion_matrix(ytest, pred_test))
```

```
training accuracy_score:  0.9835550181378476
training precision_score:  0.9932735426008968
training confusion_matrix:
 [[3624    3]
 [  65  443]]
testing accuracy_score:  0.971953578336557
testing precision_score:  0.967741935483871
testing confusion_matrix:
 [[885   4]
 [ 25 120]]
```

- Since for a Business case **'precision'** is more important than **'accuracy'** for a Spam classifier, we will go ahead with **Tf-Idf-MultinomialNB.**
- Achieved Precision: 1.0
- Accuracy 0.9661508704061895

## 6. Model Improvement

```
Trying other ML classifer algorithm
```

```
In [50]:   1  from sklearn.linear_model import LogisticRegression
           2  from sklearn.svm import SVC
           3  from sklearn.naive_bayes import MultinomialNB
           4  from sklearn.tree import DecisionTreeClassifier
           5  from sklearn.neighbors import KNeighborsClassifier
           6  from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, BaggingClassifier, ExtraTreesClassifier,Gra
           7  from xgboost import XGBClassifier
```

```
In [51]:   1  lr=LogisticRegression()
           2  svc =SVC()
           3  mnb=MultinomialNB()
           4  dt= DecisionTreeClassifier()
           5  knn=KNeighborsClassifier()
           6  rf= RandomForestClassifier()
           7  ada= AdaBoostClassifier()
           8  bg=BaggingClassifier()
           9  etc= ExtraTreesClassifier()
          10  gb= GradientBoostingClassifier()
          11  xgb= XGBClassifier()
```

```
In [52]:   1  models ={'LR': lr, 'SVC' :svc, 'MultNB':mnb, 'DT':dt, "KNN":knn,"RF":rf, "ADA":ada,"BgC":bg,"ETC":etc,"GB":gb,"XGB":
```

```
In [53]:   1  xtrain, xtest, ytrain, ytest =train_test_split(X_tv, y, test_size =0.2, random_state =42)
           2
           3  def training_model(model, xtrain, xtest, ytrain, ytest):
           4      model.fit(xtrain,ytrain)
           5      pred= model.predict(xtest)
           6      accuracy =accuracy_score(ytest,pred)
           7      precision =precision_score(ytest,pred)
           8
           9      return accuracy, precision
```

```
In [54]:   1  accuracy=[]
           2  precision =[]
           3  for i in models:
           4      a,p=training_model(models[i],xtrain, xtest, ytrain, ytest)
           5      accuracy.append(a)
           6      precision.append(p)
```

```
In [55]:    1  final =pd.DataFrame({'Classifier' : models.keys(), 'Accuracy': accuracy, 'Precision': precision})
            2  final.sort_values(by='Precision', ascending =False)
```

Out[55]:

|     | Classifier | Accuracy | Precision |
| --- | --- | --- | --- |
| 2   | MultNB | 0.966151 | 1.000000 |
| 4   | KNN | 0.890716 | 1.000000 |
| 8   | ETC | 0.979691 | 0.992063 |
| 5   | RF | 0.977756 | 0.991935 |
| 1   | SVC | 0.970986 | 0.991453 |
| 9   | GB | 0.965184 | 0.990991 |
| 0   | LR | 0.954545 | 0.980392 |
| 10  | XGB | 0.972921 | 0.953488 |
| 6   | ADA | 0.974855 | 0.940741 |
| 7   | BgC | 0.967118 | 0.930233 |
| 3   | DT | 0.970019 | 0.895833 |

### Hypertuning the best model

```
In [56]:    1  from sklearn.model_selection import GridSearchCV
```

```
In [57]:    1  from sklearn.pipeline import Pipeline
            2  from sklearn.metrics import make_scorer
```

```
In [58]:    1  # hypertuning the TF-IDF -> MNB model
            2
            3  # Define a pipeline
            4  pipeline = Pipeline([
            5      ('tfidf', TfidfVectorizer()),
            6      ('nb', MultinomialNB())
            7  ])
            8
            9  # Define the parameter grid
           10  param_grid = {
           11      'tfidf__max_df': [0.5, 0.75,1.0],
           12      'tfidf__min_df': [1, 2, 5],
           13      'tfidf__ngram_range': [(1, 1), (1, 2)],
           14      'nb__alpha': [0.01, 0.1, 1, 10]
           15  }
           16
           17  # 'tfidf__max_df': Maximum document frequency (proportion of documents in which a term appears).
           18  # 'tfidf__min_df': Minimum document frequency (minimum number of documents a term must be in).
           19  # 'tfidf__ngram_range': The lower and upper boundary of the range of n-values for different n-grams to be extracted.
           20  # 'nb__alpha': Smoothing parameter for the Naive Bayes classifier.
           21
           22  # Define the scoring dictionary
           23  scoring = {
           24      'accuracy': make_scorer(accuracy_score),
           25      'precision': make_scorer(precision_score)
           26  }
           27  # Perform GridSearchCV
           28  gv = GridSearchCV(pipeline, param_grid, cv=5, n_jobs=-1, verbose=1, scoring = scoring, refit='accuracy')
           29
           30  # Fit the model
           31  xtrain, xtest, ytrain, ytest =train_test_split(df.std_text.astype('str'), y, test_size =0.2, random_state =42)
           32  gv.fit(xtrain, ytrain)
           33
           34  # Print best parameters and best score
           35  print("Best parameters found: ", gv.best_params_)
           36  print("Best accuracy_score: ", gv.best_score_)
```

```
Fitting 5 folds for each of 72 candidates, totalling 360 fits
Best parameters found:  {'nb__alpha': 0.1, 'tfidf__max_df': 0.5, 'tfidf__min_df': 2, 'tfidf__ngram_range': (1, 1)}
Best accuracy_score:  0.9842805320435308
```

```
In [59]:    1  # precision, accuracy
            2  precision_score(ytest,gv.predict(xtest)), accuracy_score(ytest,gv.predict(xtest))
```

Out[59]: (0.9692307692307692, 0.9777562862669246)

- the accuracy has improved from 0.0.9661508704061895 to 0.9777562862669246 , but precison has dropped from 1.0 to 0.9692307692307692
- this tradeoff is not desirable as precision holds greater significance for business use in a spam classifier
- so we will keep things unchanged an procced with the **TF-IDF-> MNB** as our final model

```
In [60]:  ▶  1  X_tv= tv.fit_transform(df.std_text).toarray()
             2  xtrain, xtest, ytrain, ytest =train_test_split(X_tv, y, test_size =0.2, random_state =42)
             3  mnb.fit(xtrain, ytrain)
             4  pred_train =mnb.predict(xtrain)
             5  pred_test =mnb.predict(xtest)
             6  print('testing accuracy_score: ', accuracy_score(ytest, pred_test))
             7  print('testing precision_score: ', precision_score(ytest, pred_test))
             8  print('testing confusion_matrix: \n', confusion_matrix(ytest, pred_test))
```

```
testing accuracy_score:  0.9661508704061895
testing precision_score:  1.0
testing confusion_matrix:
 [[889   0]
 [ 35 110]]
```

## 7. Pickle

```
In [61]:  ▶  1  import pickle
             2  pickle.dump(mnb, open('model.pkl', 'wb'))
             3  pickle.dump(tv, open('vectorizer.pkl', 'wb'))
             4  pickle.dump(standardize, open('standardize.pkl','wb'))
```

```
In [62]:  ▶  1  xtest[9].shape
```

Out[62]: (7059,)

```
In [63]:  ▶  1  xtest[0].reshape(1, -1).shape
```

Out[63]: (1, 7059)

```
In [64]:  ▶  1  mnb.predict(xtest[5].reshape(1, -1)), ytest.iloc[5]
```

Out[64]: (array([1], dtype=int64), 1)

```
In [65]:  ▶  1  mnb.predict(xtest[5].reshape(1, -1))[0]
```

Out[65]: 1

```
In [ ]:  ▶  1
```

```
In [66]:  ▶  1  t= """We're happy to inform you that you're entitled to a refund for overpayment on your AMEX account. Click on this
```

```
In [67]:  ▶  1  standardize(t)
```

Out[67]: 'happi inform entitl refund overpay amex account click link link claim refund'

```
In [68]:  ▶  1  tv.transform([standardize(t)]).shape
```

Out[68]: (1, 7059)

```
In [69]:  ▶  1  mnb.predict(tv.transform([standardize(t)]))
```

Out[69]: array([0], dtype=int64)

```
In [70]:  ▶  1  mnb.predict_proba(tv.transform([standardize(t)])).max()
```

Out[70]: 0.6703360228456645