# Guided Project Report

## Applications in Natural Language Processing

Name: Shruti Verma
Course: AI and ML
(Batch 4)
Duration: 10 months

Problem Statement: Using BoW and NLTK for processing, implement a simple spam filter that marks all the spam texts as dangerous.

## Prerequisites

What things you need to install the software and how to install them:

Python 3.8 or higher versions This setup requires that your machine has latest version of python. The following url https://www.python.org/downloads/ can be referred to download python. Once you have python downloaded and installed, you will need to setup PATH variables (if you want to run python program directly, detail instructions are below in how to run software section). To do that check this: https://www.pythoncentral.io/add-python-to-path-python-is-not- recognized-as-an-internal-or-external- command/. Setting up PATH variable is optional as you can also run program without it and more instruction are given below on this topic.

Second and easier option is to download anaconda and use its anaconda prompt to run the commands. To install anaconda check this url https://www.anaconda.com/download/ You will also need to download and install below 3 packages after you install either python or anaconda from the steps above Sklearn (scikit-learn) numpy scipy if you have chosen to install python 3.8 then run below commands in command prompt/terminal to install these packages pip install -U scikit-learn pip install numpy pip install scipy if you have chosen to install anaconda then run below commands in anaconda prompt to install these packages conda install -c scikit-learn conda install -c anaconda numpy conda install -c anaconda scipy . Install nltk.

## Video Link

https://drive.google.com/file/d/1yF7Gt2BxeM40db4m_wHbI4vayuF9v7cH/view?usp=sharing

## Dataset used

Any paragraph

Education is a process of learning through which spam we acquire knowledge.It enlightens, empowers, and creates spam a positive development.

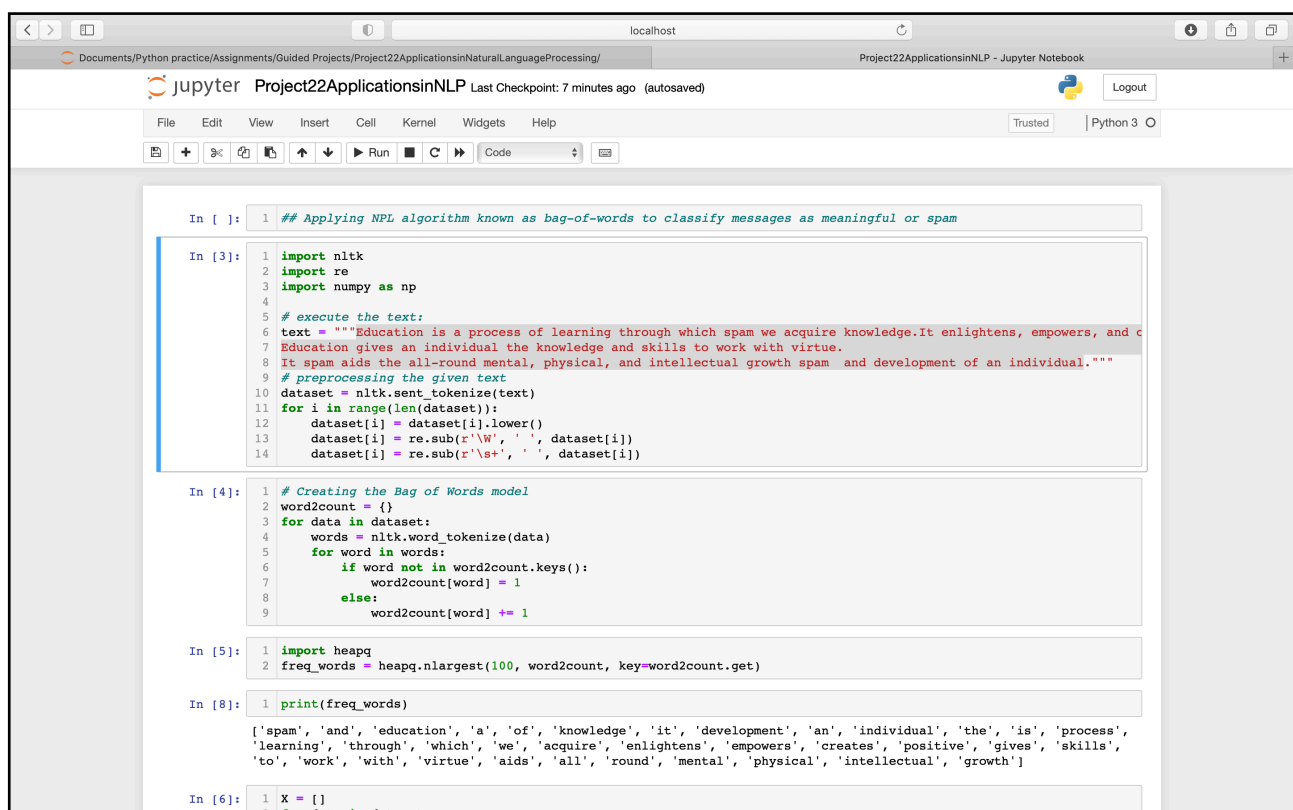Education gives an individual the knowledge and skills to work with virtue.

It spam aids the all-round mental, physical, and intellectual growth spam and development of an individual

## Method used for detection

- • Data reading
- • Preprocessing
- • Classification as spam

Importing the libraries and capturing images:

Importing necessary libraries and reading the store data

# Creating the BoW model and checking the frequency of words

```python
5  # execute the text:
6  text = """Education is a process of learning through which spam we acquire knowledge.It enlightens, empowers, and c
7  Education gives an individual the knowledge and skills to work with virtue.
8  It spam aids the all-round mental, physical, and intellectual growth spam  and development of an individual."""
9  # preprocessing the given text
10 dataset = nltk.sent_tokenize(text)
11 for i in range(len(dataset)):
12     dataset[i] = dataset[i].lower()
13     dataset[i] = re.sub(r'\W', ' ', dataset[i])
14     dataset[i] = re.sub(r'\s+', ' ', dataset[i])
```

```python
In [4]:
1  # Creating the Bag of Words model
2  word2count = {}
3  for data in dataset:
4      words = nltk.word_tokenize(data)
5      for word in words:
6          if word not in word2count.keys():
7              word2count[word] = 1
8          else:
9              word2count[word] += 1
```
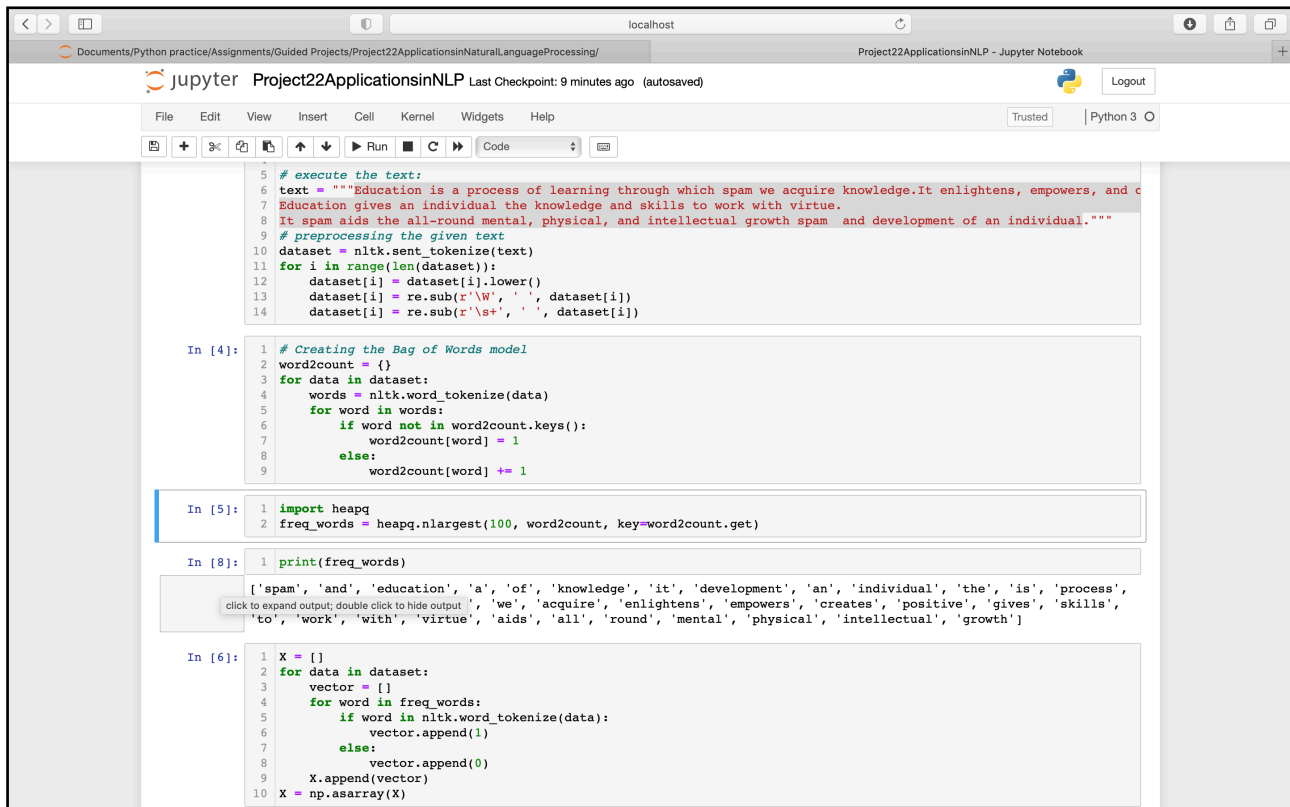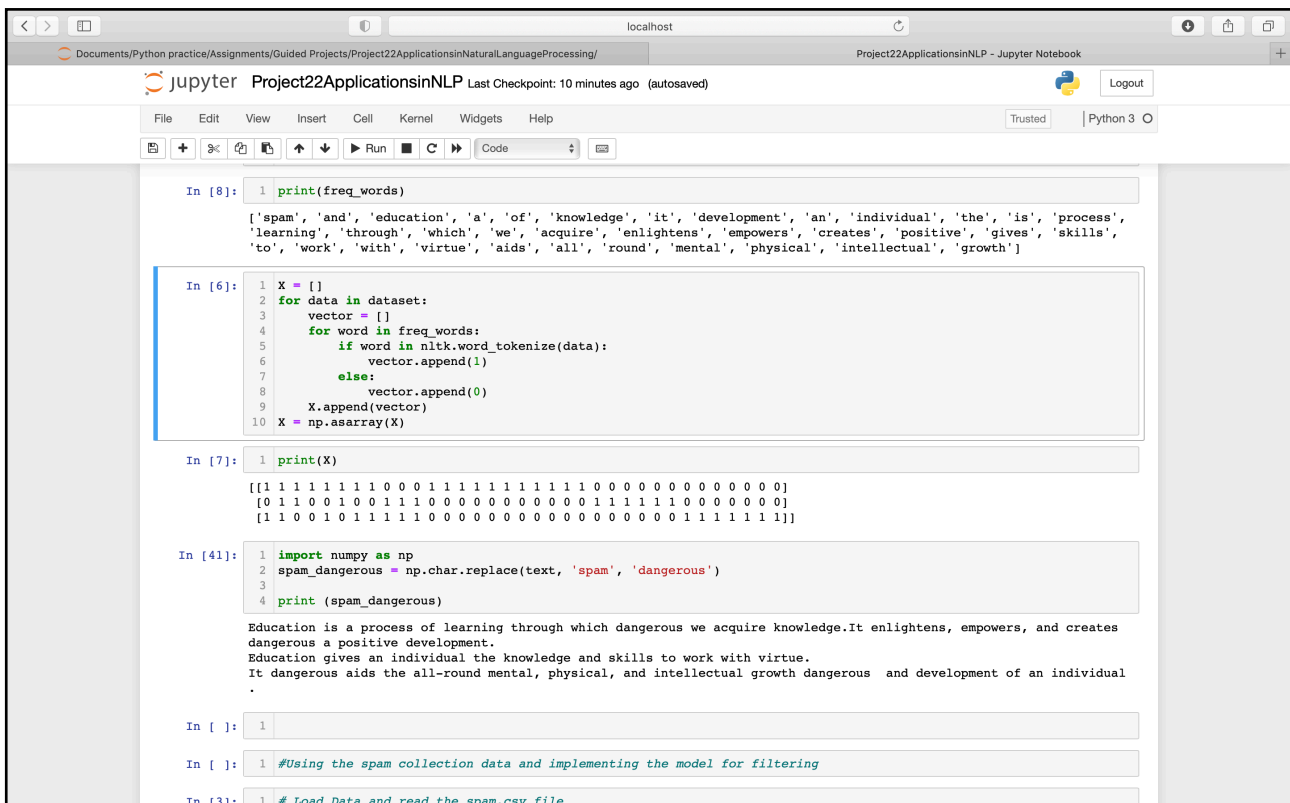
```python
In [5]:
1  import heapq
2  freq_words = heapq.nlargest(100, word2count, key=word2count.get)
```

```python
In [8]:
1  print(freq_words)
```

```
['spam', 'and', 'education', 'a', 'of', 'knowledge', 'it', 'development', 'an', 'individual', 'the', 'is', 'process',
click to expand output; double click to hide output  ', 'we', 'acquire', 'enlightens', 'empowers', 'creates', 'positive', 'gives', 'skills',
'to', 'work', 'with', 'virtue', 'aids', 'all', 'round', 'mental', 'physical', 'intellectual', 'growth']
```

```python
In [6]:
1  X = []
2  for data in dataset:
3      vector = []
4      for word in freq_words:
5          if word in nltk.word_tokenize(data):
6              vector.append(1)
7          else:
8              vector.append(0)
9      X.append(vector)
10 X = np.asarray(X)
```

# Replacing the spam word with dangerous

```python
In [8]:
1  print(freq_words)
```

```
['spam', 'and', 'education', 'a', 'of', 'knowledge', 'it', 'development', 'an', 'individual', 'the', 'is', 'process',
'learning', 'through', 'which', 'we', 'acquire', 'enlightens', 'empowers', 'creates', 'positive', 'gives', 'skills',
'to', 'work', 'with', 'virtue', 'aids', 'all', 'round', 'mental', 'physical', 'intellectual', 'growth']
```

```python
In [6]:
1  X = []
2  for data in dataset:
3      vector = []
4      for word in freq_words:
5          if word in nltk.word_tokenize(data):
6              vector.append(1)
7          else:
8              vector.append(0)
9      X.append(vector)
10 X = np.asarray(X)
```

```python
In [7]:
1  print(X)
```

```
[[1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0]
 [0 1 1 0 0 1 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 0 0 0 0 0 0]
 [1 1 0 0 1 0 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1]]
```

```python
In [41]:
1  import numpy as np
2  spam_dangerous = np.char.replace(text, 'spam', 'dangerous')
3
4  print (spam_dangerous)
```

```
Education is a process of learning through which dangerous we acquire knowledge.It enlightens, empowers, and creates
dangerous a positive development.
Education gives an individual the knowledge and skills to work with virtue.
It dangerous aids the all-round mental, physical, and intellectual growth dangerous  and development of an individual
.
```

```python
In [ ]:
1
```

```python
In [ ]:
1  #Using the spam collection data and implementing the model for filtering
```

```python
In [3]:
1  # Load Data and read the spam.csv file
```

# Load and read spam.csv

Documents/Python practice/Assignments/Guided Projects/Project22ApplicationsinNaturalLanguageProcessing/     Project22ApplicationsinNLP - Jupyter Notebook

**jupyter** Project22ApplicationsinNLP   Last Checkpoint: 13 minutes ago   (autosaved)     Logout

File   Edit   View   Insert   Cell   Kernel   Widgets   Help     Trusted   |   Python 3 ○

```python
In [ ]:  1  #Using the spam collection data and implementing the model for filtering
```

```python
In [3]:  1  # Load Data and read the spam.csv file
```

```python
In [23]:  1  import pandas as pd
          2  import numpy as np
          3  data = pd.read_csv("spam.csv", encoding = "latin-1")
          4  data = data[['v1', 'v2']]
          5  data = data.rename(columns = {'v1': 'label', 'v2': 'text'})
```

```python
In [24]:  1  #data['label']=data['label'].replace('ham', 0)
          2  #data['label']=data['label'].replace('spam', 1)
```

```python
In [25]:  1  #data.head(5)
```

Out[25]:

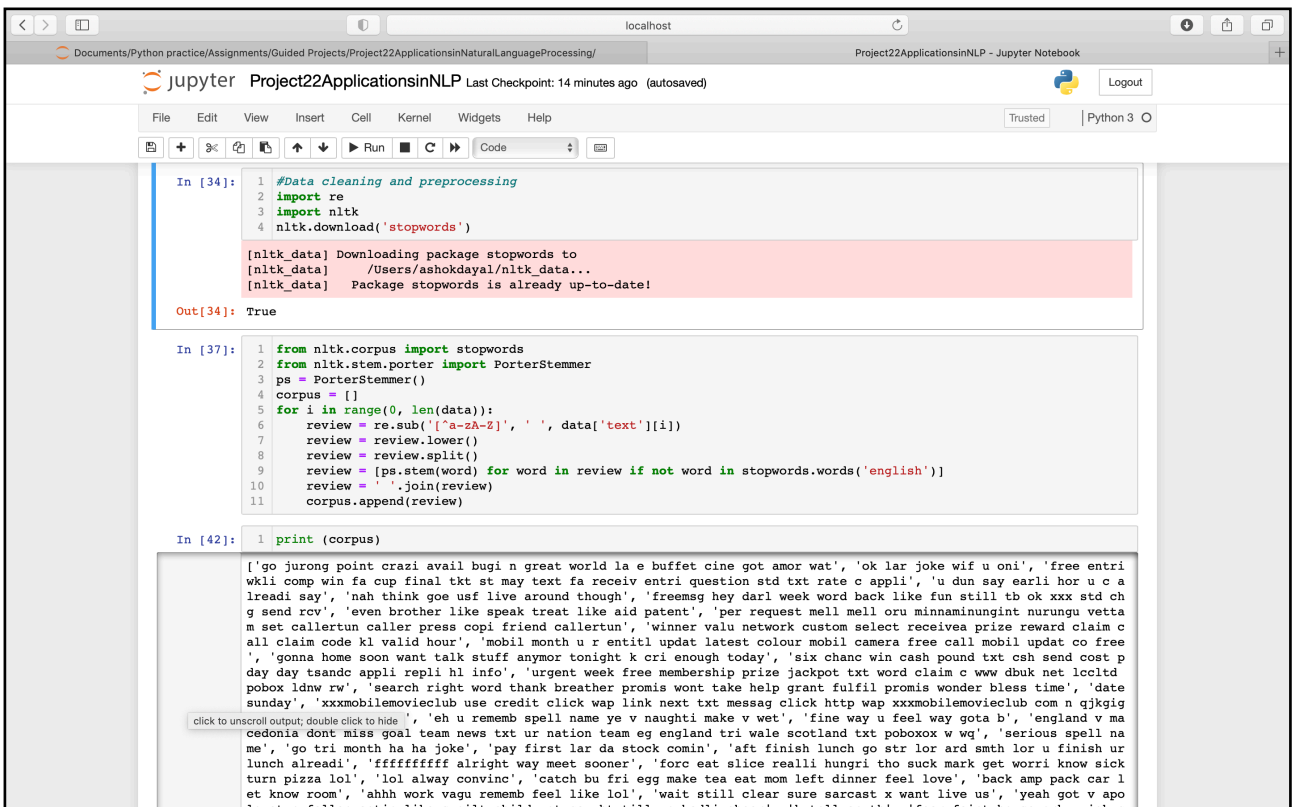|   | label | text |
|---|-------|------|
| 0 | 0 | Go until jurong point, crazy.. Available only ... |
| 1 | 0 | Ok lar... Joking wif u oni... |
| 2 | 1 | Free entry in 2 a wkly comp to win FA Cup fina... |
|   | *click to expand output; double click to hide output* | then say... |
| 4 | 0 | Nah I don't think he goes to usf, he lives aro... |

```python
In [9]:  1  # Data Pre-processing - maining removing stopwords, tokenization,lower case
```

```python
In [34]:  1  #Data cleaning and preprocessing
          2  import re
          3  import nltk
          4  nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/ashokdayal/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

Out[34]: True

# Data cleaning and preprocessing

```python
In [34]:  1  #Data cleaning and preprocessing
          2  import re
          3  import nltk
          4  nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/ashokdayal/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

Out[34]: True

```python
In [37]:   1  from nltk.corpus import stopwords
           2  from nltk.stem.porter import PorterStemmer
           3  ps = PorterStemmer()
           4  corpus = []
           5  for i in range(0, len(data)):
           6      review = re.sub('[^a-zA-Z]', ' ', data['text'][i])
           7      review = review.lower()
           8      review = review.split()
           9      review = [ps.stem(word) for word in review if not word in stopwords.words('english')]
          10      review = ' '.join(review)
          11      corpus.append(review)
```
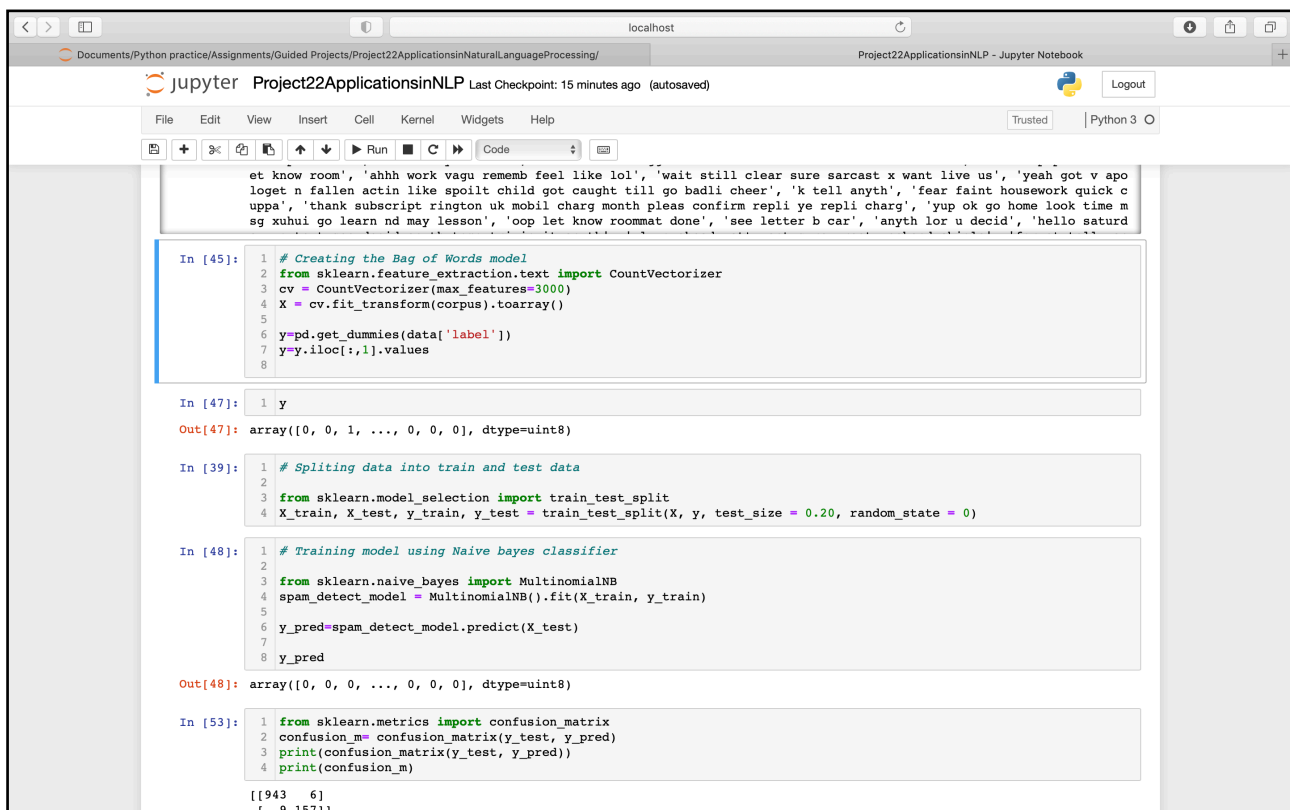
```python
In [42]:  1  print (corpus)
```

```
['go jurong point crazi avail bugi n great world la e buffet cine got amor wat', 'ok lar joke wif u oni', 'free entri
wkli comp win fa cup final tkt st may text fa receiv entri question std txt rate c appli', 'u dun say earli hor u c a
lreadi say', 'nah think goe usf live around though', 'freemsg hey darl week word back like fun still tb ok xxx std ch
g send rcv', 'even brother like speak treat like aid patent', 'per request mell mell oru minnaminungint nurungu vetta
m set callertun caller press copi friend callertun', 'winner valu network custom select receivea prize reward claim c
all claim code kl valid hour', 'mobil month u r entitl updat latest colour mobil camera free call mobil updat co free
', 'gonna home soon want talk stuff anymor tonight k cri enough today', 'six chanc win cash pound txt csh send cost p
day day tsandc appli repli hl info', 'urgent week free membership prize jackpot txt word claim c www dbuk net lccltd
pobox ldnw rw', 'search right word thank breather promis wont take help grant fulfil promis wonder bless time', 'date
sunday', 'xxxmobilemovieclub use credit click wap link next txt messag click http wap xxxmobilemovieclub com n qjkgig
 ', 'eh u rememb spell name ye v naughti make v wet', 'fine way u feel way gota b', 'england v ma
cedonia dont miss goal team news txt ur nation team eg england tri wale scotland txt poboxox w wq', 'serious spell na
me', 'go tri month ha ha joke', 'pay first lar da stock comin', 'aft finish lunch go str lor ard smth lor u finish ur
lunch alreadi', 'fffffffffff alright way meet sooner', 'forc eat slice realli hungri tho suck mark get worri know sick
turn pizza lol', 'lol alway convinc', 'catch bu fri egg make tea eat mom left dinner feel love', 'back amp pack car l
et know room', 'ahhh work vagu rememb feel like lol', 'wait still clear sure sarcast x want live us', 'yeah got v apo
loget n fallen actin like spoilt child got caught till go badli cheer', 'k tell anyth', 'fear faint housework quick c
```

## Applying BoW model and splitting data into train and test

```
et know room', 'ahhh work vagu rememb feel like lol', 'wait still clear sure sarcast x want live us', 'yeah got v apo
loget n fallen actin like spoilt child got caught till go badli cheer', 'k tell anyth', 'fear faint housework quick c
uppa', 'thank subscript rington uk mobil charg month pleas confirm repli ye repli charg', 'yup ok go home look time m
sg xuhui go learn nd may lesson', 'oop let know roommat done', 'see letter b car', 'anyth lor u decid', 'hello saturd
```

```python
In [45]:  1  # Creating the Bag of Words model
          2  from sklearn.feature_extraction.text import CountVectorizer
          3  cv = CountVectorizer(max_features=3000)
          4  X = cv.fit_transform(corpus).toarray()
          5
          6  y=pd.get_dummies(data['label'])
          7  y=y.iloc[:,1].values
          8
```
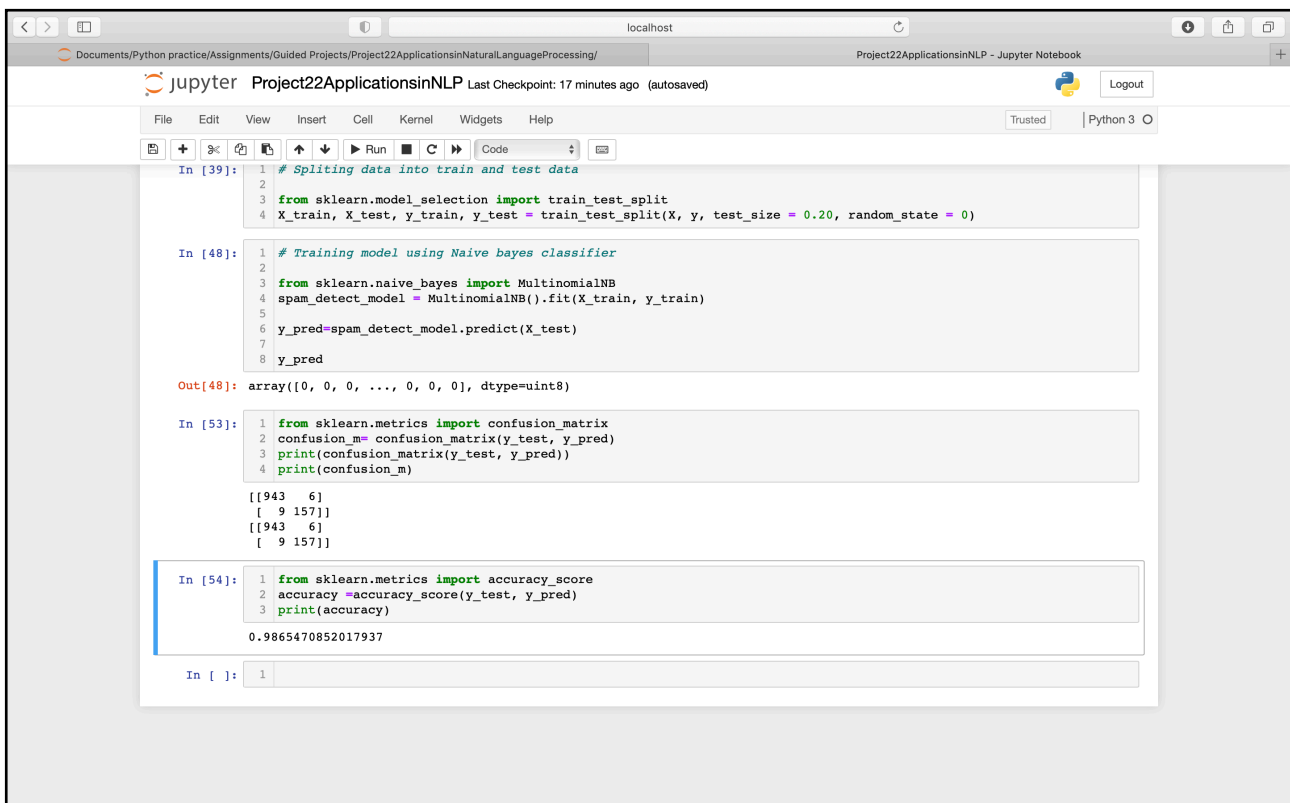
```python
In [47]:  1  y
```
Out[47]: array([0, 0, 1, ..., 0, 0, 0], dtype=uint8)

```python
In [39]:  1  # Spliting data into train and test data
          2
          3  from sklearn.model_selection import train_test_split
          4  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20, random_state = 0)
```

```python
In [48]:  1  # Training model using Naive bayes classifier
          2
          3  from sklearn.naive_bayes import MultinomialNB
          4  spam_detect_model = MultinomialNB().fit(X_train, y_train)
          5
          6  y_pred=spam_detect_model.predict(X_test)
          7
          8  y_pred
```
Out[48]: array([0, 0, 0, ..., 0, 0, 0], dtype=uint8)

```python
In [53]:  1  from sklearn.metrics import confusion_matrix
          2  confusion_m= confusion_matrix(y_test, y_pred)
          3  print(confusion_matrix(y_test, y_pred))
          4  print(confusion_m)
```
```
[[943   6]
 [  9 157]]
```

## Using Classifier and checking the accuracy

```python
In [39]:  1  # Spliting data into train and test data
          2
          3  from sklearn.model_selection import train_test_split
          4  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20, random_state = 0)
```

```python
In [48]:  1  # Training model using Naive bayes classifier
          2
          3  from sklearn.naive_bayes import MultinomialNB
          4  spam_detect_model = MultinomialNB().fit(X_train, y_train)
          5
          6  y_pred=spam_detect_model.predict(X_test)
          7
          8  y_pred
```
Out[48]: array([0, 0, 0, ..., 0, 0, 0], dtype=uint8)

```python
In [53]:  1  from sklearn.metrics import confusion_matrix
          2  confusion_m= confusion_matrix(y_test, y_pred)
          3  print(confusion_matrix(y_test, y_pred))
          4  print(confusion_m)
```
```
[[943   6]
 [  9 157]]
[[943   6]
 [  9 157]]
```

```python
In [54]:  1  from sklearn.metrics import accuracy_score
          2  accuracy =accuracy_score(y_test, y_pred)
          3  print(accuracy)
```
```
0.9865470852017937
```

```python
In [ ]:   1
```

Documents/Python practice/Assignments/Resource dumps/Unsupervised Learn... | Project13 AssociationRuleMiningMarketBasketAnalysis - Jupyter Notebook | Association_Rule_Mining-Apriori_Algorithm-Tutorial - Jupyter Notebook

**jupyter** Project13 AssociationRuleMiningMarketBasketAnalysis Last Checkpoint: a few seconds ago (autosaved)

Logout

File | Edit | View | Insert | Cell | Kernel | Widgets | Help

Trusted | Python 3

```
In [31]: 1 rules[ (rules['lift'] > 1) & (rules['confidence'] > 0.5)]
```

Out[31]:

| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | leverage | conviction |
|---|---|---|---|---|---|---|---|---|---|
| 0 | (0) | (1) | 0.001600 | 0.002000 | 0.001200 | 0.750000 | 375.000000 | 0.001197 | 3.992000 |
| 1 | (1) | (0) | 0.002000 | 0.001600 | 0.001200 | 0.600000 | 375.000000 | 0.001197 | 2.496000 |
| 2 | (0) | (5) | 0.001600 | 0.001600 | 0.001467 | 0.916667 | 572.916667 | 0.001464 | 11.980800 |
| 3 | (5) | (0) | 0.001600 | 0.001600 | 0.001467 | 0.916667 | 572.916667 | 0.001464 | 11.980800 |
| 4 | (0) | (19) | 0.001600 | 0.001200 | 0.001200 | 0.750000 | 625.000000 | 0.001198 | 3.995200 |
| 5 | (19) | (0) | 0.001200 | 0.001600 | 0.001200 | 1.000000 | 625.000000 | 0.001198 | inf |
| 6 | (1) | (5) | 0.002000 | 0.001600 | 0.001067 | 0.533333 | 333.333333 | 0.001063 | 2.139429 |
| 7 | (5) | (1) | 0.001600 | 0.002000 | 0.001067 | 0.666667 | 333.333333 | 0.001063 | 2.994000 |
| 8 | (1) | (15) | 0.002000 | 0.001467 | 0.001067 | 0.533333 | 363.636364 | 0.001064 | 2.139714 |
| 9 | (15) | (1) | 0.001467 | 0.002000 | 0.001067 | 0.727273 | 363.636364 | 0.001064 | 3.659333 |
| 10 | (1) | (18) | 0.002000 | 0.001333 | 0.001067 | 0.533333 | 400.000000 | 0.001064 | 2.140000 |
| 11 | (18) | (1) | 0.001333 | 0.002000 | 0.001067 | 0.800000 | 400.000000 | 0.001064 | 4.990000 |
| 12 | (1) | (19) | 0.002000 | 0.001200 | 0.001067 | 0.533333 | 444.444444 | 0.001064 | 2.140286 |
| 13 | (19) | (1) | 0.001200 | 0.002000 | 0.001067 | 0.888889 | 444.444444 | 0.001064 | 8.982000 |
| 14 | (19) | (5) | 0.001200 | 0.001600 | 0.001067 | 0.888889 | 555.555556 | 0.001065 | 8.985600 |
| 15 | (5) | (19) | 0.001600 | 0.001200 | 0.001067 | 0.666667 | 555.555556 | 0.001065 | 2.996400 |
| 16 | (0, 1) | (5) | 0.001200 | 0.001600 | 0.001067 | 0.888889 | 555.555556 | 0.001065 | 8.985600 |
| 17 | (0, 5) | (1) | 0.001467 | 0.002000 | 0.001067 | 0.727273 | 363.636364 | 0.001064 | 3.659333 |
| 18 | (1, 5) | (0) | 0.001067 | 0.001600 | 0.001067 | 1.000000 | 625.000000 | 0.001065 | inf |
| 19 | (0) | (1, 5) | 0.001600 | 0.001067 | 0.001067 | 0.666667 | 625.000000 | 0.001065 | 2.996800 |
| 20 | (1) | (0, 5) | 0.002000 | 0.001467 | 0.001067 | 0.533333 | 363.636364 | 0.001064 | 2.139714 |
| 21 | (5) | (0, 1) | 0.001600 | 0.001200 | 0.001067 | 0.666667 | 555.555556 | 0.001065 | 2.996400 |
| 22 | (0, 1) | (19) | 0.001200 | 0.001200 | 0.001067 | 0.888889 | 740.740741 | 0.001065 | 8.989200 |
| 23 | (0, 19) | (1) | 0.001200 | 0.002000 | 0.001067 | 0.888889 | 444.444444 | 0.001064 | 8.982000 |