

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from IPython.display import display
from sklearn.preprocessing import LabelEncoder, MinMaxScaler, StandardScaler
import pickle
```

```
In [2]: import nltk
import re
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
```

```
In [3]: df = pd.read_csv('SMSSpamCollection.csv', sep='\t', names=['label', 'message'])
df
```

```
Out[3]:
```

	label	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...
...
5567	spam	This is the 2nd time we have tried 2 contact u...
5568	ham	Will ü b going to esplanade fr home?
5569	ham	Pity, * was in mood for that. So...any other s...
5570	ham	The guy did some bitching but I acted like i'd...
5571	ham	Rofl. Its true to its name

5572 rows × 2 columns

```
In [4]: # Cleaning the messages
corpus = []
ps = PorterStemmer()

for i in range(0,df.shape[0]):

    # Cleaning special character from the message
    message = re.sub(pattern='[^a-zA-Z]', repl=' ', string=df.message[i])

    # Converting the entire message into lower case
    message = message.lower()

    # Tokenizing the review by words
    words = message.split()

    # Removing the stop words
    words = [word for word in words if word not in set(stopwords.words('english'))]

    # Stemming the words
    words = [ps.stem(word) for word in words]

    # Joining the stemmed words
    message = ' '.join(words)

    # Building a corpus of messages
    corpus.append(message)
```

```
In [5]: # Creating the Bag of Words model
from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=2500)
X = cv.fit_transform(corpus).toarray()

# Extracting dependent variable from the dataset
y = pd.get_dummies(df['label'])
y = y.iloc[:, 1].values
```

```
In [6]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=0)

# Fitting Naive Bayes to the Training set
from sklearn.naive_bayes import MultinomialNB
classifier = MultinomialNB(alpha=0.3)
classifier.fit(X_train, y_train)
```

```
Out[6]: MultinomialNB(alpha=0.3, class_prior=None, fit_prior=True)
```

```
In [7]: # Creating a pickle file for the CountVectorizer
file = open('cv-transform.pkl', 'wb')
pickle.dump(cv,file)

# Creating a pickle file for the Multinomial Naive Bayes model
filename = open('spam-sms-mnb-model.pkl', 'wb')
pickle.dump(classifier,filename )
```

```
In [8]: file.close()
filename.close()
```