

# Natural Language Processing

## Importing the libraries

```
In [0]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

## Importing the dataset

```
In [0]: dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting
= 3)
```

## Cleaning the texts

```
In [3]: import re
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
corpus = []
for i in range(0, 1000):
    review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][i])
    review = review.lower()
    review = review.split()
    ps = PorterStemmer()
    review = [ps.stem(word) for word in review if not word in set(stopwor
ds.words('english'))]
    review = ' '.join(review)
    corpus.append(review)
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
```

```
[nltk_data] Package stopwords is already up-to-date!
```

## Creating the Bag of Words model

```
In [0]: from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features = 1500)
X = cv.fit_transform(corpus).toarray()
y = dataset.iloc[:, 1].values
```

## Splitting the dataset into the Training set and Test set

```
In [0]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2
0, random_state = 0)
```

## Training the Naive Bayes model on the Training set

```
In [6]: from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(X_train, y_train)
```

```
Out[6]: GaussianNB(priors=None, var_smoothing=1e-09)
```

## Predicting the Test set results

```
In [0]: y_pred = classifier.predict(X_test)
```

## Making the Confusion Matrix

```
In [8]: from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
print(cm)
```

```
[[55 42]
 [12 91]]
```