# Spam Email Classifier

SUBMITTED TO:

SUBMITTED BY:

MR. NITIN

YASH RAJ

101511069

# NLTK used for

- Tokenizing words
- Lemmatizing words
- Removing stop words from the training data
- Using Naïve  Bayes Classifier for classification purposes

# Code:

```python
import nltk
import os
import random
from nltk import word_tokenize, WordNetLemmatizer
from nltk.corpus import stopwords
from nltk import NaiveBayesClassifier, classify

stoplist = stopwords.words('english')

def get_data(folder):
    a = []
    files = os.listdir(folder)
    for file in files:
        f = open(folder + file, 'r', encoding = "ISO-8859-1")  ## cause there are some latin words also
        a.append(f.read())
    f.close()
    return a
```

Necessary imports

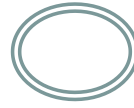Creating list of stop words from English language

Reading all the files

`

```python
def preprocess(sentence):
lemmatizer = WordNetLemmatizer()
    return [lemmatizer.lemmatize(word.lower()) for word in word_tokenize(sentence)]

def getfeatures(text):   #creating feature list
    return {word: True for word in preprocess(text) if not word in stoplist}

def print_accuracy(train_set, test_set, classifier): #finding accuracy of the training and testing data
    print ('Accuracy on the training set = ' + str(classify.accuracy(classifier, train_set)))
    print ('Accuracy of the test set = ' + str(classify.accuracy(classifier, test_set)))
    classifier.show_most_informative_features(10)
```

takes sentence as parameter and lemmatize it word by word

```python
def train(features, samples_proportion):
    train_size = int(len(features) * samples_proportion)      #finds index of data at 80%

    train_set, test_set = features[:train_size], features[train_size:]
    print ('Training set size = ' + str(len(train_set)) + ' emails')
    print ('Test set size = ' + str(len(test_set)) + ' emails')      #defining training and #testing data

    classifier = NaiveBayesClassifier.train(train_set)
    return train_set, test_set, classifier      #training using naïve Bayes #classifier
```

```python
if __name__ == "__main__":
    spam = get_data('enron1/spam/')        #reads spam emails
    ham  = get_data('enron1/ham/')                  #reads ham emails

    allemails  = [(email, 'spam') for email in spam]#creating a list of tuples having email
    allemails += [(email, 'ham') for email in ham]    with  its label
    random.shuffle(allemails)                           #shuffling the list to randomly distribute
                                                          training and testing data
    print (str(len(allemails)) + ' emails')

    allfeatures = [(getfeatures(email), label) for (email, label) in allemails]   #learning features
    print (str(len(allfeatures)) + ' feature sets')                                    from the
data

    train_set, test_set, classifier = train(allfeatures, 0.8)    #defining size of training ,testing
data
```

# Result