# SNA MINI PROJECT

Aim: Sentiment Analysis

Code:

I am going to install nltk package for this.We will do this by using command:

pip install nltk==3.3

```
C:\Users\Admin>pip install nltk==3.3
Collecting nltk==3.3
  Downloading nltk-3.3.0.zip (1.4 MB)
     ---------------------------------------- 1.4/1.4 MB 2.5 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Collecting six
  Downloading six-1.16.0-py2.py3-none-any.whl (11 kB)
Using legacy 'setup.py install' for nltk, since package 'wheel' is not installed.
Installing collected packages: six, nltk
  Running setup.py install for nltk ... done
Successfully installed nltk-3.3 six-1.16.0
```

Now as nltk package is installed successfully. We will proceed to download twitter samples.

```
>>> import nltk
>>> nltk.download('twitter_samples')
[nltk_data] Downloading package twitter_samples to
[nltk_data]     C:\Users\Admin\AppData\Roaming\nltk_data...
[nltk_data]   Package twitter_samples is already up-to-date!
True
```

This will import 3 datasets from NLTK.

```
[1]:  1  from nltk.corpus import twitter_samples
```

Next we will be creating variables for positive , negative tweets and text where strings() method will print all the tweets as strings.

```python
from nltk.corpus import twitter_samples

positive_tweets = twitter_samples.strings('positive_tweets.json')
negative_tweets = twitter_samples.strings('negative_tweets.json')
text = twitter_samples.strings('tweets.20150430-223406.json')
```

Now we will be downloading resource punkt which is a pre-trained model that helps to tokenize words and sentences.

```
>>> nltk.download('punkt')
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\Admin\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping tokenizers\punkt.zip.
True
```

We will create an object that tokenizes positive_tweets.json dataset. The tokenized() method will return special characters in this case @, _ .

```python
1  from nltk.corpus import twitter_samples
2
3  positive_tweets = twitter_samples.strings('positive_tweets.json')
4  negative_tweets = twitter_samples.strings('negative_tweets.json')
5  text = twitter_samples.strings('tweets.20150430-223406.json')
6  tweet_tokens = twitter_samples.tokenized('positive_tweets.json')
7  print(tweet_tokens[0])
```

```
['#FollowFriday', '@France_Inte', '@PKuchly57', '@Milipol_Paris', 'for', 'being', 'top', 'engaged', 'members', 'in', 'my', 'community', 'this', 'week', ':)']
```

Now we will be adding another resource: 'wordnet' and 'average_perceptron_tagger'

```
>>> nltk.download('wordnet')
[nltk_data] Downloading package wordnet to
[nltk_data]     C:\Users\Admin\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
True
>>> nltk.download('averaged_perceptron_tagger')
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     C:\Users\Admin\AppData\Roaming\nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]       date!
True
>>>
```

Now we will determine context for each word in our text. In python pos_tag function is used for that.

```
1  from nltk.tag import pos_tag
2  from nltk.corpus import twitter_samples
3  tweet_tokens=twitter_samples.tokenized('positive_tweets.json')
4  print(pos_tag(tweet_tokens[0]))
```

```
[('#FollowFriday', 'JJ'), ('@France_Inte', 'NNP'), ('@PKuchly57', 'NNP'), ('@Milipol_Paris', 'NNP'), ('for', 'IN'), ('being',
'VBG'), ('top', 'JJ'), ('engaged', 'VBN'), ('members', 'NNS'), ('in', 'IN'), ('my', 'PRP$'), ('community', 'NN'), ('this', 'D
T'), ('week', 'NN'), (':)', 'NN')]
```

We will be creating dictionaries for tweets.

```
1  def get_tweets_for_model(cleaned_tokens_list):
2      for tweet_tokens in cleaned_tokens_list:
3          yield dict([token, True] for token in tweet_tokens)
4
5  positive_tokens_for_model = get_tweets_for_model(positive_cleaned_tokens_list)
6  negative_tokens_for_model = get_tweets_for_model(negative_cleaned_tokens_list)
```

This code attaches positive or negative label to each tweet.

```
1  import random
2
3  positive_dataset = [(tweet_dict, "Positive")
4                       for tweet_dict in positive_tokens_for_model]
5
6  negative_dataset = [(tweet_dict, "Negative")
7                       for tweet_dict in negative_tokens_for_model]
8
9  dataset = positive_dataset + negative_dataset
10
11 random.shuffle(dataset)
12
13 train_data = dataset[:7000]
14 test_data = dataset[7000:]
```

We will create a model by using NaiveBayesClassifier and will be testing it by accuracy() method.

```
81]:   1  from nltk import classify
       2  from nltk import NaiveBayesClassifier
       3  classifier = NaiveBayesClassifier.train(train_data)
       4
       5  print("Accuracy is:", classify.accuracy(classifier, test_data))
       6
       7  print(classifier.show_most_informative_features(10))
```

```
Accuracy is: 0.996
Most Informative Features
                      :) = True           Positi : Negati =   1011.9 : 1.0
                follower = True           Positi : Negati =     40.4 : 1.0
                followed = True           Negati : Positi =     22.9 : 1.0
                     sad = True           Negati : Positi =     19.4 : 1.0
                 welcome = True           Positi : Negati =     15.2 : 1.0
                  arrive = True           Positi : Negati =     14.9 : 1.0
              appreciate = True           Positi : Negati =     14.1 : 1.0
                   enjoy = True           Positi : Negati =     13.5 : 1.0
                     idk = True           Negati : Positi =     12.6 : 1.0
                   didnt = True           Negati : Positi =     11.9 : 1.0
None
```

Then we will use another resource vader_lexicon. We will get back a dictionary of different scores.

```
In [45]:   1  from nltk.sentiment import SentimentIntensityAnalyzer
           2  sia = SentimentIntensityAnalyzer()
           3  sia.polarity_scores("Wow, NLTK is really powerful!")

Out[45]: {'neg': 0.0, 'neu': 0.295, 'pos': 0.705, 'compound': 0.8012}
```

We will use polarity.scores() function to classify tweets.

```
 1  tweets = [t.replace("://", "//") for t in nltk.corpus.twitter_samples.strings]
 2  from random import shuffle
 3
 4  def is_positive(tweet: str) -> bool:
 5      """True if tweet has positive compound sentiment, False otherwise."""
 6      return sia.polarity_scores(tweet)["compound"] > 0
 7
 8  shuffle(tweets)
 9  for tweet in tweets[:10]:
10      print(">", is_positive(tweet), tweet)
```

> False RT @jameschappers: Election roars to life as real voters have their say... and leave Miliband stumbling #tomorrowspaper stoday https//t.co/…
> False hey there, internet - anyone had any troubles with ordering from #fixedgearfrenzy? Been 5 days now &amp; no dispatch no tice. Getting worried :(
> True RT @KennyFarq: My column in tomorrow's Scotsman challenges the morality of the SNP's decision to vote on English schools and hospitals. #bb…
> False RT @DickMackintosh: Nick Robinson talking out of his arse on #BBCNews as usual ... nasty little Tory shite ...
> False Election 2015 Nigel Farage: Live updates as Ukip leader is grilled by voters in Birmingham http//t.co/A0eDN3ZMHR
> True So @gwatsky had a fantastic show! And I already want to buy tickets to another concert.:D
> False RT @BarrySheerman: Are there any non-Tory leaning journalists on @BBCNews? Tonight's coverage a disgrace from a broadca ster funded by publi…
> True RT @Carra23: Thought Ed Milliband came out on top tonight but the best words tonight were from the audience numerous tim es "but that wasn't…
> True RT @Tommy_Colc: Financial Times come out in support of Tories claiming Miliband is "preoccupied w/ inequality". The man who wrote it http:/…
> True RT @blairmcdougall: Salmond on Sky encouraging SNP vote in Scotland Plaid vote in Wales &amp; Green vote in England. Not

Now we will check how model performs on random tweets from twitter.

```
In [83]:    1  from nltk.tokenize import word_tokenize
            2
            3  custom_tweet = "I ordered an item from there and iits quality was worst, never going to use app again."
            4
            5  custom_tokens = remove_noise(word_tokenize(custom_tweet))
            6
            7  print(classifier.classify(dict([token, True] for token in custom_tokens)))
```

Negative

```
In [84]:    1  from nltk.tokenize import word_tokenize
            2
            3  custom_tweet = "The product was pretty good i do recommend this."
            4
            5  custom_tokens = remove_noise(word_tokenize(custom_tweet))
            6
            7  print(classifier.classify(dict([token, True] for token in custom_tokens)))
```

Positive

We will use another resource movie_reviews.

```
>>> nltk.download('movie_reviews')
[nltk_data] Downloading package movie_reviews to
[nltk_data]     C:\Users\Admin\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\movie_reviews.zip.
True
>>>
```

```
 1  positive_review_ids = nltk.corpus.movie_reviews.fileids(categories=["pos"])
 2  negative_review_ids = nltk.corpus.movie_reviews.fileids(categories=["neg"])
 3  all_review_ids = positive_review_ids + negative_review_ids
```

```
 1  from statistics import mean
 2
 3  def is_positive(review_id: str) -> bool:
 4      """True if the average of all sentence compound scores is positive."""
 5      text = nltk.corpus.movie_reviews.raw(review_id)
 6      scores = [
 7          sia.polarity_scores(sentence)["compound"]
 8          for sentence in nltk.sent_tokenize(text)
 9      ]
10      return mean(scores) > 0
```

```
 1  shuffle(all_review_ids)
 2  correct = 0
 3  for review_id in all_review_ids:
 4      if is_positive(review_id):
 5          if review_id in positive_review_ids:
 6              correct += 1
 7      else:
 8          if review_id in negative_review_ids:
 9              correct += 1
10  print(F"{correct / len(all_review_ids):.2%} correct")
```

69.10% correct

We will be using train() method that involves splitting the feature set so that one portion can be used for training and other for evaluation.

```
In [66]:  1  train_count = len(features) // 4
          2  shuffle(features)
          3  classifier = nltk.NaiveBayesClassifier.train(features[:train_count])
          4  classifier.show_most_informative_features(10)

          Most Informative Features
                     wordcount = 3            pos : neg    =      8.1 : 1.0
                     wordcount = 4            pos : neg    =      6.0 : 1.0
                     wordcount = 2            pos : neg    =      4.1 : 1.0
                     wordcount = 1            pos : neg    =      2.0 : 1.0
                     wordcount = 0            neg : pos    =      1.8 : 1.0
                 mean_positive = 0.1437       pos : neg    =      1.0 : 1.0

In [67]:  1  nltk.classify.accuracy(classifier, features[train_count:])

Out[67]:  0.6553333333333333
```