# 4222-SURYA GROUP OF INSTITUTIONS

# VIKRAVANDI.

PREPARED BY,

S.MASILAMANI,

422221106304,

ECE- DEPT,

3$^{RD}$ YEAR.

# AI_PHASE 3

## PREPROCESSING

```
import nltk
from nltk.corpus import twitter_samples
import matplotlib.pyplot as plt
import random
```
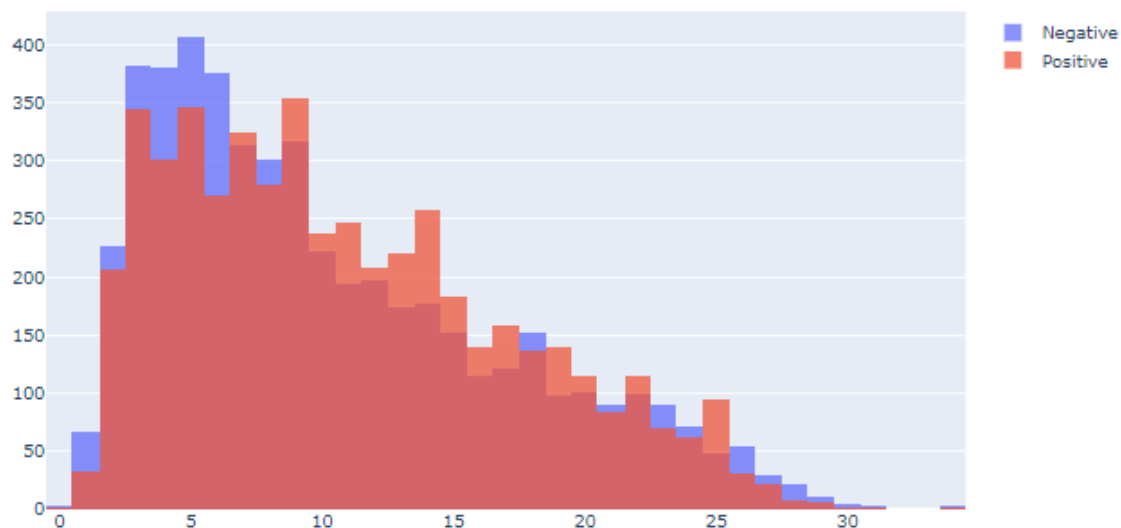
## NLTK Twitter Dataset

```
nltk.download('twitter_samples')
[nltk_data] Downloading package twitter_samples to
[nltk_data]     /usr/share/nltk_data...
[nltk_data]   Package twitter_samples is already up-to-date!
```

```
all_positive_tweets = twitter_samples.strings('positive_tweets.json')
all_negative_tweets = twitter_samples.strings('negative_tweets.json')

print('Number of positive tweets: ', len(all_positive_tweets))
print('Number of negative tweets: ', len(all_negative_tweets))

print('\nThe type of all_positive_tweets is: ', type(all_positive_tweets))
print('The type of a tweet entry is: ', type(all_negative_tweets[0]))
```
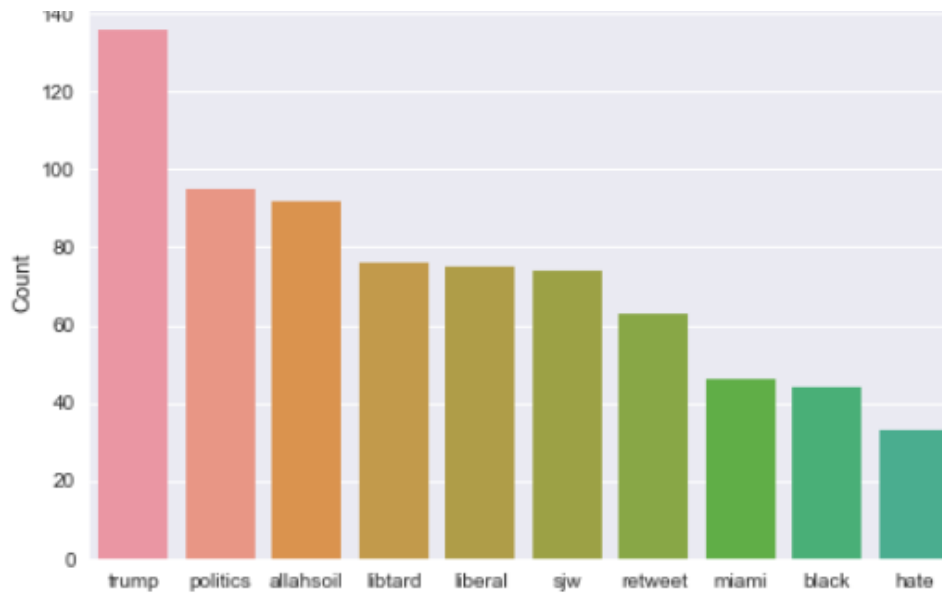


## TWEETS OF DATA

```
train = pd.read_csv('../input/twitter-tweets-data/train_tweet.csv')
test = pd.read_csv('../input/twitter-tweets-data/test_tweets.csv')
```
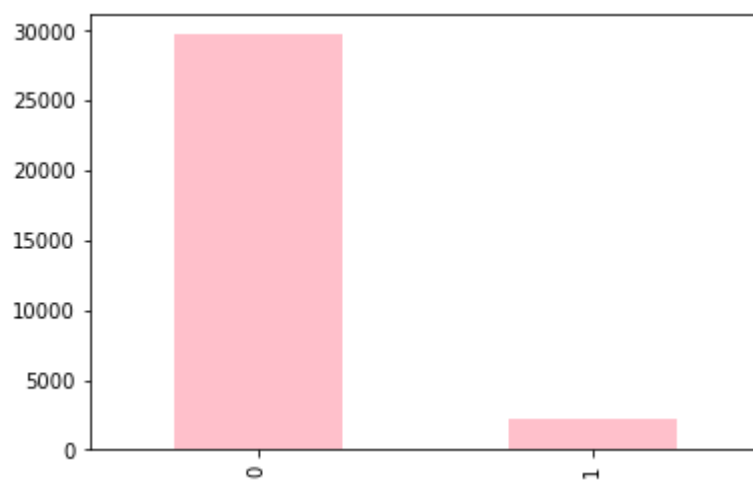
```
print(train.shape)
print(test.shape)
```
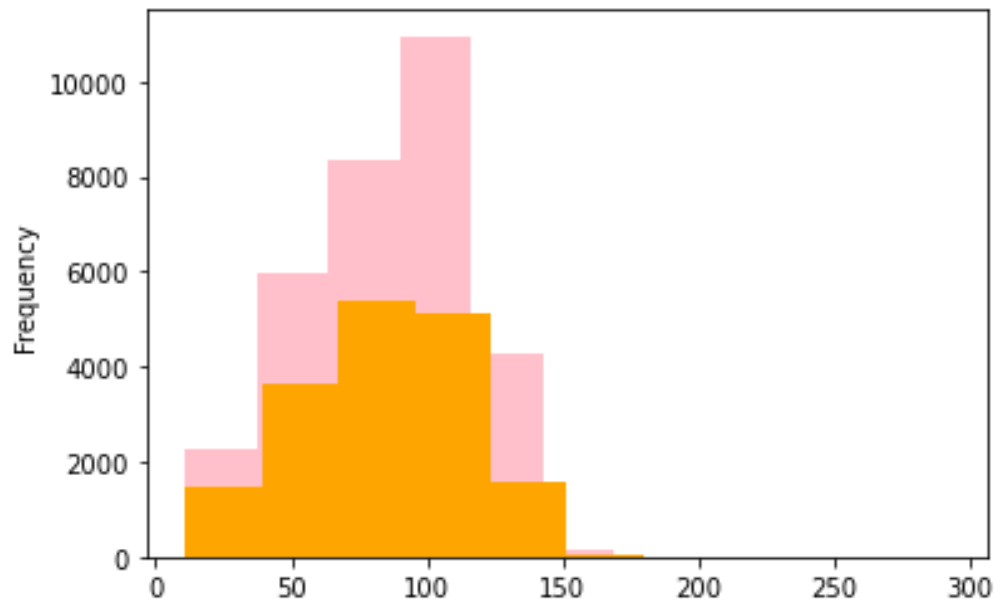
## COUNT THE TWEETS



```
train['label'].value_counts().plot.bar(color = 'pink', figsize = (6, 4))
```

## DISTRIBUTION OF TWEETS

*# checking the distribution of tweets in the data*

length_train = train['tweet'].str.len().plot.hist(color = 'pink', figsize = (6, 4))
length_test = test['tweet'].str.len().plot.hist(color = 'orange', figsize = (6, 4))



COUNTS OF TWEETS

**from** sklearn.feature_extraction.text **import** CountVectorizer

cv = CountVectorizer(stop_words = 'english')
words = cv.fit_transform(train.tweet)

sum_words = words.sum(axis=0)

words_freq = [(word, sum_words[0, i]) **for** word, i **in** cv.vocabulary_.items()]
words_freq = sorted(words_freq, key = **lambda** x: x[1], reverse = **True**)
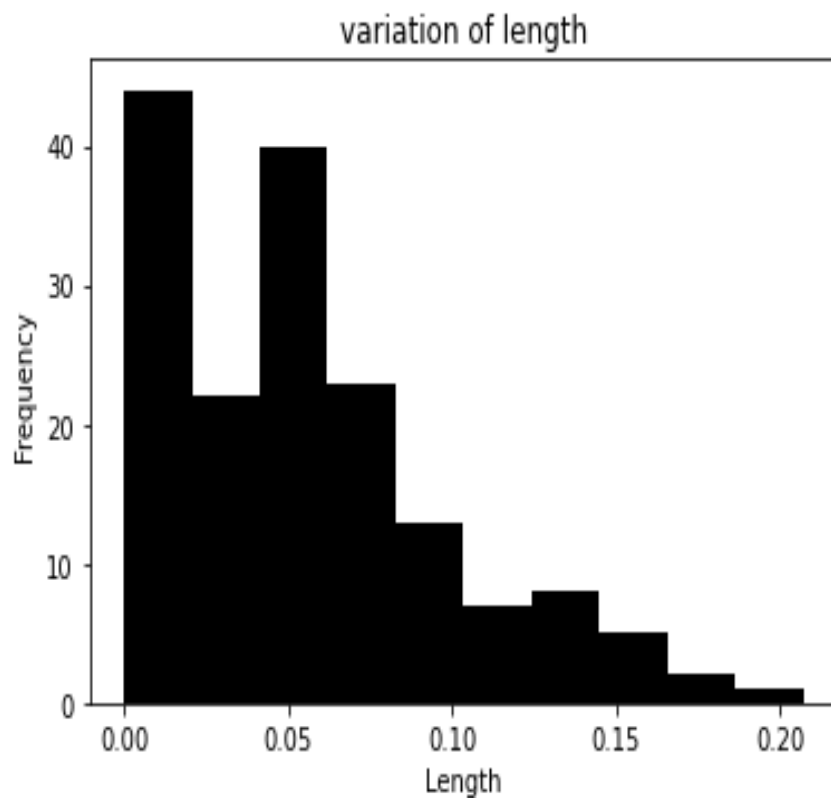
frequency = pd.DataFrame(words_freq, columns=['word', 'freq'])

frequency.head(30).plot(x='word', y='freq', kind='bar', figsize=(15, 7), color = 'blue')
plt.title("Most Frequently Occuring Words - Top 30")

NEGATIVE TWEETS

```
negative_words =' '.join([text for text in train['tweet'][train['label'] == 1]])

wordcloud = WordCloud(background_color = 'red', width=800, height=500, random_state = 0,
max_font_size = 110).generate(negative_words)
plt.figure(figsize=(10, 7))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.title('The Negative Words')
plt.show()
```



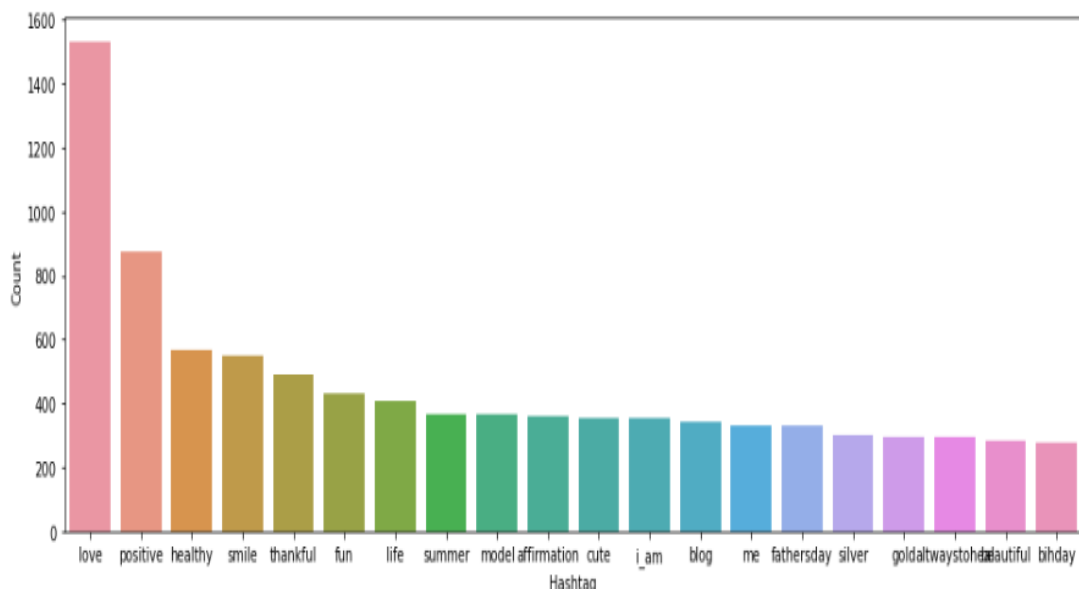*# collecting the hashtags*

```
def hashtag_extract(x):
    hashtags = []

    for i in x:
        ht = re.findall(r"#(\w+)", i)
        hashtags.append(ht)

    return hashtags
```

```
a = nltk.FreqDist(HT_negative)
d = pd.DataFrame({'Hashtag': list(a.keys()),
          'Count': list(a.values())})

# selecting top 20 most frequent hashtags
d = d.nlargest(columns="Count", n = 20)
plt.figure(figsize=(16,5))
ax = sns.barplot(data=d, x= "Hashtag", y = "Count")
ax.set(ylabel = 'Count')
plt.show()
```
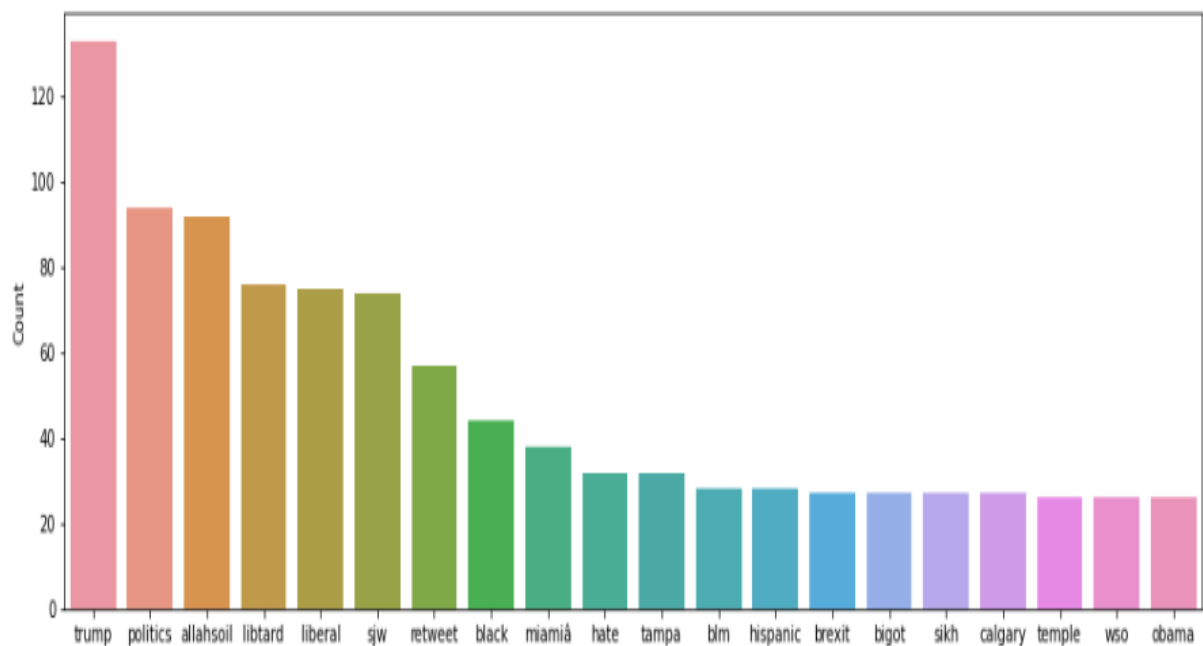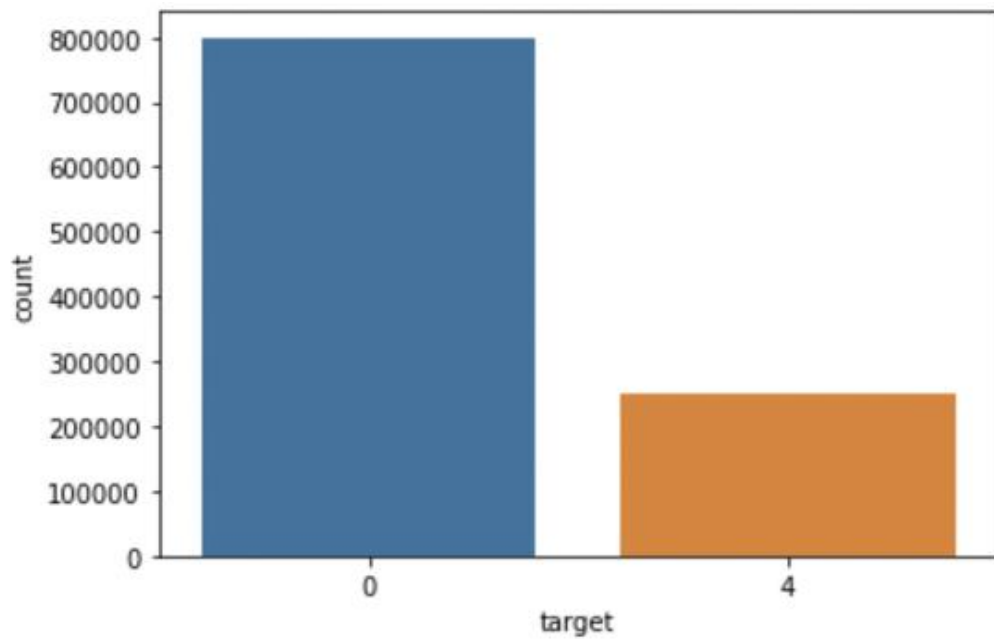


```
model_w2v.wv.most_similar(positive = "dinner")
```

CONCLUSION

It is a Natural Language Processing Problem where Sentiment Analysis is done by Classifying the Positive tweets from negative tweets by machine learning models for classification, text mining, text analysis, data analysis and data visualizatio