

SENTIMENT ANALYSIS

Presented By

BALARAM REDDY (RA2111032010012)

SANJAY G (RA2111032010058)



ABSTRACT

Sentiment analysis is the process of identifying and extracting subjective information in text, such as opinions, emotions, and attitudes. It is a powerful tool that can be used to gain insights into human sentiment.

The project aims to develop a sentiment analysis model using logistic regression and natural language processing (NLP) techniques.

The model is expected to be able to predict the sentiment of tweets and posts with high accuracy, which can be used to improve a variety of applications, such as customer service, marketing, and public relations.



OBJECTIVE

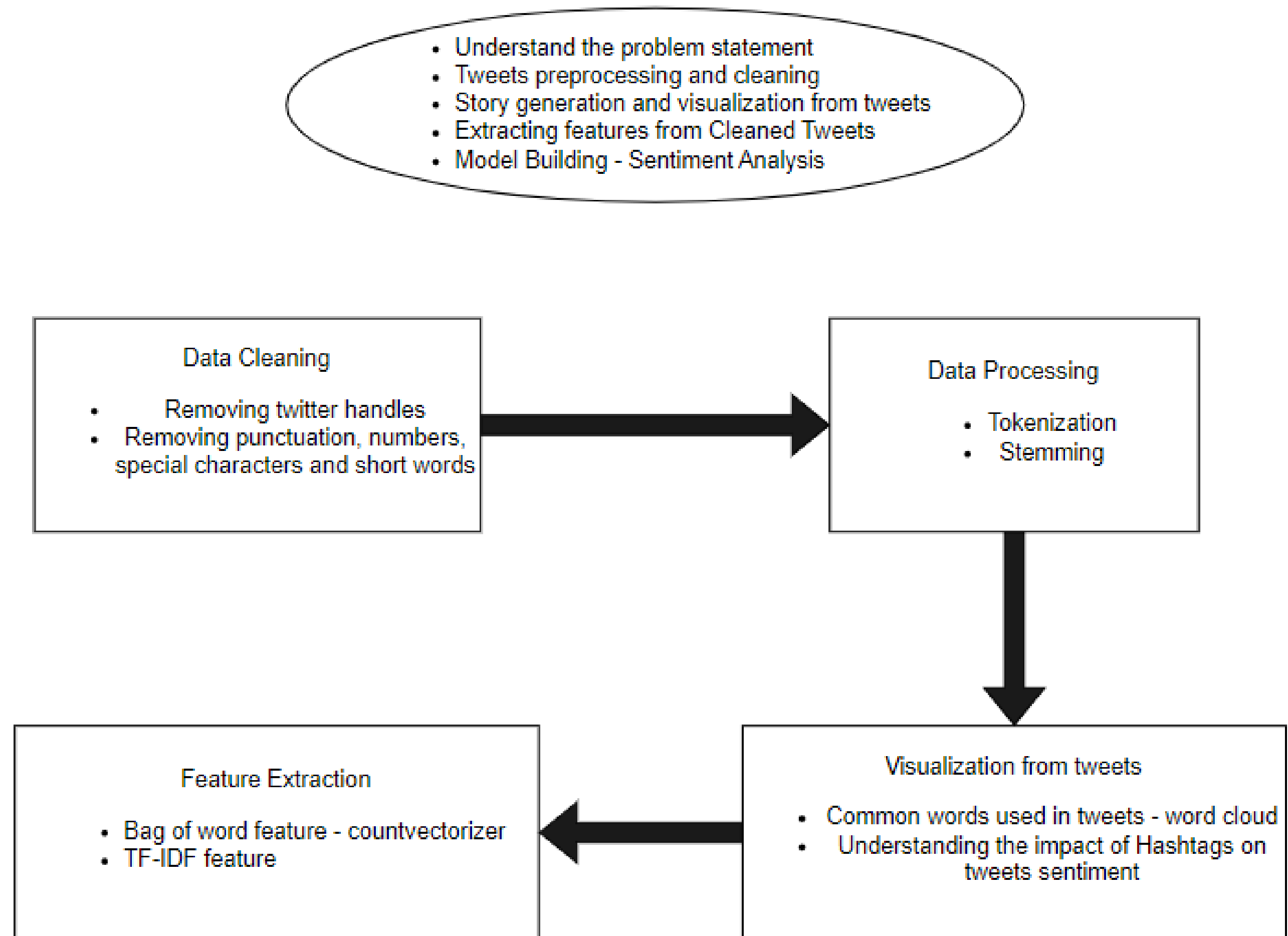
Sentiment analysis is an analytical technique that uses statistics, natural language processing, and machine learning to determine the emotional meaning of communications OR tweets.

The specific objectives of this project are to:

- Identifying the sentiment of text
- Measuring the intensity of sentiment
- Categorizing sentiment
- Identifying the target of sentiment
- Understanding the context of sentiment



SYSTEM FLOW



```
[ ] import re
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import string
import nltk
import warnings
warnings.filterwarnings("ignore", category=DeprecationWarning)

%matplotlib inline
```

```
▶ train = pd.read_csv('/content/train.csv')
test = pd.read_csv('/content/test.csv')
train.head()
```

```
➞
```

	id	label	tweet
0	1	0	@user when a father is dysfunctional and is s...
1	2	0	@user @user thanks for #lyft credit i can't us...
2	3	0	bihday your majesty
3	4	0	#model i love u take with u all the time in ...
4	5	0	factsguide: society now #motivation

```

▶ #data cleaning
combi = train.append(test, ignore_index=True)
def remove_pattern(input_txt, pattern):
    r = re.findall(pattern, input_txt)
    for i in r:
        input_txt = re.sub(i, '', input_txt)

    return input_txt

#removing twitter handles(@user)
combi['tidy_tweet'] = np.vectorize(remove_pattern)(combi['tweet'], "@[\w]*")

```

```

↳ <ipython-input-9-ad3e0576b06c>:2: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.
    combi = train.append(test, ignore_index=True)

```

```

[ ] # remove special characters, numbers, punctuations
combi['tidy_tweet'] = combi['tidy_tweet'].str.replace("[^a-zA-Z#]", " ")

```

```

<ipython-input-10-7266f6bf6f90>:2: FutureWarning: The default value of regex will change from True to False in a future version.
    combi['tidy_tweet'] = combi['tidy_tweet'].str.replace("[^a-zA-Z#]", " ")

```

```

[ ] #remove short word
combi['tidy_tweet'] = combi['tidy_tweet'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>3]))
combi['tidy_tweet'].head()

```

```

0    when father dysfunctional selfish drags kids i...
1    thanks #lyft credit cause they offer wheelchai...
2                                bihday your majesty
3                                #model love take with time
4                                factsguide society #motivation
Name: tidy_tweet, dtype: object

```

#Tokenization

```
tokenized_tweet = combi['tidy_tweet'].apply(lambda x: x.split())
tokenized_tweet.head()
```

```
0    [when, father, dysfunctional, selfish, drags, ...
1    [thanks, #lyft, credit, cause, they, offer, wh...
2                [bihday, your, majesty]
3                [#model, love, take, with, time]
4                [factsguide, society, #motivation]
Name: tidy_tweet, dtype: object
```

[] #Stemming

```
from nltk.stem.porter import *
stemmer = PorterStemmer()
```

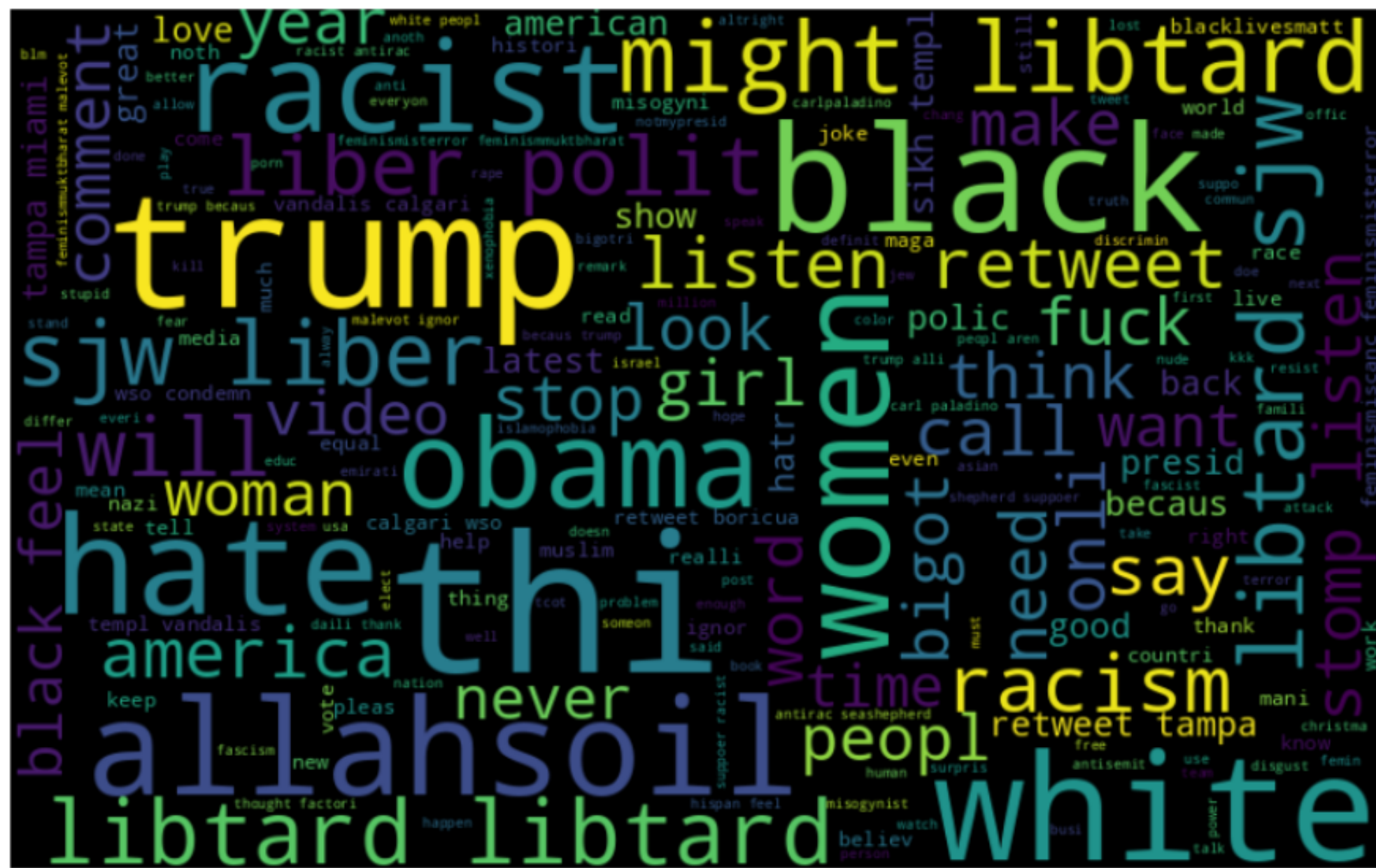
```
tokenized_tweet = tokenized_tweet.apply(lambda x: [stemmer.stem(i) for i in x]) # stemming
tokenized_tweet.head()
```

```
0    [when, father, dysfunct, selfish, drag, kid, i...
1    [thank, #lyft, credit, caus, they, offer, whee...
2                [bihday, your, majesti]
3                [#model, love, take, with, time]
4                [factsguid, societi, #motiv]
Name: tidy_tweet, dtype: object
```

[] #stitch

```
for i in range(len(tokenized_tweet)):
    tokenized_tweet[i] = ' '.join(tokenized_tweet[i])
```

```
combi['tidy_tweet'] = tokenized_tweet
```

```
▶ # function to collect hashtags
def hashtag_extract(x):
    hashtags = []
    # Loop over the words in the tweet
    for i in x:
        ht = re.findall(r"#(\w+)", i)
        hashtags.append(ht)

    return hashtags
```

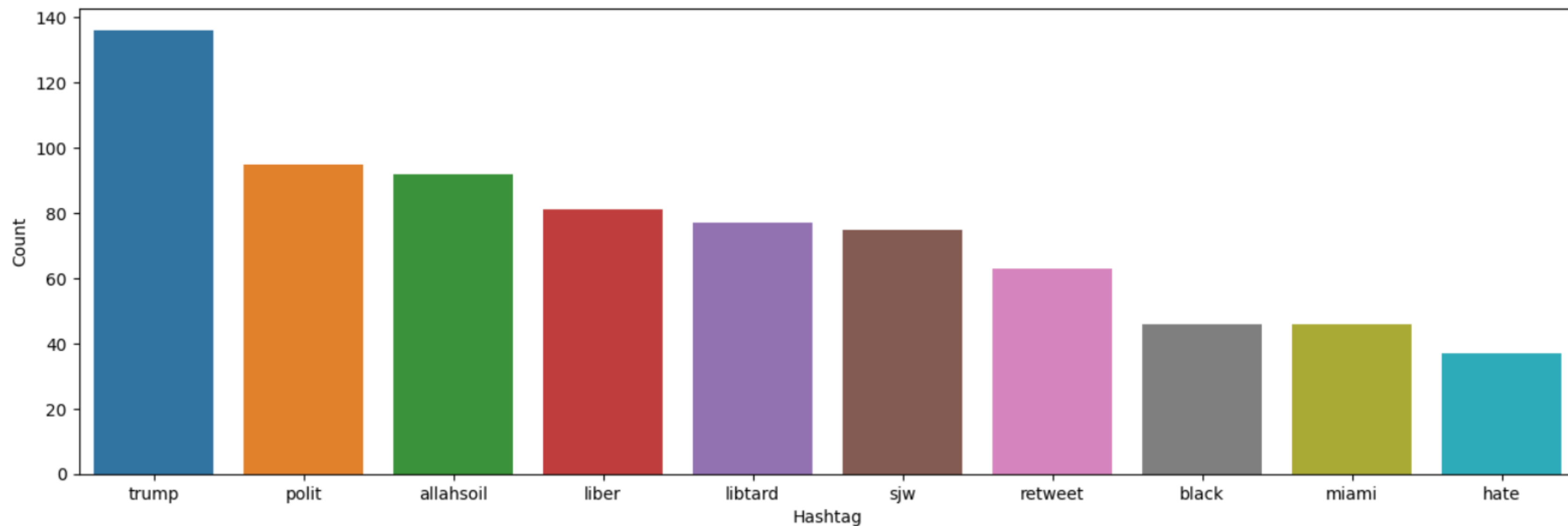
```
▶ # extracting hashtags from non racist/sexist tweets

HT_regular = hashtag_extract(combi['tidy_tweet'][combi['label'] == 0])

# extracting hashtags from racist/sexist tweets
HT_negative = hashtag_extract(combi['tidy_tweet'][combi['label'] == 1])

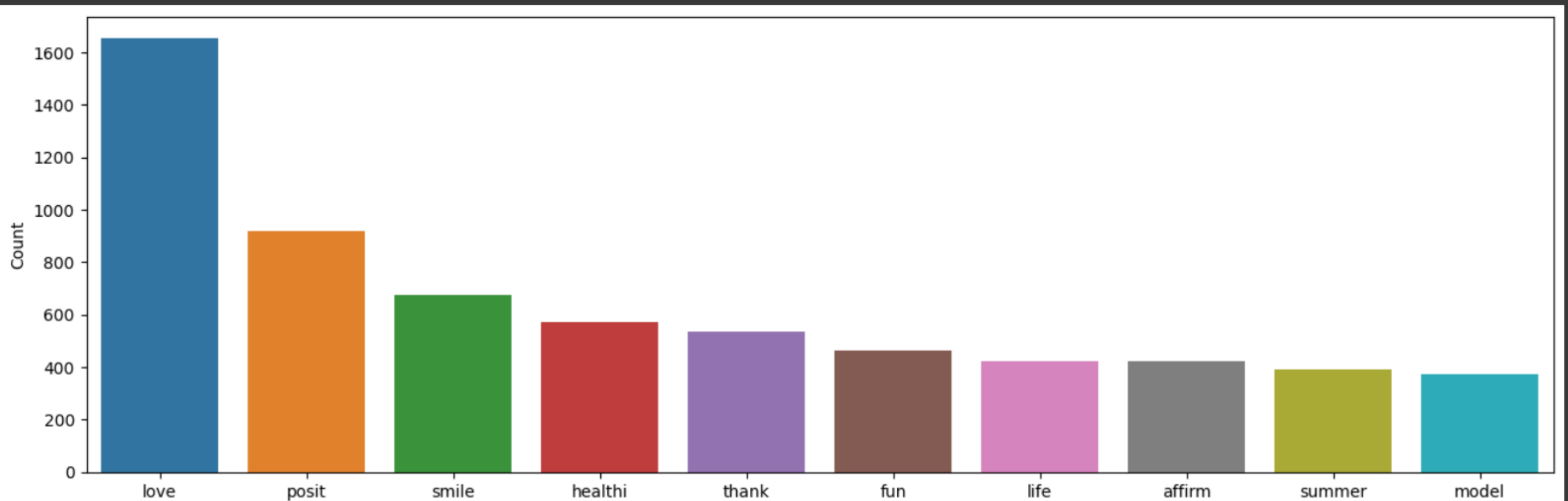
# unnesting list
HT_regular = sum(HT_regular,[])
HT_negative = sum(HT_negative,[])
```

```
#for Racist/Sexist Tweets
b = nltk.FreqDist(HT_negative)
e = pd.DataFrame({'Hashtag': list(b.keys()), 'Count': list(b.values())})
# selecting top 10 most frequent hashtags
e = e.nlargest(columns="Count", n = 10)
plt.figure(figsize=(16,5))
ax = sns.barplot(data=e, x= "Hashtag", y = "Count")
ax.set(ylabel = 'Count')
plt.show()
```




```
#non Racist/Sexist Tweets
a = nltk.FreqDist(HT_regular)
d = pd.DataFrame({'Hashtag': list(a.keys()),
                  'Count': list(a.values())})

# selecting top 10 most frequent hashtags
d = d.nlargest(columns="Count", n = 10)
plt.figure(figsize=(16,5))
ax = sns.barplot(data=d, x= "Hashtag", y = "Count")
ax.set(ylabel = 'Count')
plt.show()
```



```
#CountVectorizer
```

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
bow_vectorizer = CountVectorizer(max_df=0.90, min_df=2, max_features=1000, stop_words='english')
```

```
# bag-of-words feature matrix
```

```
bow = bow_vectorizer.fit_transform(combi['tidy_tweet'])
```

```
bow
```

```
<49159x1000 sparse matrix of type '<class 'numpy.int64'>'
      with 191502 stored elements in Compressed Sparse Row format>
```

```
#TF-IDF Features
```

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
tfidf_vectorizer = TfidfVectorizer(max_df=0.90, min_df=2, max_features=1000, stop_words='english')
```

```
# TF-IDF feature matrix
```

```
tfidf = tfidf_vectorizer.fit_transform(combi['tidy_tweet'])
```

```
tfidf
```

```
<49159x1000 sparse matrix of type '<class 'numpy.float64'>'
      with 191502 stored elements in Compressed Sparse Row format>
```

```
#Building model using Bag-of-Words features(CountVectorizer)
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import f1_score

train_bow = bow[:31962,:]
test_bow = bow[31962:,:]

# splitting data into training and validation set
xtrain_bow, xvalid_bow, ytrain, yvalid = train_test_split(train_bow, train['label'], random_state=10, test_size=0.3)

lreg = LogisticRegression()
lreg.fit(xtrain_bow, ytrain) # training the model

prediction = lreg.predict_proba(xvalid_bow) # predicting on the validation set
prediction_int = prediction[:,1] >= 0.3 # if prediction is greater than or equal to 0.3 than 1 else 0
prediction_int = prediction_int.astype(np.int)

f1_score(yvalid, prediction_int) # calculating f1 score
```

0.5641447368421053

```
#Building model using TF-IDF features

train_tfidf = tfidf[:31962,:]
test_tfidf = tfidf[31962:,:]

xtrain_tfidf = train_tfidf[ytrain.index]
xvalid_tfidf = train_tfidf[yvalid.index]

lreg.fit(xtrain_tfidf, ytrain)

prediction = lreg.predict_proba(xvalid_tfidf)
prediction_int = prediction[:,1] >= 0.3
prediction_int = prediction_int.astype(np.int)
f1_score(yvalid, prediction_int)
```

0.5726643598615917



**THANK
YOU**