

In [4]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import re
import string
import nltk
import warnings
%matplotlib inline

warnings.filterwarnings('ignore')
```

In [5]:

```
#add Data set
df =pd.read_csv('train_E6oV3lV.csv')
df.head()
```

Out[5]:

	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

In [6]:

```
# datatype info
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 31962 entries, 0 to 31961
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  -
0    id      31962 non-null     int64
1   label   31962 non-null     int64
2   tweet   31962 non-null     object
dtypes: int64(2), object(1)
memory usage: 749.2+ KB
```

In [9]:

```
#remove pattern
def remove_pattern(input_txt, pattern):
    r = re.findall(pattern, input_txt)
    for word in r:
        input_txt = re.sub(word, "", input_txt)
    return input_txt
```

In [10]:

df.head()

Out[10]:

	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

In [11]:

```
# remove twitter handles (@user)
df['clean_tweet'] = np.vectorize(remove_pattern)(df['tweet'], "@[\w]*")
```

In [12]:

df.head()

Out[12]:

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

In [13]:

```
# remove special characters, numbers and punctuations
df['clean_tweet'] = df['clean_tweet'].str.replace("[^a-zA-Z#]", " ")
df.head()
```

Out[13]:

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

In [14]:

```
# remove short words
df['clean_tweet'] = df['clean_tweet'].apply(lambda x: " ".join([w for w in x.split() if len(w) > 3]))
df.head()
```

Out[14]:

	id	label	tweet	clean_tweet
0	1	0	@user when a father is dysfunctional and is s...	when father dysfunctional selfish drags kids i...
1	2	0	@user @user thanks for #lyft credit i can't us...	thanks #lyft credit can't cause they don't off...
2	3	0	bihday your majesty	bihday your majesty
3	4	0	#model i love u take with u all the time in ...	#model love take with time urð±!!! ðððððððð...
4	5	0	factsguide: society now #motivation	factsguide: society #motivation

In [15]:

```
# individual words considered as tokens
tokenized_tweet = df['clean_tweet'].apply(lambda x: x.split())
tokenized_tweet.head()
```

Out[15]:

```
0    [when, father, dysfunctional, selfish, drags, ...
1    [thanks, #lyft, credit, can't, cause, they, do...
2                                [bihday, your, majesty]
3    [#model, love, take, with, time, urð±!!!, ð...
4                                [factsguide:, society, #motivation]
Name: clean_tweet, dtype: object
```

In [16]:

```
import nltk
from nltk.stem import PorterStemmer
```

In [17]:

```
ps = PorterStemmer()
tokenized_tweet = tokenized_tweet.apply(lambda sentence: [(word) for word in sentence])
tokenized_tweet.head()
```

Out[17]:

```
0    [when, father, dysfunctional, selfish, drags, ...
1    [thanks, #lyft, credit, can't, cause, they, do...
2                                [bihday, your, majesty]
3    [#model, love, take, with, time, urð±!!!, ð...
4                                [factsguide:, society, #motivation]
Name: clean_tweet, dtype: object
```

```
# combine words into single sentence
for i in range(len(tokenized_tweet)):
    tokenized_tweet[i] = " ".join(tokenized_tweet[i])

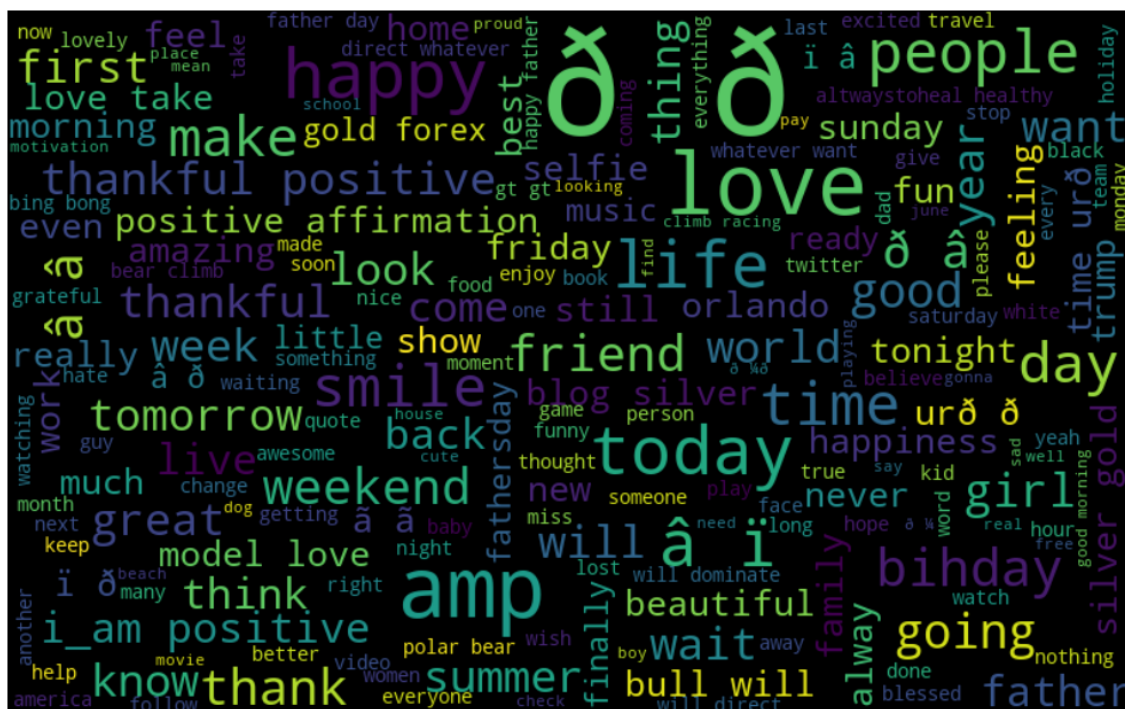
df['clean_tweet'] = tokenized_tweet
df.head()
```

	id	label	tweet	clean_tweet
0	1	0	@user when a father is dysfunctional and is s...	when father dysfunctional selfish drags kids i...
1	2	0	@user @user thanks for #lyft credit i can't us...	thanks #lyft credit can't cause they don't off...
2	3	0	bihday your majesty	bihday your majesty
3	4	0	#model i love u take with u all the time in ...	#model love take with time urð□□±!!! ð□□□ð□□ð...
4	5	0	factsguide: society now #motivation	factsguide: society #motivation

```
# visualize the frequent words
all_words = " ".join([sentence for sentence in df['clean_tweet']])

from wordcloud import WordCloud
wordcloud = WordCloud(width=800, height=500, random_state=42, max_font_size=100).generate

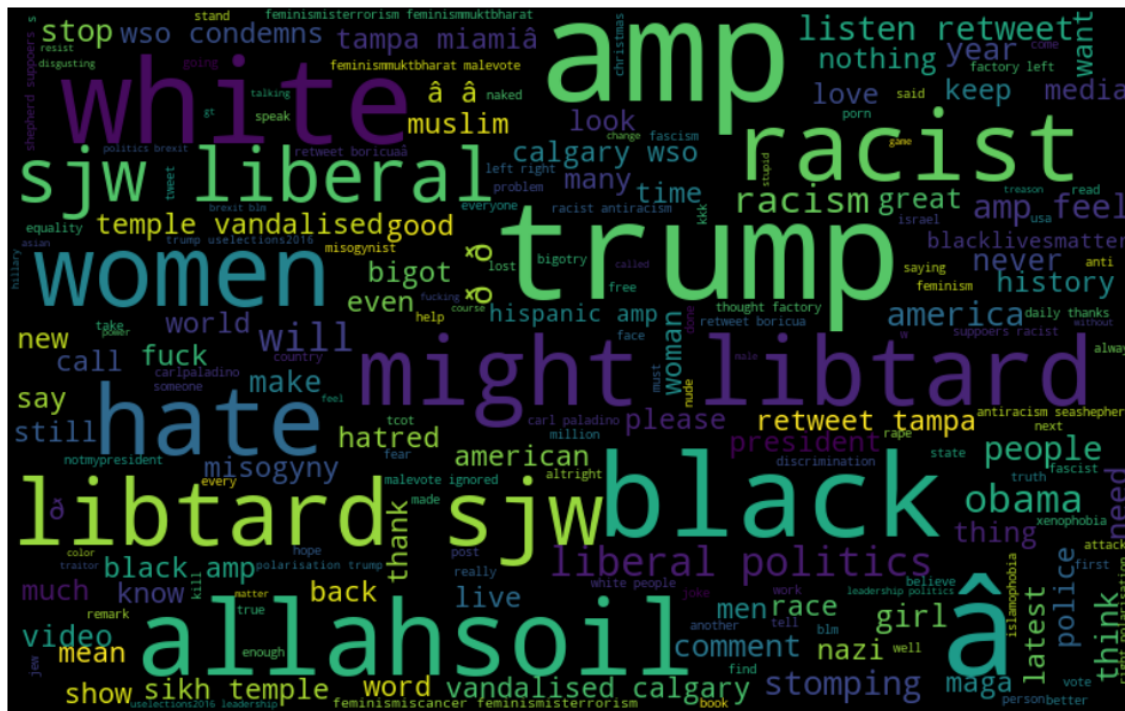
# plot the graph
plt.figure(figsize=(15,8))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.show()
```




```
# frequent words visualization for -ve
all_words = " ".join([sentence for sentence in df['clean_tweet'][df['label']==1]])

wordcloud = WordCloud(width=800, height=500, random_state=42, max_font_size=100).generate

# plot the graph
plt.figure(figsize=(15,8))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off')
plt.show()
```



```
# extract the hashtag
def hashtag_extract(tweets):
    hashtags = []
    # Loop words in the tweet
    for tweet in tweets:
        ht = re.findall(r"#(\w+)", tweet)
        hashtags.append(ht)
    return hashtags
```

```
# extract hashtags from non-racist/sexist tweets
ht_positive = hashtag_extract(df['clean_tweet'][df['label']==0])
```

```
# extract hashtags from racist/sexist tweets
ht_negative = hashtag_extract(df['clean_tweet'][df['label']==1])
```

In [25]:

```
ht_positive[:5]
```

Out[25]:

```
[['run'], ['lyft', 'disappointed', 'getthanked'], [], ['model'], ['motivati  
on']]
```

In [26]:

```
# unnest list  
ht_positive = sum(ht_positive, [])  
ht_negative = sum(ht_negative, [])
```

In [27]:

```
ht_positive[:5]
```

Out[27]:

```
['run', 'lyft', 'disappointed', 'getthanked', 'model']
```

In [28]:

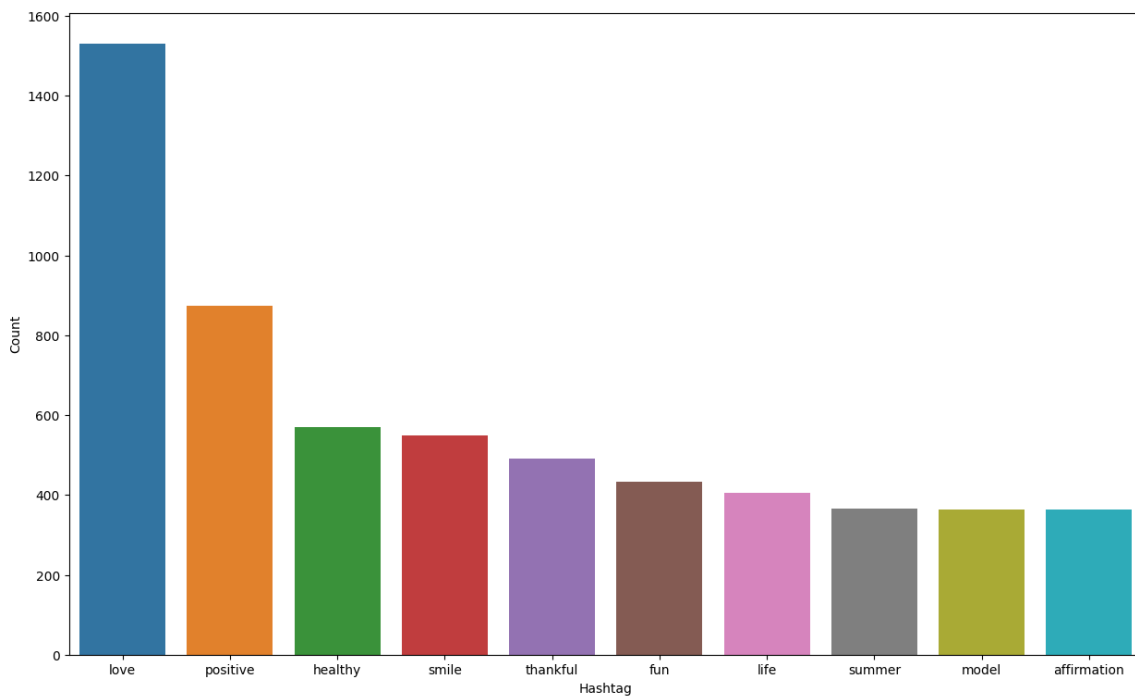
```
freq = nltk.FreqDist(ht_positive)  
d = pd.DataFrame({'Hashtag': list(freq.keys()),  
                  'Count': list(freq.values())})  
d.head()
```

Out[28]:

	Hashtag	Count
0	run	33
1	lyft	2
2	disappointed	1
3	getthanked	2
4	model	364

In [29]:

```
# select top 10 hashtags
d = d.nlargest(columns='Count', n=10)
plt.figure(figsize=(15,9))
sns.barplot(data=d, x='Hashtag', y='Count')
plt.show()
```



In [30]:

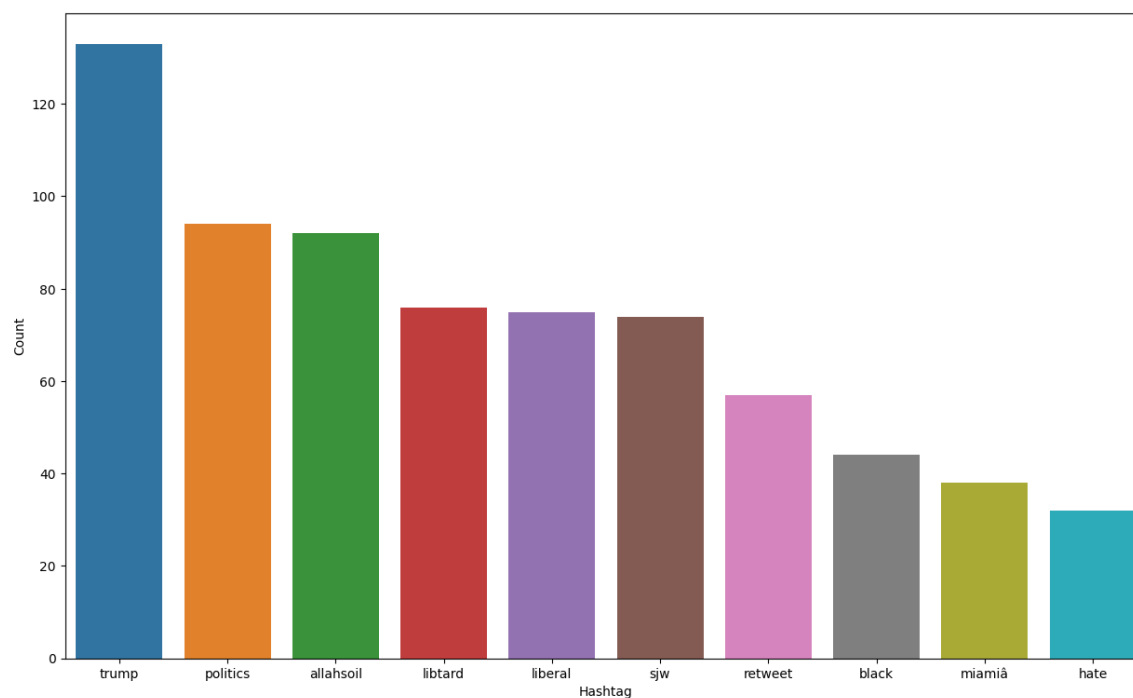
```
freq = nltk.FreqDist(ht_negative)
d = pd.DataFrame({'Hashtag': list(freq.keys()),
                  'Count': list(freq.values())})
d.head()
```

Out[30]:

	Hashtag	Count
0	cnn	9
1	michigan	2
2	tcot	14
3	australia	6
4	opkillingbay	2

In [31]:

```
# select top 10 hashtags  
d = d.nlargest(columns='Count', n=10)  
plt.figure(figsize=(15,9))  
sns.barplot(data=d, x='Hashtag', y='Count')  
plt.show()
```



In []:

In []: