

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
plt.style.use('ggplot')
import seaborn as sns
import numpy as np

In [2]: os.environ['KAGGLE_USERNAME'] = "shenshui1992" # username from the json file
os.environ['KAGGLE_KEY'] = "94596f3e366e032f965e53754f9524b" # key from the json file

In [3]: !kaggle datasets download -d manchunhui/us-election-2020-tweets

us-election-2020-tweets.zip: skipping, found more recently modified local copy (use --force to force download)

Download datasets from kaggle

In [4]: !unzip us-election-2020-tweets.zip

Archive: us-election-2020-tweets.zip
  replace hashtag_donaldtrump.csv? [y], [n], [A], [I], [N], [e], [r], [name]: A
  inflating: hashtag_donaldtrump.csv
  inflating: hashtag_joe Biden.csv

unzip the zipfile just download

In [5]: df_j=pd.read_csv('hashtag_joe Biden.csv',lineterminator='\n')

put csv file in data frame

In [6]: df_j.sample(5)

Out [6]:
   created_at      tweet_id      tweet      likes      retweet_count      source      user_id      user_name      user_screen_name      user
0  2020-11-06  13245176e+18  #Biden coming thru in the polls right now! 1.0      0.0      Twitter for iPhone      1.809547e+08  AkeemRakim...  ZamRakim  Akeem alwi
1  2020-11-06  13247031e+18  The @Election Choo Crew is a Nice Break F... 15.0      17.0      Twitter Web App      1.991527e+08  The Dollar Vigilante  DollarVigilante  The C with...
2  2020-11-08  1325231e+18  Debuted with 62% Scooby from the stream! 0.0      0.0      Twitter for Android      3.253185e+09  These Are Special Times - CelticDon  CelticSuperG  https://cc...
3  2020-11-07  13251201e+18  Likes @ZDF With them Landspage! 1.0      0.0      Twitter for Android      1.260546e+18  Siegfried Goldstein  Siegfried_GSt  Mann Kin
4  2020-11-04  1323375e+18  Trump secures another state now it's Texas! 2.0      1.0      Twitter for Android      1.267470e+18  Adrian Kore Sports  adrian_kore  i vntofgre

In [7]: df_j['target']="joebiden"

add a column name "target"

In [8]: df_j = df_j.drop(columns = ['created_at', 'tweet_id', 'user_id', 'user_name', 'user_screen_name', 'user_description', 'user_join_date', 'collected_at'])

drop useless columns

In [9]: df_j=df_j.rename(columns={"tweet": "text", "who": "target"})

In [10]: df_t=pd.read_csv('hashtag_donaldtrump.csv',lineterminator='\n')

put csv file in data frame

In [11]: df_t.sample(5)

Out [11]:
   created_at      tweet_id      tweet      likes      retweet_count      source      user_id      user_name      user_screen_name      user
0  2020-10-29  13218352e+18  Ma guarda quant'cast stene...a che f... 1.0      0.0      Twitter for Android      1.282230e+18  Antonella Starace  AntonellaStaras3  f...
1  2020-11-06  13245020e+18  Ruista wants to undermine confidence in our el... 1.0      0.0      Twitter Web App      1.062068e+09  TacosTheCrow  tacosthecrow  li...
2  2020-11-04  13239209e+18  Zet er vermenen nu remand van het #CDA bij #W... 3.0      0.0      Twitter for iPhone      1.082735e+18  Klaaskaaal@klaaskaaal  Klaaskaaal  Klaaskaaal
3  2020-11-03  1323615e+18  Left vote for #ElectionDay 4.0      1.0      Twitter for iPhone      2.864475e+09  Assad Ad Awan rk  markhor_2018  Ph...
4  2020-10-18  13178956e+18  A Trump saying the quiet part out loud... 27.0      27.0      Twitter for iPhone      9.792121e+17  Save Our Country  BigBlueWaveUSA

In [12]: df_t = df_t.drop(columns = ['created_at', 'tweet_id', 'user_id', 'user_name', 'user_screen_name', 'user_description', 'user_join_date', 'collected_at'])

drop useless columns

In [13]: df_t['target']="donaldtrump"

add a column name "target"

In [14]: df_t=df_t.rename(columns={"tweet": "text"})

rename column "tweet" to "text"

In [15]: df_j_t = pd.concat([df_t, df_j]).reset_index()

combine two dataframe together

In [16]: df_j_t = df_j_t.drop(columns = ['index'])

drop "index" column

In [17]: df_j_t.head(5)

Out [17]:
   text      likes      retweet_count      source      user_followers_count      user_location      lat      long      city      country
0  #Elecciones2020 En #Florida: #JoeBiden dice ... 0.0      0.0      TweetDeck      1860.0      Philadelphia, PA / Miami, FL  25.774270  -80.193660  NaN      United States of America
1  Usa 2020, trump contro facebook e twitt... 26.0      9.0      Social MediaSet      1067661.0      NaN      NaN      NaN      NaN      NaN
2  #Trump: As a student i used to hear for years... 2.0      1.0      Twitter Web App      1185.0      Portland      45.520247  -122.674195  Portland      United States of America
3  2 hours since last tweet from #Trump Maybe he... 0.0      0.0      Trumpyweeler      32.0      NaN      NaN      NaN      NaN      NaN
4  You get a tie! And you get a tie! #Trump's ra... 4.0      3.0      Twitter for iPhone      5393.0      Washington DC  38.894992  -77.036556  Washington      United States of America

In [18]: df_j_t.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1747805 entries, 0 to 1747804
Data columns (total 14 columns):
 #   Column      Dtype
---  ---
0   text      object
1   likes      float64
2   retweet_count      float64
3   source      object
4   user_followers_count      float64
5   user_location      object
6   lat      float64
7   long      float64
8   city      object
9   country      object
10  continent      object
11  state      object
12  state_code      object
13  target      object
dtypes: float64(6), object(9)
memory usage: 186.7+ MB

In [19]: df_j_t.head(5)

Out [19]:
   text      likes      retweet_count      source      user_followers_count      user_location      lat      long      city      country
0  #Elecciones2020 En #Florida: #JoeBiden dice ... 0.0      0.0      TweetDeck      1860.0      Philadelphia, PA / Miami, FL  25.774270  -80.193660  NaN      United States of America
1  Usa 2020, trump contro facebook e twitt... 26.0      9.0      Social MediaSet      1067661.0      NaN      NaN      NaN      NaN      NaN
2  #Trump: As a student i used to hear for years... 2.0      1.0      Twitter Web App      1185.0      Portland      45.520247  -122.674195  Portland      United States of America
3  2 hours since last tweet from #Trump Maybe he... 0.0      0.0      Trumpyweeler      32.0      NaN      NaN      NaN      NaN      NaN
4  You get a tie! And you get a tie! #Trump's ra... 4.0      3.0      Twitter for iPhone      5393.0      Washington DC  38.894992  -77.036556  Washington      United States of America

In [20]: df_j_t.describe()

Out [20]:
   likes      retweet_count      user_followers_count      lat      long
count      1747805e+00      1.747805e+06      2.538003e+04      35.434583      -41.083772
mean      8.670996e+00      1.890808e+00      2.538003e+04      35.434583      -41.083772
std      2.860510e+02      7.101567e+01      3.527233e+05      18.425141      67.666098
min      0.000000e+00      0.000000e+00      0.000000e+00      -90.000000      -175.202642
25%      0.000000e+00      0.000000e+00      7.500000e+01      31.816038      -87.086700
50%      0.000000e+00      0.000000e+00      4.350000e+02      39.783730      -74.009615
75%      2.000000e+00      0.000000e+00      2.072002e+03      45.520247      6.130161
max      1.657020e+05      6.347300e+04      8.241710e+07      90.000000      179.048837

In [21]: df_j_t.columns

Out [21]:
Index(['text', 'likes', 'retweet_count', 'source', 'user_followers_count', 'user_location', 'lat', 'long', 'city', 'country', 'continent', 'state', 'state_code', 'target'],
      dtype='object')

In [22]: df_text_target = df_j_t.drop(columns = ['likes', 'retweet_count', 'source', 'user_followers_count', 'user_location', 'lat', 'long', 'city', 'country', 'continent', 'state', 'state_code'])

create a new dataframe only have text and target

In [23]: df_text_target.head(5)

Out [23]:
   text      target
0  #Elecciones2020 En #Florida: #JoeBiden dice ...  donaldtrump
1  Usa 2020, Trump contro Facebook e Twitt...  donaldtrump
2  #Trump: As a student i used to hear for years...  donaldtrump
3  2 hours since last tweet from #Trump! Maybe he...  donaldtrump
4  You get a tie! And you get a tie! #Trump's ra...  donaldtrump

In [24]: df_text_target.count()

Out [24]:
text      1747805
target    1747805
dtype: int64

In [25]: x = df_text_target["target"].value_counts()
plt.grid()
sns.boxplot(x=index, y=x)
plt.gca().set_ylabel('samples')
plt.title("distribution")

/usr/local/lib/python3.6/dist-packages/seaborn/decorators.py:43: FutureWarning: Pass the following v
ariables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data',
and passing other arguments without an explicit keyword will result in an error or misinterpretation.
FutureWarning:

Out [25]:
Text(0.5, 1.0, 'distribution')

In [26]: plt.grid()
plt.hist(df_text_target[df_text_target["target"] == "donaldtrump"]["text"].str.len(), color="b")
plt.title("Trump tweets length")

Out [26]:
Text(0.5, 1.0, 'Trump tweets length')

In [27]: plt.grid()
plt.hist(df_text_target[df_text_target["target"] == "joebiden"]["text"].str.len(), color="r")
plt.title("Biden tweets length")

Out [27]:
Text(0.5, 1.0, 'Biden tweets length')

In [28]: sns.catplot(y="continent", data=df_j_t, hue="target", kind="count")

<Figure size 720x360 with 0 Axes>

Most unique words in each category

In [29]: features_categorical = ['source', 'user_location', 'city', 'country', 'state', 'target']
for i in features_categorical:
    print(df_j_t[i].value_counts()[1:])
    print("=====")

Twitter Web App      561380
Twitter for iPhone   518843
Twitter for Android  488212
Twitter for iPad     61362
TweetDeck            29968
Name: source, dtype: int64
=====
United States      39253
USA                 16105
Washington, DC     10831
India               10680
California, USA    9904
Name: user_location, dtype: int64
=====
New York           37749
London             19204
Los Angeles        16838
Washington         16475
Paris              16173
Name: city, dtype: int64
=====
United States of America  332495
United States            61905
United Kingdom           58051
India                    40091
Germany                  35379
Name: country, dtype: int64
=====
California           56966
New York             33886
England              40854
Texas                30682
Florida              29852
Name: state, dtype: int64
=====
donaldtrump      970919
joebiden          776886
Name: target, dtype: int64
=====

In [30]: import plotly.graph_objects as go

The top 10 sources of tweets

In [31]: df_j_t.query('target == "joebiden").groupby(by="source").count().text.sort_values(ascending=False)[1:]
x= df_j_t.query('target == "joebiden").groupby(by="source").count().text.sort_values(ascending=False)
df_j_t.query('target == "donaldtrump").groupby(by="source").count().text.sort_values(ascending=False)[1:]
x2= df_j_t.query('target == "donaldtrump").groupby(by="source").count().text.sort_values(ascending=False)
df_j_t.index
fig_1 = go.Figure((go.Bar(x=x, y=y, name="Joe Biden"),
                    go.Bar(x=x2, y=y2, name="donaldtrump")))
fig_1.update_layout(title_text="Top 10 sources")
fig_1.update_xaxes(title="source")
fig_1.update_yaxes(title="tweet count")
fig_1.show(renderer="colab")

In [32]: y = df_j_t.query('target == "joebiden" & (country == "United States of America)').groupby(
    by="state").count().text.sort_values(ascending=False)
x = df_j_t.query('target == "joebiden" & (country == "United States of America)').groupby(
    by="state").count().text.sort_values(ascending=False).index
y2 = df_j_t.query('target == "donaldtrump" & (country == "United States of America)').groupby(
    by="state").count().text.sort_values(ascending=False).index
fig_2 = go.Figure((go.Bar(x=x, y=y, name="Joe Biden"),
                    go.Bar(x=x2, y=y2, name="donaldtrump")))
fig_2.update_layout(title_text="Tweets count for each state of America")
fig_2.update_xaxes(title="source")
fig_2.update_yaxes(title="tweets count")
fig_2.show(renderer="colab")

In [33]: import re

Remove URL

In [34]: # Function for url's
import re
def remove_urls(text):
    url_pattern = re.compile(r'(https?://\S+|www.\S+)')
    return url_pattern.sub(r'link', text)

In [35]: df_text_target['text'] = df_text_target['text'].apply(remove_urls)

Lowercase

In [36]: df_text_target['text'] = df_text_target['text'].str.lower()
display(df_text_target['text']).head(5)

0  #elecciones2020 | en #florida: #joebiden dice ...
1  usa 2020, trump contro facebook e twitter: cop...
2  #trump: as a student i used to hear for years...
3  2 hours since last tweet from #trump! maybe he...
4  you get a tie! and you get a tie! #trump's ra...
Name: text, dtype: object

In [37]: not_list = ["don", "ain", "ain't", "aren", "arent", "aren't", "cannot", "cant", "can't", "couldn", "cou
ldn't", "couldn't", "didn't", "doesn", "doesn't", "don", "don't", "hadn", "hasn", "hasn't", "has
n't", "haven", "haven't", "mightn", "mightn't", "n", "n't", "mustn", "mustn't", "needn", "needn't", "nt", "shouldn", "shouldn't",
"wasn", "wasn't", "wasn't", "don't"]

In [38]: def before_lowercase(text):
    text = re.sub(r'["'+not_list+"']', "", text)
    return text

In [39]: df_text_target['text'] = df_text_target['text'].apply(before_lowercase)
df_text_target['text'].head(5)

0  #elecciones2020 | en #florida: #joebiden dice ...
1  usa 2020, trump contro facebook e twitter: co...
2  #trump: as a student i used to hear for years...
3  2 hours since last tweet from #trump! maybe he...
4  you get a tie! and you get a tie! #trump's ra...
Name: text, dtype: object

remove Emojis&Emotions

In [40]: def remove_emoji(text):
    emoji_pattern = re.compile("["
        u'\U0001F600-\U0001F64F'  # emoticons
        u'\U0001F300-\U0001F5FF'  # symbols & pictographs
        u'\U0001F680-\U0001F7FF'  # transport & map symbols
        u'\U0001F1E0-\U0001F1FF'  # flags (iOS)
        u'\U000020AC-\U0001F251'
        "]+", flags=re.UNICODE)
    return emoji_pattern.sub(r'', text)

In [41]: df_text_target['text'] = df_text_target['text'].apply(lambda x:remove_emoji(x))

In [42]: df_text_target.head(5)

Out [42]:
   text      target
0  #elecciones2020 | en #florida: #joebiden dice ...  donaldtrump
1  usa 2020, trump contro facebook e twitt...  donaldtrump
2  #trump: as a student i used to hear for years...  donaldtrump
3  2 hours since last tweet from #trump! maybe he...  donaldtrump
4  you get a tie! and you get a tie! #trump's ra...  donaldtrump

In [43]: from nltk.corpus import stopwords

remove stopwords

In [44]: import nltk
nltk.download('stopwords')

[nltk data] Downloading package stopwords to /root/nltk_data...
[nltk data] Unzipping corpora/stopwords.zip.

In [45]: my_stopwords=set(stopwords.words('english'))

In [46]: print(len(my_stopwords))

179

In [47]: def remove_sw(text,s_list):
    for i in text.split():
        if i not in my_stopwords:
            return s
        else:
            s.append(s)
    return s

In [48]: b=[]
for i in df_text_target['text']:
    b.append(remove_sw(t, my_stopwords))

In [49]: df_text_target['text2']= b
df_text_target['text2'].head(5)

Out [49]:
0  #elecciones2020. | en, #florida, #joebiden...
1  [usa, 2020., trump, contro, facebook, e, twit...
2  [#trump, student, used, hear, years, ten, y...
3  [2, hours, since, last, tweet, #trump], mayb...
4  [get, tie!, get, tie!, #trump, 's, rally, #iow...
Name: text2, dtype: object

In [50]: #combine individual words
def combine_text(input):
    combined = ' '.join(input)
    return combined
df_text_target['text'] = df_text_target['text2'].apply(combine_text)
df_text_target['text']

Out [50]:
0  #elecciones2020 | en #florida: #joebiden dice ...
1  usa 2020, trump contro facebook e twitter: co...
2  #trump: student used hear years ten years heard...
3  2 hours since last tweet #trump! maybe busy trem...
4  get tie get tie trump 's rally #iowa link

1747800 stop laying cmm paris lonnot nott give fuck B...
1747801 ex hñh(u va pu pñc #vri oav touc onaboc tou ...
1747802 l'otan va sortir de sa lèthargie et redevenir ...
1747803 "congiuntinfuoriregione" sono felice per #jil...
1747804 ik moet zeggen dat ik Biden "the lesser two e...
Name: text, Length: 1747805, dtype: object

Remove punctuation

In [51]: import string

In [52]: def remove_punctuation(text):
    table = str.maketrans('', '', string.punctuation)
    return text.translate(table)

In [53]: df_text_target['text'] = df_text_target['text'].apply(lambda x: remove_punctuation(x))

In [54]: df_text_target['text2']

Out [54]:
0  elecciones2020 en florida joebiden dice que n...
1  usa 2020 trump contro facebook e twitter: co...
2  #trump: as a student i used to hear years heard...
3  2 hours since last tweet trump maybe busy trem...
4  get tie get tie trump 's rally #iowa link

1747800 stop laying cmm paris lonnot nott give fuck B...
1747801 ex hñh(u va pu pñc #vri oav touc onaboc tou ...
1747802 l'otan va sortir de sa lèthargie et redevenir ...
1747803 "congiuntinfuoriregione" sono felice per jil...
1747804 ik moet zeggen dat ik Biden the lesser two evi...
Name: text2, Length: 1747805, dtype: object

Sentiment Analysis - what's the attitude of tweeter users toward Trump and Biden?

Explore the Data

In [55]: temp = df_text_target.groupby('target').count().['text'].reset_index().sort_values(by='text',ascending=F
alse)
temp

Out [55]:
   target      text
0  donaldtrump  970919
1  joebiden      776886

In [56]: plt.figure(figsize=(9,6))
sns.countplot(y=df_j_t.country, order = df_j_t.country.value_counts().iloc[:20].index)
plt.title("Top 20 locations")
plt.show()

Top 20 locations

country
United States of America 30000
United States 25000
United Kingdom 10000
India 8000
Germany 7000
France 6000
Canada 5000
Italy 4000
Australia 3000
Mexico 2000
Turkey 1000
The Netherlands 1000
Brazil 1000
Pakistan 1000
Spain 1000
Ireland 1000
Netherlands 1000
Colombia 1000
Argentina 1000
Venezuela 1000
count

Most Common Words

In [57]: from collections import Counter
cnt = Counter()
for text in df_text_target['text'].values:
    for word in text.split():
        cnt[word] += 1
cnt.most_common(10)

Out [57]:
[('trump', 1232839),
 ('link', 1134227),
 ('biden', 818980),
 ('joebiden', 430727),
 ('de', 220338),
 ('election2020', 181076),
 ('donaldtrump', 173666),
 ('vote', 163702),
 ('realnotadtrump', 140039),
 ('la', 130200)]

In [58]: df_text_target['target']=df_text_target['target'].map({'joebiden':0, "donaldtrump":1})

Analysis for Trump

In [59]: def clean(text):
    text = str(text).lower()
    text = re.sub('[\.\?\!\,\']', '', text)
    text = re.sub('https?://\S+|www.\S+', '', text)
    text = re.sub('<.+>', '', text)
    text = re.sub('<.+>', '', text)
    text=re.sub(r'8[A-Za-z0-9+]+',text)
    text=re.sub(r'8+',text)
    text=re.sub(r'(\S+)',text)
    text=re.sub(r'["'+not_list+"']', '', text)
    return text

In [60]: def getSubjectivity(text):
    return TextBlob(text).sentiment.subjectivity
def getPolarity(text):
    return TextBlob(text).sentiment.polarity
def getAnalysis(score):
    if score < 0:
        return 'negative'
    elif score==0:
        return 'neutral'
    else:
        return 'positive'

In [61]: Trump_Tweets = df_j_t.query('target == "donaldtrump"').sort_values('user_followers_count',ascending =
False)

In [62]: Trump_Tweets = Trump_Tweets.dropna().loc[Trump_Tweets.country == 'United States of America']
Trump_Tweets.reset_index(inplace = True, drop = True)

In [63]: Trump_Tweets['ClearTweet'] = Trump_Tweets['text'].apply(clean)

In [64]: from textblob import TextBlob

In [65]: Trump_Tweets['subjectivity']= Trump_Tweets['ClearTweet'].apply(getSubjectivity)
Trump_Tweets['polarity']= Trump_Tweets['ClearTweet'].apply(getPolarity)
Trump_Tweets['analysis']= Trump_Tweets['polarity'].apply(getAnalysis)
Trump_Tweets.head()

Out [65]:
   text      likes      retweet_count      source      user_followers_count      user_location      lat      long      city      country      contin
0  donaldtrump  970919
1  joebiden      776886

In [66]: Without a doubt the #JoeBiden
1  Woah. Have you read the article? No! He...
The latest episode of #SNL tackled #DonaldTrum...
#NatalieMorales explains why some Latine voter...
#NatalieMorales Actress explained why Latine...

0  206.0      12.0      Twitter Web App      5747472.0      Los Angeles, CA  34.053691  -118.242760  Los Angeles      United States of America      Ne Ameri
1  685.0      291.0      Twitter for iPhone      3750110.0      Los Angeles      34.053691  -118.242760  Los Angeles      United States of America      Ne Ameri
2  149.0      15.0      SocialFlow      3264802.0      Hollywood, CA  34.098003  -118.329523  Los Angeles      United States of America      Ne Ameri
3  83.0      10.0      SocialFlow      3264585.0      Hollywood, CA  34.098003  -118.329523  Los Angeles      United States of America      Ne Ameri
4  99.0      21.0      SocialFlow      3264524.0      Hollywood, CA  34.098003  -118.329523  Los Angeles      United States of America      Ne Ameri

In [66]: Trump_Tweets.polarity = Trump_Tweets.polarity.apply(lambda x: getAnalysis(x))
```











```
'kayleighmcenany': 405838,
'post': 581131,
'conservatives': 171531,
'budgetvins': 343106,
'rumpcoupput': 750319,
'notaidtrump2020': 523562,
'kag2020landsidedevictory': 400195,
'breakout': 123112,
'dems': 202990,
'hoax': 342791,
'pomey': 713965,
'died': 212214,
'infect': 363884,
'infected': 363897,
'side': 410570,
'immune': 358025,
'son': 683861,
'ason': 90203,
'stupidity': 701906,
'aplying': 55748,
'potchimout2020': 798865,
'codyosaid': 742626,
'hillarycyclinoton': 340066,
'pointot': 575795,
'organization': 544660,
'shares': 668059,
'talk': 714500,
'pomey': 713965,
'doublestandards': 223940,
'forward': 285990,
'fucking': 292611,
'appen': 329464,
'tried': 746563,
'twice': 763426,
'used': 823113,
'times': 735406,
'absolutely': 31153,
'appene': 339979,
'tries': 746577,
'steal': 695335,
'urges': 775899,
'california': 134604,
'continue': 170172,
'official': 535017,
'top': 227116,
'boxes': 121093,
'spite': 689881,
'threat': 733223,
'yahoo': 826836,
'electionfraud': 241826,
'rumpcriminal': 750444,
'kag2020landsidede': 456869,
'senatemajldr': 662393,
'messing': 487552,
'comple': 739571,
'nasareth': 508694,
'haiprofthedog': 327616,
'howyourmessingwithasnonofabitch': 528284,
'quirlarsson': 523042,
'tonyloami': 740254,
'blackabbath': 113077,
'eddievanhalen': 236323,
'usenate': 779161,
'uscongress': 777510,
'rymes': 530576,
'louisiana': 449467,
'cnn': 143020,
'tbn': 717995,
'pomey': 713965,
'rssist': 624570,
'theview': 730888,
'dnc': 227116,
'dnc': 220105,
'washingtonpost': 806464,
'msnbc': 497484,
'google': 314482,
'bing': 110277,
'nancypelosi': 505649,
'chuckschumer': 134329,
'dobbs': 220386,
'limbaugh': 441094,
'trumpcrump': 377452,
'apology': 66222,
'cenotralparkfive': 144807,
'rumpadministration': 749199,
'giving': 59666,
'trumpisweak': 753436,
'rumpfailure': 751323,
'uspol': 79717,
'knew': 414681,
'trying': 759514,
'information': 364627,
'investigative': 371085,
'journalists': 394773,
'considered': 171696,
'harmful': 339981,
'worthy': 822342,
'silenced': 673134,
'pres': 586206,
'wins': 817626,
'leftists': 433060,
'investigated': 371039,
'leading': 431667,
'balden': 86910,
'uspresidentialelections2020': 778949,
'electionsdialogue': 777815,
'elections2020': 242190,
'electionday': 241674,
'secretlyfilmed': 659407,
'documentary': 220821,
'details': 208395,
'colossal': 163060,
'failure': 268000,
'experts': 264371,
'featured': 273616,
'real': 743671,
'file': 277485,
'clear': 158120,
'disposal': 217315,
'quest': 624895,
'needed': 510203,
'save': 651225,
'hundreds': 349009,
'ousted': 547286,
'role': 636047,
'health': 341199,
'humanservicesdepartmentot': 348591,
'objecting': 532351,
'reponse': 626329,
'front': 291616,
'leadership': 431570,
'lay': 430917,
'wars': 138888,
'push': 600016,
'kinds': 411781,
'sate': 426835,
'elapse': 619623,
'run': 640284,
'basemenot': 50875,
'chiaraferragni': 150091,
'barackobama': 89459,
'trust': 758997,
'presidentot': 587160,
'candidates': 136964,
'straight': 699888,
'play': 573046,
'guess': 322382,
'today': 737922,
'day': 193385,
'sate': 332204,
'level': 436685,
'aburdity': 31373,
'paying': 559144,
'attention': 77418,
'seen': 660025,
'crazy': 180002,
'weeds': 763204,
'obamacare': 531637,
'union': 772221,
'votes': 798367,
'christian': 153495,
'grace': 316939,
'votechmailot': 800090,
'notwinlow': 527224,
...}
```

```
In [85]: import nltk
from nltk.corpus import stopwords
from sklearn.feature_extraction.text import CountVectorizer

stopwords = stopwords.words('english')

print(stopwords)

count_vector = CountVectorizer(token_pattern=r'\w{1,}', ngram_range=(1, 2), stop_words = stopwords)

['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", "your", "yours", "yourself", "yourselves", 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 's', 'what', 'which', 'who', 'whom', 'this', 'that', 'that'll', 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'r', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 'a', 'as', 'at', 'can', 'could', 'will', 'just', 'don', 'don't', 'should', "should've", 'now', 'd', 'll', 'm', 'o', 's', 't', 've', 'y', 'ain', 'aren', 'aren't', 'couldn', 'couldn't', 'didn', 'didn't', 'doesn', 'doesn't', 'hadn', 'hadn't', 'hasn', 'hasn't', 'haven', 'haven't', 'isn', 'isn't', 'ma', 'mightn', "mightn't", 'mustn', 'mustn't', 'needn', 'needn't', 'shan', 'shan't', 'shouldn', "shouldn't", 'wasn', 'wasn't', 'weren', 'weren't', 'won', 'won't', 'wouldn', "wouldn't"]
```

```
In [86]: from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline

clf = LogisticRegression()
pipe = Pipeline([
    ('count_vector', CountVectorizer()),
    ('clf', LogisticRegression())
])

pipe.fit(X_train,y_train)

/usr/local/lib/python3.6/dist-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning:
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
```

```
Out[86]: Pipeline(memory=None,
               steps=[('count_vector',
                       CountVectorizer(analyzer='word', binary=False,
                                       decode_error='strict',
                                       dtype=class 'numpy.int64', encoding='utf-8',
                                       input='content', lowercase=True, max_df=1.0,
                                       max_features=None, min_df=1,
                                       ngram_range=(1, 1), preprocessor=None,
                                       stop_words=None, strip_accents=None,
                                       token_pattern='(?u)\\b\\w+\\b',
                                       tokenizer=None, vocabulary=None)),
                      ('clf',
                       LogisticRegression(C=1.0, class_weight=None, dual=False,
                                           fit_intercept=True, intercept_scaling=1,
                                           l1_ratio=None, max_iter=100,
                                           multi_class='auto', n_jobs=None,
                                           penalty='l2', random_state=None,
                                           solver='lbfgs', tol=0.0001, verbose=0,
                                           warm_start=False))])
```

```
In [87]: from sklearn import metrics
predicted = pipe.predict(X_test)
```

```
In [88]: print('accuracy :', metrics.accuracy_score(predicted, y_test))

accuracy : 0.8300501770787443
```