

Created by Daniel Silalahi (Student Number:21073058)

## Download the Data

```

1 %%capture
2 !pip install spacy
3 !pip install scattertext
4 !pip install tika
5 !pip install spacytextblob
6 !pip install unicode
7
8
9 import spacy
10 import json
11 import pylab
12 from IPython.core.display import display, HTML
13 import nltk
14 from tika import parser
15 import numpy as np
16 import pandas as pd
17 import matplotlib.pyplot as plt
18 from spacytextblob.spacytextblob import SpacyTextBlob
19
20 %matplotlib inline
21 pylab.rcParams['figure.figsize'] = (10., 8.)
22 nlp = spacy.load("en_core_web_sm")
23 nlp.add_pipe('spacytextblob')

```

```

1 #Make a ./data/project directory
2 !mkdir data
3 !mkdir data/project

```

## Insert and process link to data

```
1 !curl 'filename' -o data/project/'filename'
```

% Total	% Received	% Xferd	Average Speed	Time	Time	Time	Current
			Dload	Upload	Total	Spent	Left
0	0	0	0	0	0	0	0

0curl: (6) Could not resolve host: filename

## Install tweepy

```
1 !pip install tweepy
```

```

Requirement already satisfied: tweepy in /usr/local/lib/python3.10/dist-packages (4.14.0)
Requirement already satisfied: oauthlib<4,>=3.2.0 in /usr/local/lib/python3.10/dist-packages (from tweepy) (3.2.2)
Requirement already satisfied: requests<3,>=2.27.0 in /usr/local/lib/python3.10/dist-packages (from tweepy) (2.32.3)
Requirement already satisfied: requests-oauthlib<2,>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from tweepy) (1.3.1)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27.0->tweepy) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27.0->tweepy) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27.0->tweepy) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27.0->tweepy) (2024.8.30)

```

## Twitter Data

```

1 import tweepy
2 import pandas as pd
3
4 consumer_key = "g1HqYDd0uqJrhEfHhpkDRxbUb"
5 consumer_secret = "KbhPjZiLmjFySqWh6WNw5FGFgLJR2NiMhc0VZemfyZXECqJkj"
6 access_key = "1128328751157133313-dtgTxSnjRaTYywoJtm4I8taUjyhs"
7 access_secret = "Zu0nAivKugIvQs4HKm82qSpBhzA0KVkdJjq8NirZ43RLA"
8
9 # Twitter authentication
10 auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
11 auth.set_access_token(access_key, access_secret)
12
13 # Creating an API object
14 api = tweepy.API(auth)
15
16 # create a query for tweets, but exclude retweets and replies
17 tweets = tweepy.Cursor(api.search,
18                         q="roevwade exclude:replies exclude:retweets",
19                         tweet_mode='extended').items()
20
21 list = []
22 for tweet in tweets:
23     text = tweet._json["full_text"]
24
25     refined_tweet = {'text' : text,
26                     'date' : tweet.created_at}
27
28     list.append(refined_tweet)
29
30 df = pd.DataFrame(list)
31 print(df.head())
32

```



```
-----
AttributeError                                Traceback (most recent call last)
<ipython-input-5-c0870c86613f> in <cell line: 17>()
    15
    16 # create a query for tweets, but exclude retweets and replies
--> 17 tweets = tweepy.Cursor(api.search,
    18                        q="roevwade exclude:replies exclude:retweets",
    19                        tweet_mode='extended').items()

AttributeError: 'API' object has no attribute 'search'
```

Clean text (remove emojis or any non ASCII characters) and remove date, as it is only a short time period (1 week)

```
1 import re
2
3 # Create a new dataframe text column
4 df_cleaned = pd.DataFrame(columns=['Text'])
5
6 # Clean up values of df['text'] and insert into df_cleaned['Text']
7 #https://stackoverflow.com/questions/36340627/remove-non-ascii-characters-from-pandas-column
8 df_cleaned['Text'] = df['text'].str.encode('ascii', 'ignore').str.decode('ascii')
9
10 # Lowercase all text to make analysis easier
11 df_cleaned['Text'] = df_cleaned['Text'].str.lower()
12
13
14 # clean other parts of text, such as mentions, hashtags, links, etc
15 # https://docs.python.org/3/library/re.html
16 # https://medium.com/@oscar.sefa/twitter-sentiment-analysis-using-python-for-beginners-1ee1bc15dc86
17 def cleanTweets(text):
18     text = re.sub('@[A-Za-z0-9_]+', '', text) #removes @mentions
19     text = re.sub('#', '', text) #removes hashtag symbol
20     text = re.sub('https?:\/\/\S+', '', text) #Removes links
21     text = re.sub('\n', ' ', text) # Removes any blank lines of space within tweets,
22     return text # makes one big paragraph instead of multiple smaller ones
23
24
25 df_cleaned['Text'] = df_cleaned['Text'].apply(cleanTweets) #apply cleanTweet function to the tweet
26 print('Original Data')
27 print(df.head())
28 print('\nCleaned Data')
29 print(df_cleaned.head())
```



Original Data

	text	date
0	Watching the New Amsterdam episode on RoeVWade...	2023-01-03 18:39:20
1	Check out our episode on #RoeVWade & #Chap...	2023-01-03 18:32:45
2	The same legal arguments used to overturn #R...	2023-01-03 17:45:09

```

3 (OPINION/@tweetmattingly) Churches were always... 2023-01-03 17:29:30
4 States enacted 50 #abortion restrictions in 20... 2023-01-03 16:38:57

```

Cleaned Data

Text

```

0 watching new amsterdam episode feeling physica...
1 check episode roevwade chappelle tomorrow spec...
2 legal arguments used overturn roevwade also th...
3 (opinion/) churches always active abortion deb...
4 states enacted 50 abortion restrictions 2022, ...

```

Create a list that contains full sentences of abortion in tweet

```

1 #https://docs.python.org/3/library/collections.html
2 from collections import Counter
3
4 # Find most common words
5 word_count = Counter(" ".join(df_cleaned["Text"]).split()).most_common(20)
6
7 # Create a dataframe with most common words
8 frequent_words = pd.DataFrame(word_count, columns = ['Word', 'Frequency'])
9
10 print(frequent_words)

```

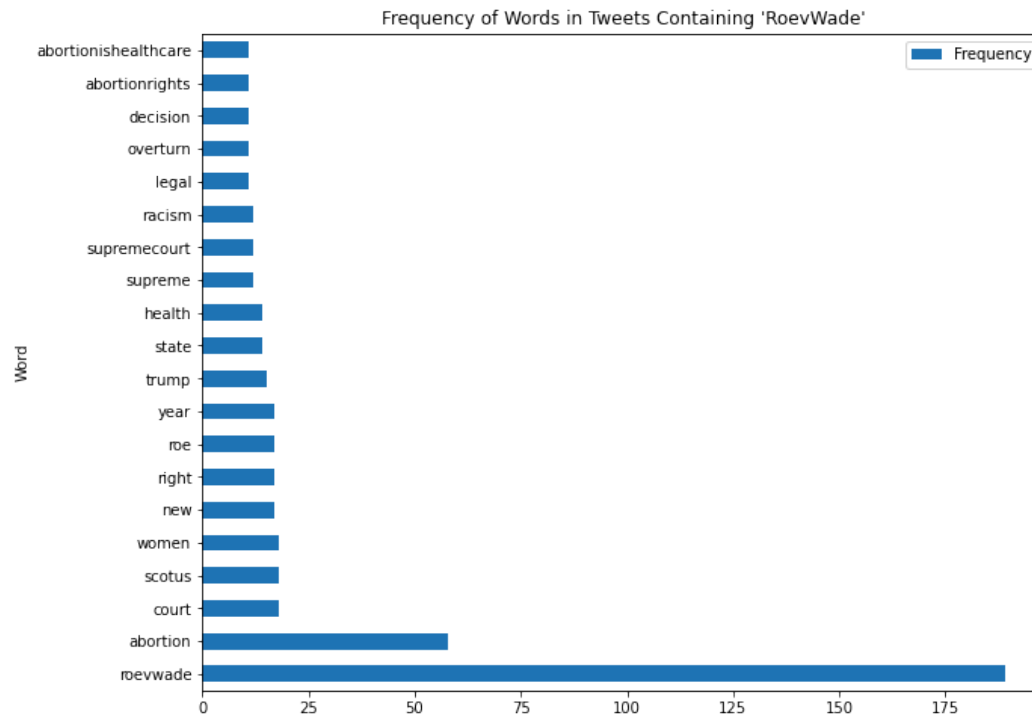
	Word	Frequency
0	roevwade	189
1	abortion	58
2	court	18
3	scotus	18
4	women	18
5	new	17
6	right	17
7	roe	17
8	year	17
9	trump	15
10	state	14
11	health	14
12	supreme	12
13	supremecourt	12
14	racism	12
15	legal	11
16	overturn	11
17	decision	11
18	abortionrights	11
19	abortionishealthcare	11

Creating a bar chart displaying most common words

```

1 #https://matplotlib.org/stable/gallery/lines_bars_and_markers/barh.html
2 import matplotlib.pyplot as plot
3
4 frequent_words.plot.barh(x="Word", y="Frequency", title="Frequency of Words in Tweets Containing 'RoevWade'");
5
6 plot.show(block=True)

```



## Sentiment Analysis

Sentiment analysis is the computational study of people's opinions, sentiments, emotions, appraisals, and attitudes towards entities such as products, services, organizations, individuals, issues, events, topics, and their attributes. Let's study the sentiment of the tweets in this dataset.

spacytextblob performs sentiment analysis using the TextBlob library. Adding spacytextblob to a spaCy nlp pipeline creates a new extension attribute for the Doc.

```

1 #Create a new dataframe for sentiment analysis column
2 df_sentimentAnalysis = pd.DataFrame(columns=['Text'])
3 df_sentimentAnalysis['Text'] = df['text']
4
5 df_sentimentAnalysis['Text'] = df_sentimentAnalysis['Text'].str.encode('ascii', 'ignore').str.decode('ascii')
6 df_sentimentAnalysis['Text'] = df_sentimentAnalysis['Text'].apply(cleanTweets)
7

```

```

8 df_sentimentAnalysis['Polarity'] = df_sentimentAnalysis['Text'].apply(lambda x: nlp(x)._.blob.polarity)
9 df_sentimentAnalysis['Subjectivity'] = df_sentimentAnalysis['Text'].apply(lambda x: nlp(x)._.blob.subjectivity)
10
11 # Function to return positive, neutral, negative according to the polarity
12 def analysis(score):
13     if score>0:
14         return 'Positive'
15     elif score ==0:
16         return 'Neutral'
17     else:
18         return 'Negative'
19
20 #apply function to the 'Polarity' column and store values in a new column named 'Analysis'
21 df_sentimentAnalysis['Analysis'] = df_sentimentAnalysis['Polarity'].apply(analysis)
22
23 # create a new df for subjectivity >= 0.5
24 df_subjective = df_sentimentAnalysis[df_sentimentAnalysis['Subjectivity'] >= 0.5]
25
26 print('Sentiment Analysis')
27 print(df_sentimentAnalysis.head())
28 print('\nSentiment Analysis for Subjectivity>=0.5')
29 print(df_subjective.head())

```



```

-----
NameError                                Traceback (most recent call last)
<ipython-input-1-0c37ea5a0996> in <module>
      1 #Create a new dataframe for sentiment analysis column
----> 2 df_sentimentAnalysis = pd.DataFrame(columns=['Text'])
      3 df_sentimentAnalysis['Text'] = df['text']
      4
      5 df_sentimentAnalysis['Text'] = df_sentimentAnalysis['Text'].str.encode('ascii', 'ignore').str.decode('ascii')

NameError: name 'pd' is not defined

```

Create a pie chart of Polarity

```

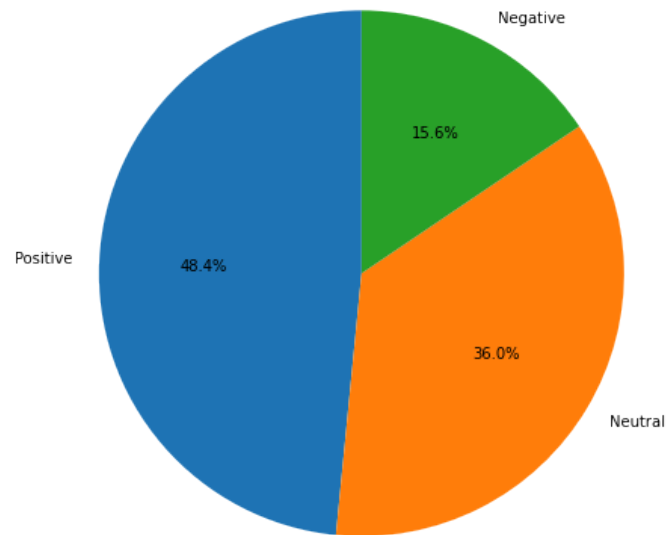
1 #https://matplotlib.org/stable/gallery/pie_and_polar_charts/pie_features.html
2
3
4 # Pie chart for analysis:
5 pie_chart_values = df_sentimentAnalysis['Analysis'].value_counts()
6 print(pie_chart_values)
7 pie_chart_values.plot(kind = 'pie', label='', autopct='%1.1f%%', startangle=90)
8 plt.title('Pie Chart for Polarity of All Tweets')
9 plt.show()
10
11 # create a pie chart for analysis but only for subjectivity >= 0.5
12 pie_chart_values = df_subjective['Analysis'].value_counts()

```

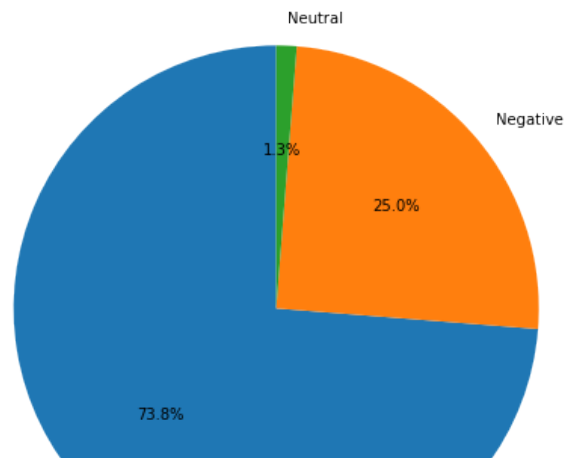
```
13 print(pie_chart_values)
14 pie_chart_values.plot(kind = 'pie', label='', autopct='%1.1f%%', startangle=90)
15 plt.title('Pie Chart for Polarity of Tweets with Subjectivity >= 0.5')
16 plt.show()
17
```

```
Positive    109  
Neutral     81  
Negative     35  
Name: Analysis, dtype: int64
```

Pie Chart for Polarity of All Tweets



```
Positive    59  
Negative    20  
Neutral      1  
Name: Analysis, dtype: int64
```

Pie Chart for Polarity of Tweets with Subjectivity  $\geq 0.5$ 

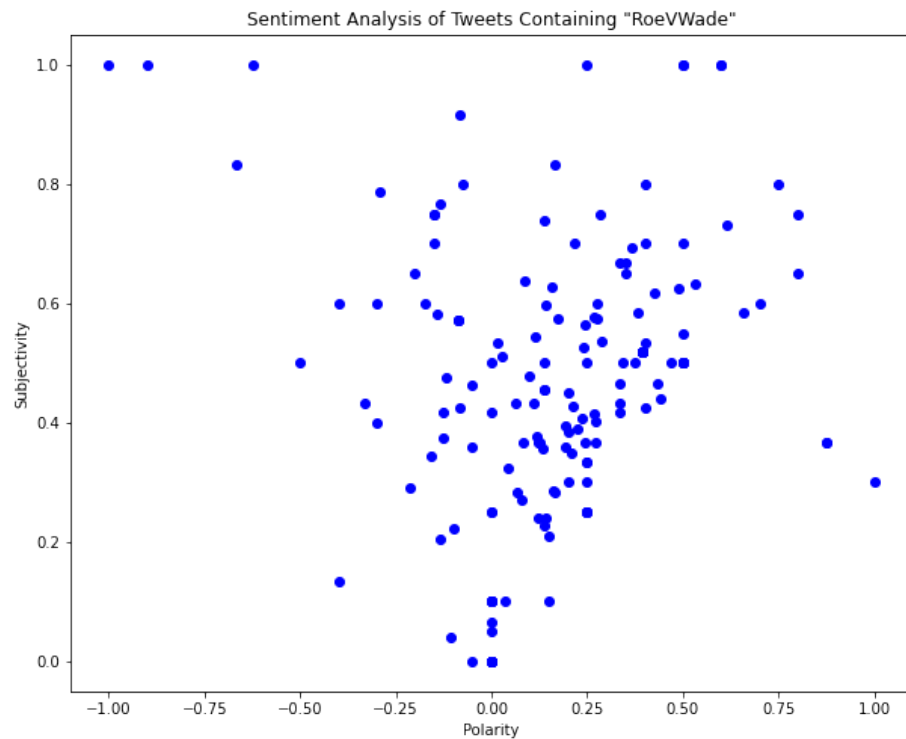


Positive



### Creating a Scatterplot of Polarity and Subjectivity

```
1 #create a loop that repeats for amount of rows the df_sentimentAnalysis has
2 for row in range(df_sentimentAnalysis.shape[0]):
3     plt.scatter(df_sentimentAnalysis['Polarity'][row], df_sentimentAnalysis['Subjectivity'][row], color='blue')
4 plt.title('Sentiment Analysis of Tweets Containing "RoeVWade"')
5 plt.xlabel('Polarity')
6 plt.ylabel('Subjectivity')
7 plt.show()
```



To get a sense of what the nlp model defines as high or low polarity, we will print the 5 tweets of lowest polarity and 5 tweets with highest polarity

```
1 # get the 5 lowest polarity values
2 lowest5 = df_subjective.nsmallest(5, 'Polarity')
3
```