Introduction As a resource for social data, Twitter's platform has been used to measure the quality of life through sentiment analysis. This capstone project explores another methodological technique of using specific keyword terms to determine dominant topics, word patterns, and sentiment leanings in a geographical area. Focusing on New York City and Los Angeles for comparative analysis, the keyword term "why" will be used to build a Python analysis around topic modeling and sentiment analysis. With this approach, the analysis reveals social and cultural differences, the overall sentiment of tweets, and areas of interest to tweeters. Contents 1. Install Libraries 2. Import Python Libraries 3. Data Setup (import, query, convert JSON to DataFrame, and clean Twitter data) 4. Categorizing Sentiment on Tweets 5. Exploratory Analysis 6. Topic Modeling 7. Topic Bubble Map 8. Accuracy of Sentiment Analysis **Install libraries** pip install searchtweets-v2pip install nltknltk.download('stopwords')nltk.download('punkt')#to import nltk SentimentIntensityAnalyzer() nltk.download('vader_lexicon')pip install textblobpip install contractionspip install wordcloudpip install sklearnpip install plotlypip install geopandas **Python Libraries** In [1]: import searchtweets as twitter import pandas as pd from pandas.io.json import json normalize import numpy as np import re import contractions import string import textwrap import matplotlib.pyplot as plt from wordcloud import WordCloud, STOPWORDS from PIL import Image import plotly.graph objs as go import plotly.express as px import geopandas as gpd import nltk from nltk.corpus import stopwords from nltk.probability import FreqDist from nltk.stem import PorterStemmer, WordNetLemmatizer from nltk.sentiment.vader import SentimentIntensityAnalyzer from nltk.tokenize import word tokenize from nltk.util import ngrams from textblob import TextBlob import sklearn from sklearn import metrics from sklearn.feature extraction.text import TfidfVectorizer from sklearn.decomposition import LatentDirichletAllocation from sklearn.preprocessing import normalize from sklearn.model selection import train test split from sklearn.linear model import LogisticRegression from sklearn.naive bayes import MultinomialNB from sklearn.metrics import confusion_matrix, classification_report import warnings from warnings import simplefilter warnings.filterwarnings('ignore') simplefilter(action='ignore', category=FutureWarning) **Data Setup Import Twitter Data** In [2]: #file contains API and Bearer tokens search_args = twitter.load_credentials('twitter_keys.yaml', yaml_key='search_tweets_v2', env overwrite=False) **Query Twitter Data** In [3]: query whynyc = twitter.gen request parameters(query = 'why (place:01a9a39529b27f36 OR place:011add077f4d2da3 OR place:00c39537733fa112 OR place:002e24c67 results_per_call = 500, start time = '2021-01-01', end time = '2022-01-01', tweet_fields = 'id, created_at, text, geo', granularity='' #140k whynyc tweets = twitter.collect results(query_whynyc, $max_tweets=140000,$ result_stream_args=search_args In [4]: query_whyla = twitter.gen_request_parameters(query = 'why (place:3b77caf94bfc81fe) -is:retweet -is:nullcast lang:en', results_per_call = 500, start_time = '2021-01-01', end time = '2022-01-01', tweet fields = 'id, created at, text, geo', granularity='' whyla tweets = twitter.collect results(query_whyla, max tweets=90000, result stream args=search args **Convert JSON to Dataframe** whynyc df = pd.json normalize(whynyc tweets, record path=['data']) In [5]: whyla df = pd.json normalize(whyla tweets, record path=['data']) **Data Cleaning** In [6]: def find_links(tweet): #function extracts the links return re.findall('(http\S+|bit.ly/\S+)', tweet) #def find retweeted(tweet): # function finds and extracts retweeted twitter handles #return re.findall('(?<=RT\s)(@[A-Za-z0-9]+[A-Za-z0-9-_]+)', tweet) def find mentioned(tweet): #function finds and extracts the twitter handles of people mentioned return re.findall('(?<!RT\s)(@[A-Za-z0-9]+[A-Za-z0-9-]+)', tweet) def find_hashtags(tweet): #This function will extract hashtags **return** re.findall('(#[A-Za-z0-9]+[A-Za-z0-9-_]+)', tweet) In [7]: # make new columns for retweeted usernames, mentioned usernames and hashtags whynyc df['links'] = whynyc df.text.apply(find links) #whynyc_df['retweeted'] = whynyc_df.text.apply(find_retweeted)
whynyc_df['mentioned'] = whynyc_df.text.apply(find_mentioned) whynyc_df['hashtags'] = whynyc_df.text.apply(find_hashtags) whyla_df['links'] = whyla_df.text.apply(find_links) #whyla_df['retweeted'] = whyla_df.text.apply(find_retweeted) whyla_df['mentioned'] = whyla_df.text.apply(find_mentioned) whyla_df['hashtags'] = whyla_df.text.apply(find_hashtags) In [8]: #removing links, mentions, hashtags from corpora column stopwords = nltk.corpus.stopwords.words('english') lemmatizer = WordNetLemmatizer() punctuation = string.punctuation #'!'\$%&\'()*+,-./:;<=>?[\\]^``{|}~•@' symbol = '-...«>""'' #for symbols not captured in punctuation def clean_df(tweet): #remove parts of a tweet tweet = re.sub(r'http\S+', '', tweet) tweet = re.sub(r'bit.ly/\S+', '', tweet) $\#tweet = re.sub('(RT \setminus S@[A-Za-z0-9]+[A-Za-z0-9-_]+)', '', tweet)$ tweet = $re.sub('(@[A-Za-z0-9]+[A-Za-z0-9-_]+)', '', tweet)$ #tweet = re.sub('(#[A-Za-z0-9]+[A-Za-z0-9-]+)', '', tweet)#removing these that showup after data cleaning processing tweet = re.sub('&', '&', tweet) tweet = re.sub('\n', '', tweet) #lower-case characters tweet = tweet.lower() #remove contractions tweet = contractions.fix(tweet) #remove numbers tweet = re.sub('([0-9]+)', '', tweet)#remove punctuation tweet = re.sub('['+ string.punctuation +']+', ' ', tweet) #remove symbols not captured in punctuation tweet = re.sub('['+ symbol +']+', ' ', tweet) #remove whitespace tweet = $re.sub(r'^\s+|\s+$', '', tweet)$ tweet = $re.sub(r'\s+', '', tweet)$ #tokenize words and remove stopwords tweet token list = [word for word in tweet.split(' ')#] if word not in stopwords] # remove stopwords #apply word lemmatization tweet_token_list = [lemmatizer.lemmatize(word) if '#' not in word else word for word in tweet_token_list] tweet = ' '.join(tweet_token_list) return tweet #create a new column for the cleaned text column. whynyc_df['corpora'] = whynyc_df.text.apply(clean_df) whyla_df['corpora'] = whyla_df.text.apply(clean_df) In [9]: #pull list of columns list(whynyc df) Out[9]: ['id', 'text', 'created at', 'geo.place id', 'geo.coordinates.type', 'geo.coordinates.coordinates', 'withheld.copyright', 'withheld.country codes', 'withheld.scope', 'links', 'mentioned', 'hashtags', 'corpora'] In [10]: #reorder columns in dataframe whynyc = whynyc_df[['id', 'created_at', 'text', 'corpora', 'geo.place_id']] whyla = whyla_df[['id', 'created_at', 'text', 'corpora', 'geo.place id']] #created for a one-time analysis #adding this to whynyc and whyla dataframes may cause kernel to crash due to size whynyc pot = whynyc df[['mentioned', 'hashtags', 'links']] whyla pot = whyla df[['mentioned', 'hashtags', 'links']] **Categorizing Sentiment on Tweets** In [11]: sid = SentimentIntensityAnalyzer() In [12]: def add sentiment(why df): #pulling polarity scores from NLTK library sa nltk list = [] for i in why df['text']: sa_nltk_list.append((sid.polarity_scores(str(i)))['compound']) #why df['score'] = pd.Series(sa nltk list, dtype='float64') #pulling subjectivity and polarity scores from TextBlob def subjectivity(text): return TextBlob(text).sentiment.subjectivity why df['subjectivity'] = why df['text'].apply(subjectivity) #Create a function to get the polarity sa tb list = [] def polarity(text): return TextBlob(text).sentiment.polarity sa_tb_list = why_df['text'].apply(polarity) #average of NLTK and TextBlob's polarity scores via Numpy avg = np.mean(np.array([sa nltk list, sa tb list]), axis=0) why df['score'] = pd.DataFrame(avg) #Categorizing sentiment scores def sentiment category(sentiment): label = '' if(sentiment>0): label = 'positive' elif(sentiment == 0): label = 'neutral' else: label = 'negative' why_df['sentiment'] = why_df['score'].apply(sentiment_category) return why df whynyc = add sentiment(whynyc) whyla = add sentiment(whyla) **Exploratory Analysis** Due to links, mentions, and hahtags accounting for less than 1/3 of the total tweets queried, this portion will only provide a basic idea of how much relevance the parts of a tweet (below) play a role. In a future project, a network analysis will come into play for this part.discourse. 1. Links (either to an image or a website) 2. Mentioned 3. Hashtags % of tweets (NYC/LA) **New York City** Los Angeles Links accounts for 41,566 tweets 24,497 tweets (30%/27%) Mentioned tweeters account for 46,067 tweets 28,535 tweets (33%/32%) 9,180 tweets Hashtags account for 5,537 tweets (7%/6%)len(whynyc) In [13]: 138886 Out[13]: len(whyla) In [14]: 89344 Out[14]: whynyc_pot['links'].value_counts() In [15]: 97361 Out[15]: [https://t.co/CiyzXjOgTz] 16 13 [https://t.co/UW7TkMtUIM] [https://t.co/SPAVAK2eTW] 12 [https://t.co/Ef0jFoFAbo] 11 [https://t.co/D3QdYw2qGe] [https://t.co/slR3zFrGH5] [https://t.co/9ttmJbtdm5] [https://t.co/dPxLlQNyXM] [https://t.co/n0Mf9AhPBK] Name: links, Length: 41321, dtype: int64 In [16]: whynyc_pot['mentioned'].value counts() 79110 [] Out[16]: [@NYCTSubway] 312 [@YouTube] 174 [@MTA] 129 [@nypost] 101 [@VP, @lizshuler, @AFLCIO] 1 [@PaulMcG1994, @DiMarco Pattaya] 1 [@operaqueenie] 1 [@dcboyisangry, @jondunnsays, @JoyceWhiteVance] 1 [@cbaibix] Name: mentioned, Length: 45930, dtype: int64 whynyc pot['hashtags'].value counts() In [17]: 126520 [] Out[17]: [#RHOP] [#Yankees] 65 [#LGM] 62 [#RHOA] 54 [#jasonblack] [#AmcManhattan] [#WWNXT, #TakeOver] [#RecallRonDeSantis, #RecallAbbott] [#Stupidity] Name: hashtags, Length: 9164, dtype: int64 In [18]: whyla pot['links'].value counts() 64696 Out[18]: [https://t.co/H4JxLhbZow] 49 [https://t.co/QNaC8UUGfn] 19 [https://t.co/S5YKpBwbVm] [https://t.co/Ob6sf4xrpK] [https://t.co/r4To8Aon3G] [https://t.co/foF3Df4BcH] [https://t.co/91jGhMWQyl] [https://t.co/BRiGLEVktk, https://t.co/6zIFe7BZjU] [https://t.co/dqnarpbqIs] Name: links, Length: 24498, dtype: int64 In [19]: whyla pot['mentioned'].value counts() 54450 [] Out[19]: [@YouTube] 188 [@GavinNewsom] 44 [@Dodgers] 43 [@thehill] 43 [@Nothing22280589, @londongal 28, @OccupyDemocrats] [@amandaapittman] [@RealRaider2055] [@mikejschaefer] [@BookSyrup] Name: mentioned, Length: 28525, dtype: int64 In [20]: whyla pot['hashtags'].value counts() 82567 [] Out[20]: [#Dodgers] 69 48 [#FreeBritney] [#RHOBH] [#BB23] 25 [#COVID, #vaxxing, #masking, #science] [#OscarIsaacs, #foxstudios, #xmenmovies] [#DemonsSouls, #ps5] [#ClimateAction, #StandWithGavin, #VoteNoOnRecall] [#pointlessskip] Name: hashtags, Length: 5539, dtype: int64 **Text Analysis** In [21]: **def** format why df(why df): #convert created at column to datetime why df['datetime'] = pd.to datetime(why df['created at'], errors='coerce') #create a day column why df['day'] = why df['datetime'].dt.date #create a month column why df['month'] = why df['datetime'].dt.month #break up text column into length why df['length']=why df['text'].apply(lambda x:len(x.split())) return why df whynyc = format why df(whynyc) whyla = format why df(whyla)In [22]: whynyc['length'].describe() count 138886.000000 Out[22]: 22.437315 mean std 14.227934 1.000000 min 25% 11.000000 18.000000 50% 75% 31.750000 109.000000 Name: length, dtype: float64 whyla['length'].describe() 89344.000000 count Out[23]: 21.401986 mean std 13.847466 min 1.000000 25% 11.000000 50% 17.000000 75% 29.000000 109.000000 max Name: length, dtype: float64 In [24]: sentiment_colors = { 'positive': '#2A9D8F', 'neutral': '#847979', 'negative': '#F4A259'} px.histogram(whynyc, x='length', color='sentiment', color_discrete_map = sentiment_colors) sentiment 6000 positive negative neutral 5000 4000 3000 2000 1000 20 80 100 40 length px.histogram(whyla, x='length', color='sentiment', color discrete map = sentiment colors) sentiment positive 4000 negative neutral 3500 3000 2500 2000 1500 1000 500 100 length In [27]: #convert dataframe to lists whynyc_list = whynyc['corpora'].values.tolist() whynyc_list = ' '.join(whynyc_list).lower() whyla_list = whyla['corpora'].values.tolist() whyla_list = ' '.join(whyla_list).lower() In [28]: #create a frequency distribution and graph it fdist whynyc = FreqDist(word tokenize(whynyc list)) plt.figure(figsize=(10, 4)) fdist_whynyc.plot(30, cumulative=False) plt.show() 20000 18000 16000 14000 12000 10000 8000 6000 pinow know peed think love year thing make Samples In [29]: fdist_whyla = FreqDist(word_tokenize(whyla_list)) plt.figure(figsize=(10, 4)) fdist_whyla.plot(30,cumulative=False) plt.show() 14000 12000 10000 8000 6000 4000 2000 plood g Samples **Text Analysis - Word Clouds** def create_wordcloud(file_name, list_name): #pull the image file mask = np.array(Image.open(file_name)) #function converts RGB values from 0 (black) to white (255) def transform_zeros(val): **if** val == 0: return 255 else: return val #map and create a mask for image maskable_image = np.ndarray((mask.shape[0],mask.shape[1]), np.int32) for i in range(len(mask)): maskable image[i] = list(map(transform zeros, mask[i])) wordcloud = WordCloud(width = 3000, height = 2000,#random_state=1, background_color='white', colormap='twilight_r', contour_width = 1, contour_color = '#111954', collocations=True, stopwords = STOPWORDS, mask=maskable_image) .generate(list_name) def plot_cloud(wordcloud): # Set figure size plt.figure(figsize=(15, 7)) # Display image plt.imshow(wordcloud) # No axis details plt.axis('off'); return plot_cloud(wordcloud) create wordcloud('new-york-city.png', whynyc list) create_wordcloud('los-angeles.png', whyla_list) go said N-grams def get ngrams(text, ngram from=2, ngram to=2, n=None, max features=20000): vectorizer = TfidfVectorizer(ngram_range = (ngram_from, ngram_to), max_features = max_features, stop words='english').fit(text) bag of words = vectorizer.transform(text) sum_words = bag_of_words.sum(axis = 0) words_freq = [(word, sum_words[0, i]) for word, i in vectorizer.vocabulary_.items()] words_freq = sorted(words_freq, key = lambda x: x[1], reverse = True) return words_freq[:n] In [32]: def ngrams_table(why_df): unigrams = pd.DataFrame(get_ngrams(why_df['corpora'], ngram_from=1, ngram_to=1, n=15)) bigrams = pd.DataFrame(get_ngrams(why_df['corpora'], ngram_from=2, ngram_to=2, n=15)) trigrams = pd.DataFrame(get_ngrams(why_df['corpora'], ngram_from=3, ngram_to=3, n=15)) quadgrams = pd.DataFrame(get_ngrams(why_df['corpora'], ngram_from=4, ngram_to=4, n=15)) ngrams = pd.concat([unigrams, bigrams, trigrams, quadgrams], axis = 1) ngrams.columns = ['unigrams', 'frequency', 'bigrams', 'frequency', 'trigrams', 'frequency', 'quadgrams', 'f ngrams return ngrams ngrams_nyc = ngrams_table(whynyc) ngrams_la = ngrams_table(whyla) ngrams nyc In [33]: Out[33]: trigrams frequency unigrams frequency bigrams frequency quadgrams frequency 0 3529.134437 151.470388 like look like 663.753302 new york city new york new york 156.424098 news network elected 2155.107287 new york 642.976160 new york new 115.670007 31.975772 people official network elected official 507.129674 2 1956.359802 feel like york new york 112.996076 31.975772 know elected official silent 1510.708458 31.975772 398.085112 brooklyn new york 87.817685 want make sense obvious official silent obvious 1222.646746 31.975772 lol want know 314.946995 make make sense 78.753550 miscarriage silent obvious miscarriage 5 1179.093271 215.473221 56.847758 31.975772 time people like manhattan new york obvious miscarriage justice 6 1170.154142 36.189982 31.975772 year old 214.833320 got gt gt gt miscarriage justice social need 1164.732432 sound like 195.317658 idk feel like 33.265314 31.975772 security 8 think 1153.762136 year ago 187.132847 really want know 31.393852 justice social security irs 31.975772 social security irs social medium 178.868629 1112.067675 news network elected 31.205158 31.975772 make administration security irs administration 10 1060.079614 black people 160.050376 network elected official 31.205158 31.975772 understand say 1031.997556 159.638005 31.205158 irs administration ssi veteran elected official silent people administration ssi veteran 12 1013.232131 135.994399 official silent obvious 31.205158 31.975772 deserve silent obvious 943.120180 133.882774 31.975772 13 31.205158 ssi veteran deserve date reason acting like miscarriage obvious miscarriage veteran deserve date expect 14 good 901.036878 people think 133.700021 31.205158 31.975772 justice In [34]: ngrams la Out[34]: unigrams frequency bigrams frequency trigrams frequency quadgrams frequency 0 2499.620279 55.803544 like look like 423.889082 los angeles california 248.712247 It It It It graphically audiovisual cosmos graphically 1 1344.985599 feel like 353.107215 100.248971 32.907557 people face race audiovisual 2 306.655385 know 1303.699416 los angeles It It It 57.832755 audiovisual face race age 32.907557 3 lol 1107.359648 make sense 237.456033 make make sense 55.159320 face race age nationality 32.907557 graphically audiovisual 4 1004.126975 angeles california 178.055481 31.199730 32.907557 want race age nationality exact age nationality exact 5 862.895499 want know 165.832352 audiovisual face race 31.199730 32.907557 got location cosmos graphically 6 781.008396 145.078652 31.199730 32.700657 people like face race age audiovisual face nationality exact location 7 27.344379 31.199730 think 755.096647 sound like 138.559658 race age nationality everybody los angeles hollywood 8 748.275575 131.546071 age nationality exact 31.199730 13.905531 time year old california nationality exact 9 31.199730 11.941129 make 727.414897 social medium 118.785631 al haqq nur graphically location haqq nur graphically exact location 10 idk 715.780248 year ago 110.466998 25.998278 11.941129 everybody audiovisual understand cosmos graphically 24.236062 98.619342 idk feel like 11.901923 11 shit 713.766732 audiovisual body people graphically 12 683.379039 88.351261 guest catch live weeknight 9.367749 need gt gt gt 19.437893 audiovisual 13 say 670.838977 people think 83.344778 make feel like 17.547511 catch live weeknight et 9.199203 cosmos 14 going 659.337095 82.137094 today feel like 17.176137 south los angeles california 9.098979 graphically **Sentiment EDA** In [35]: whynyc['score'].describe() 138886.000000 Out[35]: 0.030624 0.331123 -0.991750 25% -0.186600 50% 0.000000 75% 0.248782 0.982300 Name: score, dtype: float64 whynyc['sentiment'].describe() In [36]: 138886 Out[36]: unique 3 positive top 62204 Name: sentiment, dtype: object fig, ax = plt.subplots() whynyc['sentiment'].value counts().plot(ax=ax, kind='bar', xlabel='numbers', ylabel='frequency') 60000 50000 40000 30000 20000 10000 0 positive negative neutra numbers In [38]: fig, ax = plt.subplots() whyla['sentiment'].value counts().plot(ax=ax, kind='bar', xlabel='numbers', ylabel='frequency') plt.show() 40000 35000 30000 25000 20000 15000 10000 5000 positive negative numbers In [39]: **def** pol and sub of tweets(title, why df): # plot the polarity and subjectivity fig = px.scatter(why df,x='score', y='subjectivity', color = 'sentiment', color discrete map = sentiment colors, size='subjectivity', hover_name=why_df.text.apply(lambda txt: '<br'.join(textwrap.wrap(txt, width=35))))</pre> #add a vertical line at x=0 for Netural Reviews fig.update layout(title=title, shapes=[dict(type= 'line', yref= 'paper', y0= 0, y1= 1, xref= 'x', x0= 0, x1= 0)])return fig.show() pol and sub of tweets ('Subjectivity and Polarity Scores of Tweets in NYC', whynyc) pol and sub of tweets ('Subjectivity and Polarity Scores of Tweets in LA', whyla)

