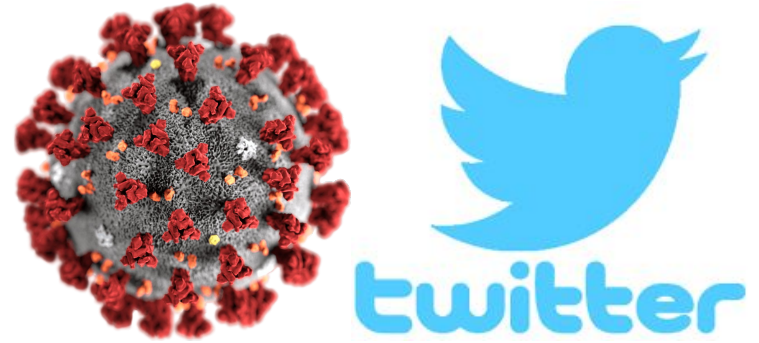


#Covid19



Investigating the relationship between social media activity and the coronavirus

Group 6: Andrew Cummings,
Oliver Posewitz, Ian Ustanik

Agenda

- Data Acquisition
 - Twitter
 - Covid-19
- Data Preprocessing
 - Compilation
 - Cleansing
 - Aggregation
- Methods of Analysis
 - Statistical
 - NLP - Sentiment
 - Regression/Location

```
1 def tweet_search(hashtag):
2     query = hashtag
3     count = 1000
4     searched_tweets = [status for status in tweepy.Cursor(api.search, q=query).items(count)]
5
6     tweet_list = list(searched_tweets)
7     tweet_json = [tweet._json for tweet in tweet_list]
8     return tweet_json
9
10 def tweets_to_df(tweet_json):
11     ID = []
12     tweets = []
13     location = []
14     time = []
15
16     for tweet in tweet_json:
17         ID.append(tweet['id'])
18         tweets.append(tweet['text'])
19         location.append(tweet['user']['location'])
20         time.append(tweet['created_at'])
21
22     data_tuples = list(zip(ID, tweets, location, time))
23     tweets_df = pd.DataFrame(data_tuples, columns=['ID', 'tweets', 'location', 'time'])
24     return tweets_df
```

Data Acquisition

Twitter Data

- Tweepy request to Twitter API
- Use three 'bins' of keywords:
 - #covid19
 - #maskup
 - #travel
- Request limitations → Pagination
- Retrieved ID, tweets, location, and time
- 21,000 tweets total over 7-day period

	ID	tweets	location	time
0	1328125627812114433	RT @ArTalks: GodSpeed 🙏🙏🙏🙏🙏🙏🙏🙏 @SpaceX ...	Los Angeles Universe Omniverse	Sun Nov 15 23:59:59 +0000 2020
1	1328125567271538688	As if we needed evidence that the Pandemic is ...	#notinabubble	Sun Nov 15 23:59:45 +0000 2020
2	1328125516222754818	RT @nonsequiteuse: Hey @HiltonHouston & @H...	Houston, Texas	Sun Nov 15 23:59:33 +0000 2020
3	1328125492545798144	RT @cartoonnetwork: Rep your fav #CartoonNetwo...	someone's mind	Sun Nov 15 23:59:27 +0000 2020
4	1328125490255835136	Waiting to take a pandemic seriously until you...	Norman, OK	Sun Nov 15 23:59:26 +0000 2020
...
995	1328063679926038531	RT @Alt_ReddTruq: #MaskUp	Florida, USA	Sun Nov 15 19:53:50 +0000 2020
996	1328063582974734339	Do you think racist folks will finally get #BL...	Texas, USA	Sun Nov 15 19:53:27 +0000 2020
997	1328063420739031043	RT @cartoonnetwork: Rep your fav #CartoonNetwo...	Naahhhh, have a cawfee -w-	Sun Nov 15 19:52:48 +0000 2020
998	1328063393438322688	For those of you that need a reminder. #maskup...	New York, NY	Sun Nov 15 19:52:41 +0000 2020
999	1328063376891801600	RT @cartoonnetwork: Rep your fav #CartoonNetwo...		Sun Nov 15 19:52:37 +0000 2020

Covid-19 Data

- Center for Disease Control website



Data Preprocessing

Cleansing and Purging

- Used conditionals to tag records with location data like state name in Twitter data, formatted time attribute
- Concatenate each bin of tweets into one dataframe

location	location	location
Universe Omniverse	Miami, FL	FL
#notinabubble	Houston, TX	TX
Houston, Texas	Montville, NJ	NJ
someone's mind	MD	MD
Norman, OK	The Woodlands, TX	TX

#stitch csv's together

```
import os

directory = 'travel'

li = []

for csv_file in os.listdir(directory):
    mask_df = pd.read_csv(directory+'/' + csv_file, index_col=None, header = 0)
    li.append(mask_df)

travel_df = pd.concat(li, axis=0, ignore_index=True)

# travel_df.to_csv('travel.csv')
```

Combining Different Data Sets

- Data was initially stored in many separate csv's
- Iterate over csv's using os package to create one file for each bin

covid1913	11/20/2020 3:43 PM	Microsoft Excel C...	186 KB
covid1914	11/20/2020 3:43 PM	Microsoft Excel C...	187 KB
covid1915	11/20/2020 3:43 PM	Microsoft Excel C...	185 KB
covid1916	11/20/2020 3:43 PM	Microsoft Excel C...	186 KB
covid1917	11/20/2020 3:43 PM	Microsoft Excel C...	189 KB
covid1918	11/20/2020 3:43 PM	Microsoft Excel C...	189 KB
covid1919	11/20/2020 3:43 PM	Microsoft Excel C...	190 KB
mask13	11/20/2020 3:41 PM	Microsoft Excel C...	194 KB
mask14	11/20/2020 3:41 PM	Microsoft Excel C...	198 KB
mask15	11/20/2020 3:41 PM	Microsoft Excel C...	191 KB
maskup16	11/20/2020 3:41 PM	Microsoft Excel C...	184 KB
maskup17	11/20/2020 3:41 PM	Microsoft Excel C...	190 KB
maskup18	11/20/2020 3:41 PM	Microsoft Excel C...	193 KB
maskup19	11/20/2020 3:41 PM	Microsoft Excel C...	188 KB
travel13	11/20/2020 3:44 PM	Microsoft Excel C...	203 KB
travel14	11/20/2020 3:44 PM	Microsoft Excel C...	205 KB
travel15	11/20/2020 3:44 PM	Microsoft Excel C...	205 KB
travel16	11/20/2020 3:44 PM	Microsoft Excel C...	206 KB
travel17	11/20/2020 3:44 PM	Microsoft Excel C...	208 KB
travel18	11/20/2020 3:44 PM	Microsoft Excel C...	186 KB
travel19	11/20/2020 3:44 PM	Microsoft Excel C...	184 KB

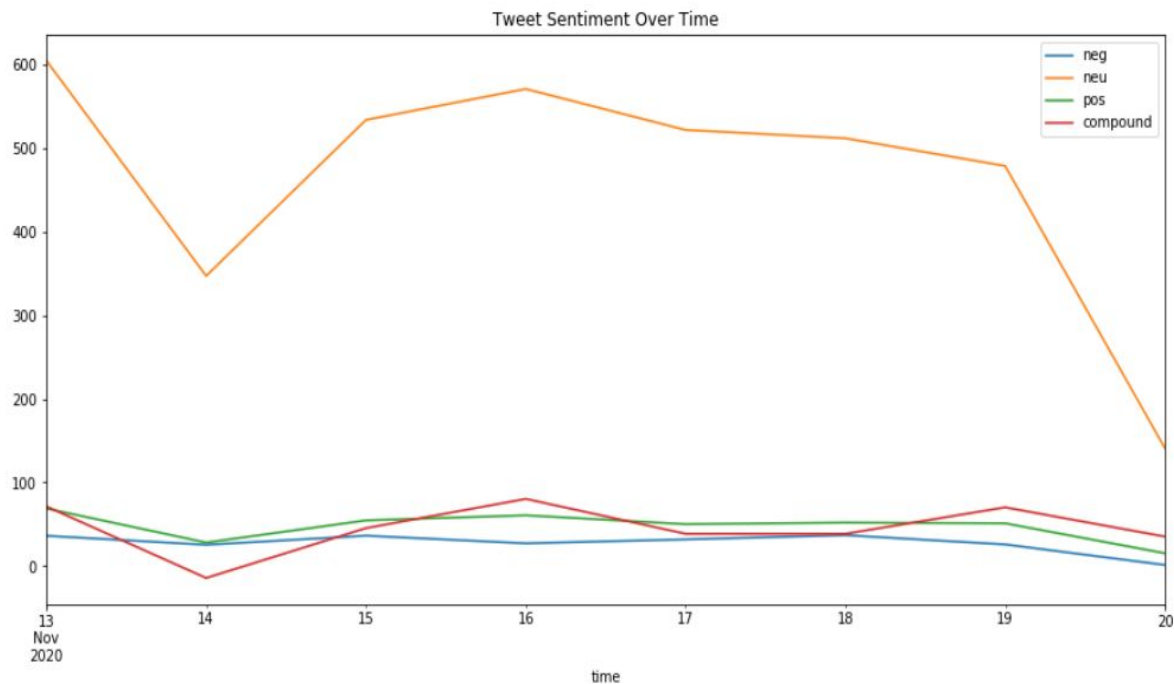
all



11/20/2020 5:05 PM Microsoft Excel Co... 4,222 KB

Aggregation and Summarization

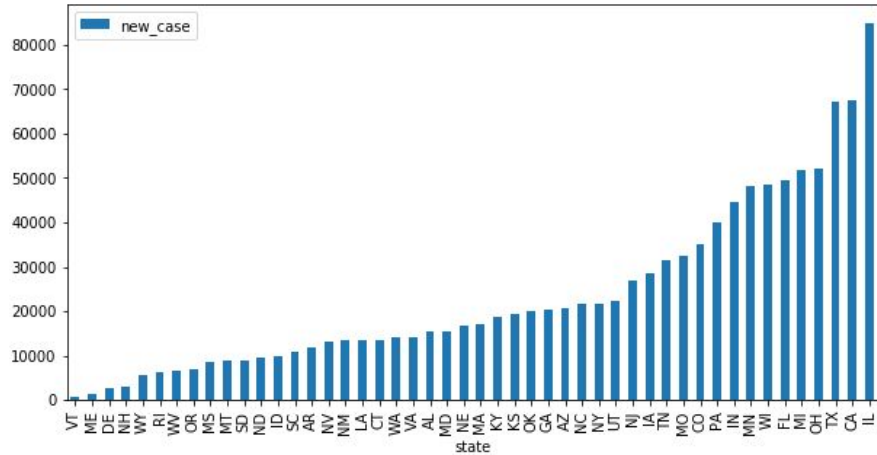
- Dimension:
 - 21,000 rows
 - 9 columns
- Continental US
 - Lower 48
- 7-day period
 - (11/13-11/19)
- 4,314 tweets
 - Location formatted



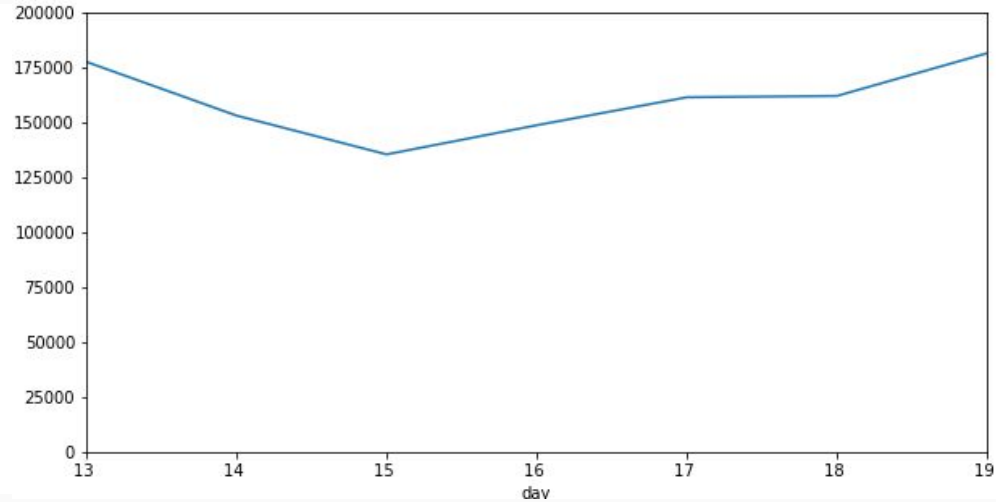
Methods of Analysis

What is the current state of Covid-19?

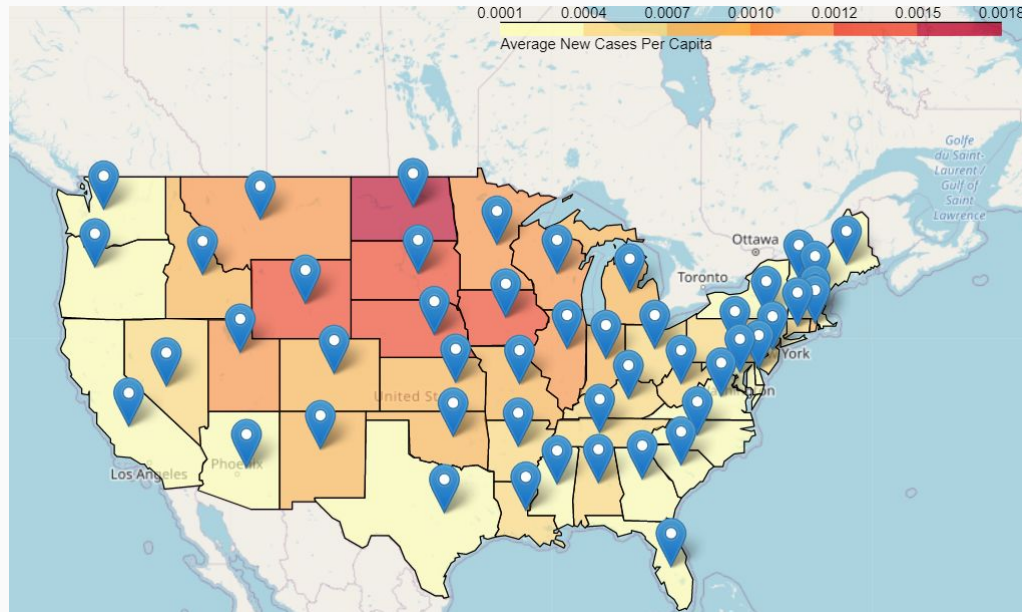
New Cases By State



New Cases By Day



Supporting Visualization



Conclusion

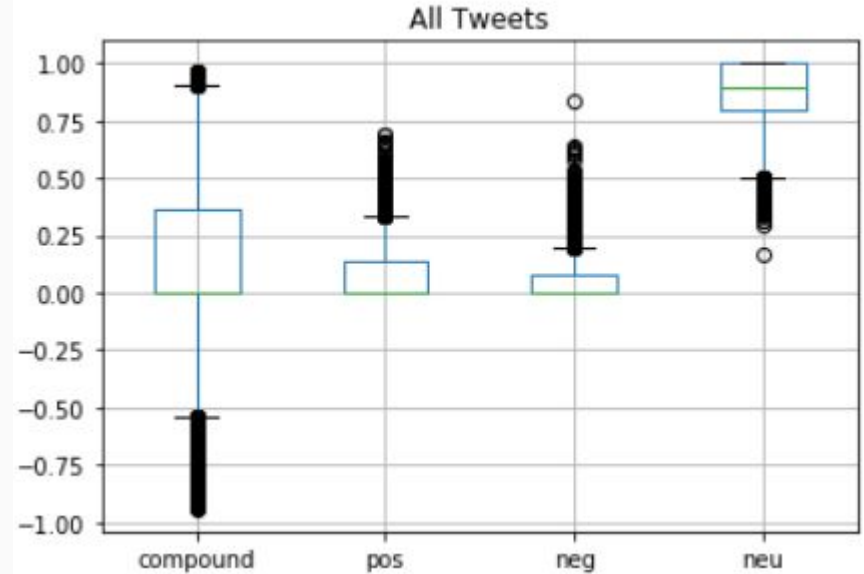
- Top 3 states for new cases: Illinois, California, Texas

	tot_cases	conf_cases	new_case	tot_death	conf_death	prob_death	new_death
state							
IL	4099809	4099809.0	84841	78778	75762.0	3016.0	802
CA	7200493	0.0	67658	127996	0.0	0.0	358
TX	7216532	0.0	67042	137620	0.0	0.0	966

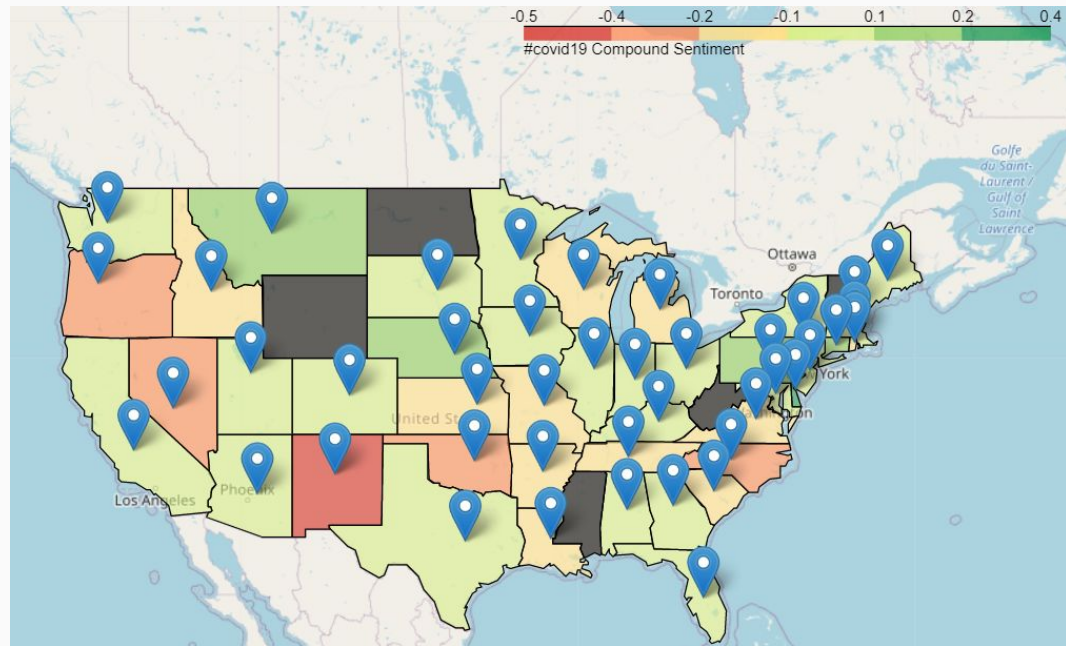
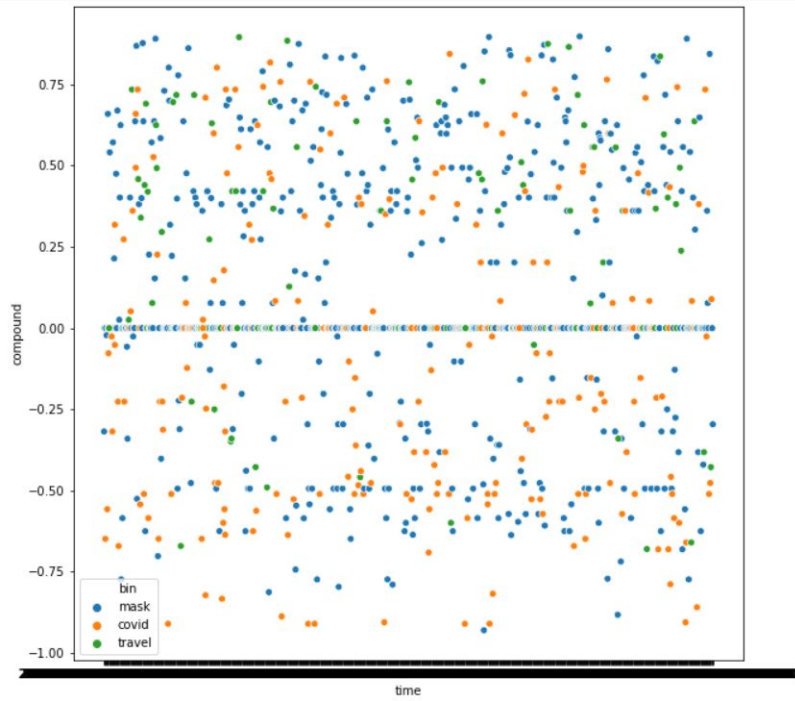
Focusing on new cases, we can look at distribution by state and by day in order to identify trends that may correlate with parallel trends with social media.

What is the sentiment of coronavirus-related tweets?

- Used VADER model for sentiment analysis - NLTK
- Returned positive, neutral, negative, and compound scores
- Applied on all hashtags and for all locations



Supporting Visualization



Conclusion

- While #travel had the highest compound score:

- #covid19 → Highest negative
- #maskup → Highest positive

- New Mexico and North Carolina

- High % change of new cases
- High negative score for #covid19

- Sentiment scores somewhat consistent over 7-day period
- Correlation warrants further investigation

	bin	neg	neu	pos	compound
2	travel	0.020455	0.893188	0.086358	0.170609
1	mask	0.049105	0.854347	0.096553	0.105901
0	covid	0.081858	0.845110	0.073033	-0.029543

Does a correlation exist between social media sentiment and Covid-19 in the United States?

- Carried out regression and location analysis

```
#define response variable
y = data['compound']

#define explanatory variable
x = data[['time']]

#add constant to predictor variables
x = sm.add_constant(x)

#fit linear regression model
model = sm.OLS(y, x).fit()

#view model summary
print(model.summary())
```

```
=====
                        OLS Regression Results
=====
Dep. Variable:          compound    R-squared:                0.003
Model:                  OLS        Adj. R-squared:             0.003
Method:                 Least Squares    F-statistic:            957.7
Date:                  Mon, 23 Nov 2020    Prob (F-statistic):     6.04e-210
Time:                  14:15:08    Log-Likelihood:         -1.6245e+05
No. Observations:      303000    AIC:                    3.249e+05
Df Residuals:          302998    BIC:                    3.249e+05
Df Model:               1
Covariance Type:       nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	-8408.2516	271.705	-30.946	0.000	-8940.786	-7875.718
time	0.0114	0.000	30.946	0.000	0.011	0.012

```
=====
Omnibus:                7583.644    Durbin-Watson:           0.006
Prob(Omnibus):           0.000    Jarque-Bera (JB):        3814.192
Skew:                    -0.005    Prob(JB):                 0.00
Kurtosis:                 2.450    Cond. No.                  2.67e+11
=====
```

Regression Examples

OLS Regression Results

```

=====
Dep. Variable:          new_case    R-squared:                0
Model:                  OLS         Adj. R-squared:           0
Method:                 Least Squares   F-statistic:              4
Date:                   Mon, 23 Nov 2020   Prob (F-statistic):       8.34
Time:                   14:26:13         Log-Likelihood:           -1.5837
No. Observations:      173619          AIC:                      3.167
Df Residuals:          173616          BIC:                      3.167
Df Model:               2
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.
const	1.773e+07	1.89e+06	9.360	0.000	1.4e+07	2.14
time	-24.0374	2.568	-9.359	0.000	-29.071	-19
compound	14.4045	12.976	1.110	0.267	-11.027	39

```

=====
Omnibus:                112920.449    Durbin-Watson:           0
Prob(Omnibus):           0.000        Jarque-Bera (JB):        1818704
Skew:                    2.919         Prob(JB):                0
Kurtosis:                17.742        Cond. No.:               2.63
=====

```

OLS Regression Results

```

=====
Dep. Variable:          compound    R-squared:                0
Model:                  OLS         Adj. R-squared:           0
Method:                 Least Squares   F-statistic:              4
Date:                   Mon, 23 Nov 2020   Prob (F-statistic):       4.40
Time:                   14:30:41         Log-Likelihood:           -1.624
No. Observations:      303000          AIC:                      3.24
Df Residuals:          302997          BIC:                      3.24
Df Model:               2
Covariance Type:       nonrobust
=====

```

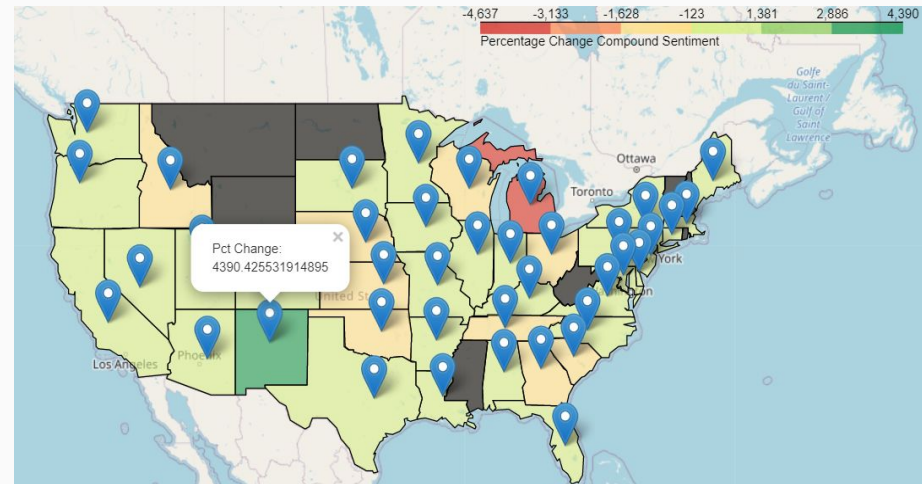
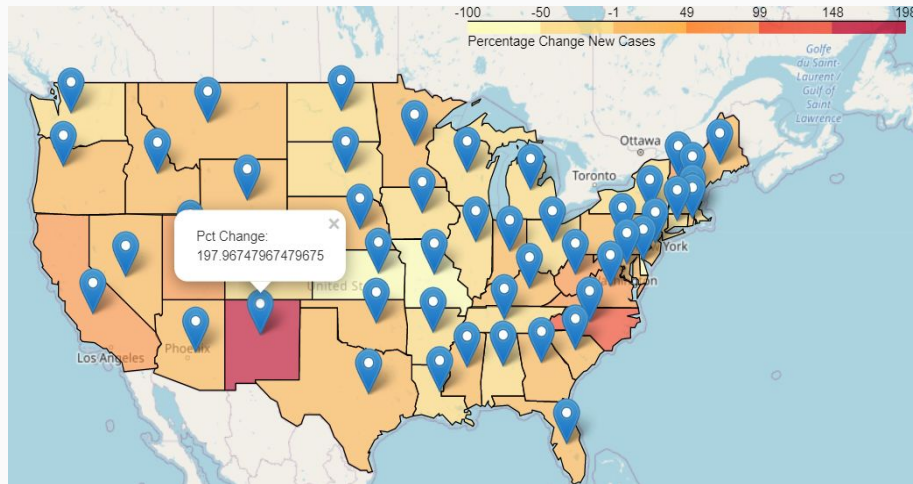
	coef	std err	t	P> t	[0.025	0.
const	-8402.0063	271.711	-30.923	0.000	-8934.552	-786
time	0.0114	0.000	30.923	0.000	0.011	0.011
new_case	-8.856e-07	3.15e-07	-2.815	0.005	-1.5e-06	-2.6

```

=====
Omnibus:                7574.269    Durbin-Watson:           0
Prob(Omnibus):           0.000        Jarque-Bera (JB):        381
Skew:                    -0.005         Prob(JB):                0
Kurtosis:                2.451        Cond. No.:               2.6
=====

```

Supporting Visualization

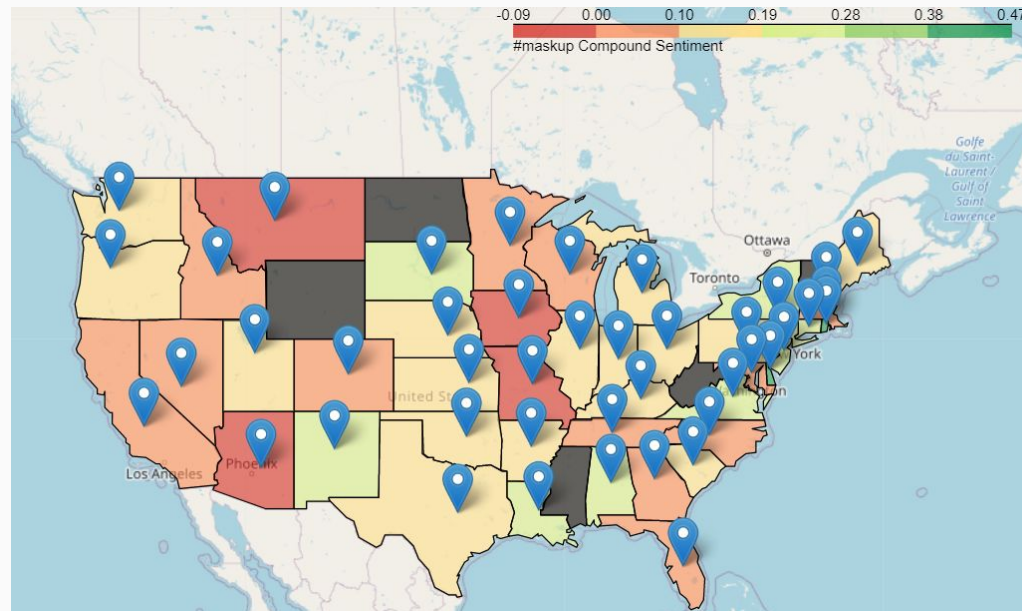


Conclusion

- Limited by time frame of publicly accessible tweets
- Data thoroughly tested for any correlation
- Process projected on a larger scale would yield more statistically significant results

Thank You!

Appendix I



Appendix II

