

A decorative graphic on the left side of the slide consisting of two overlapping parallelograms. The front one is blue and the back one is a light green. They are positioned diagonally, with the blue one partially covering the green one.

NYC Covid-19 Cases and Twitter

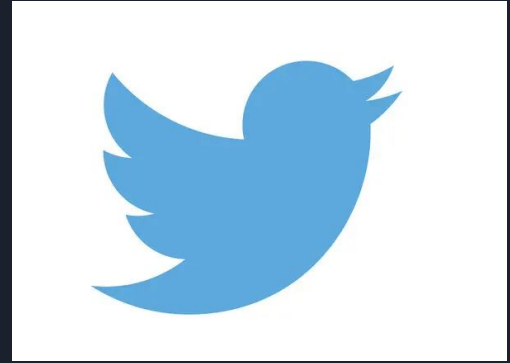
Eric Cusick



Business Case

We're looking to see if there is a correlation of NYC Covid cases rate to keywords found in tweets from verified user in New York City.

The goal is to use NLP and Random Forest Regressor modeling to see any words that have a higher correlation with new Covid cases.



Twint API

Twint API is used to gather all of the tweets relating to Covid in NYC.

Input parameters to obtain the desire tweets:

- Location : New York City
- Date : 01/01/2020 - 02/14/2021
- Verified Users
- Likes Count greater than 50

All desire tweets were compile into a csv file named 'covidtweets.csv'

The following search words were used to obtain the tweets:

- Covid
- Corona
- Coronavirus
- Mask
- Vaccine
- Quarantine





NLP

- Approximate 3000 tweets were obtained from using Twint API
- Removed hashtag, url links, @username, punctuations, and digits
- Tokenization
- Stopwords
- Lemmatization
- Count Vectorizer

Before Cleaning

```
'Here's a fun video by Newark's own @DJLILMAN973 Ft. our Mayor @rasjbaraka reminding all of Newark to Mask Up. #MaskUpNewark h  
ttps://t.co/2BuHyG7KCD'
```

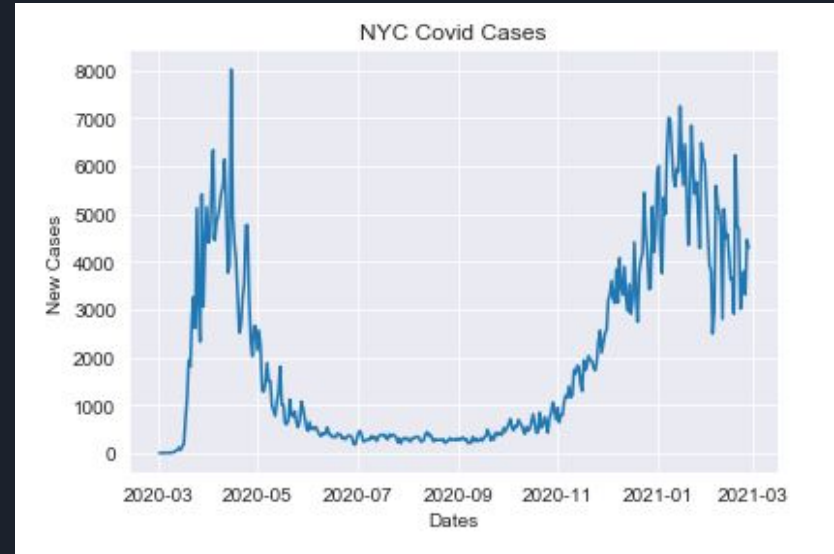
After Cleaning

```
'Heres fun video Newarks Ft Mayor remind Newark Mask Up MaskUpNewark'
```

NYC Cases EDA

NYC Covid Dataset were obtained and modify from the NYTime's [Github](#) page where there is ongoing reports of Covid.

The
New York
Times





Modeling Process Random Forest Regressor

Created two dataset for modeling, original and likes count.

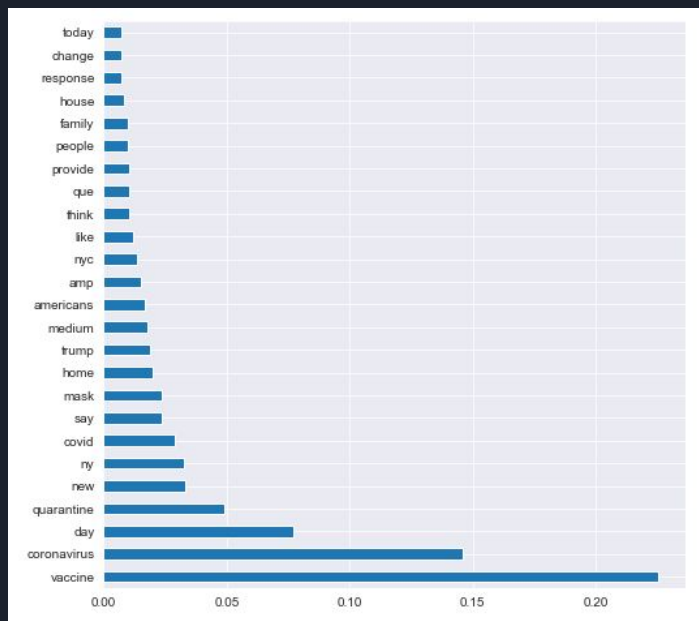
- Original dataset is there as a control model taking in a combine tweets for each individual days to correlate with NYC Covid cases
- The likes count dataset is the same as original but with one exception. The number of likes that a tweet has is the value of those words.



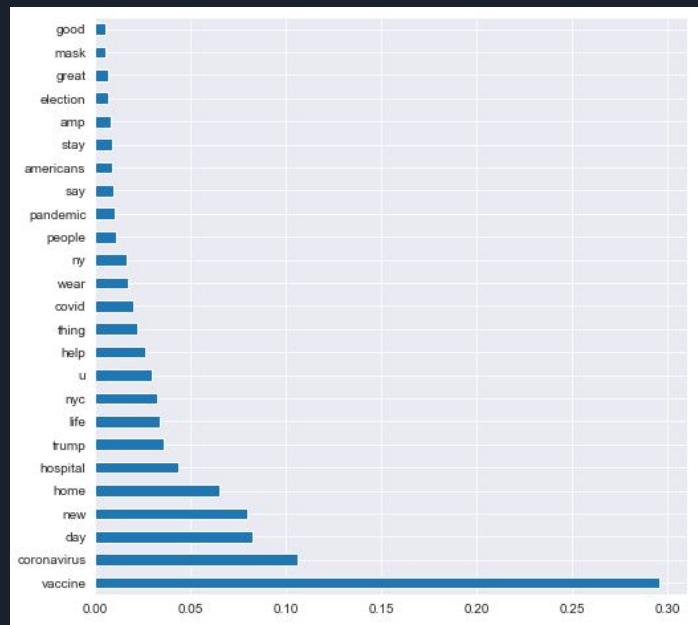
Final Model

Feature Importances of both Random Forest Regressor Models

Original Dataset



Likes Count Dataset



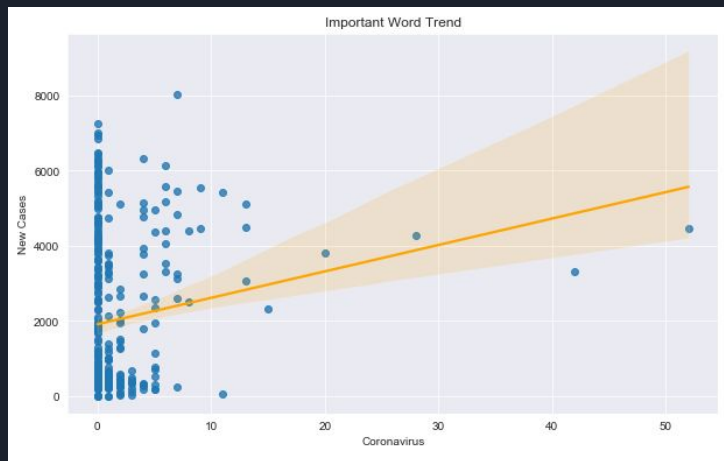
Final Model Cont.

Original:

- Training R2: 0.61
- Test R2: 0.28
- Train RMSE: 1308.28
- Test RMSE: 1593.89

Likes:

- Training R2: 0.52
- Testing R2: 0.27
- Train RMSE: 1437.65
- Test RMSE: 1606.66



Recommendation

After looking at the features importance graphs, there is a lot of health based words and some political words as well that correlate with new covid cases. The verified user should be careful when it comes to writing and phrasing a tweets. As they could influence the behaviors of their followers.



Future Work

Some further work may include:

- Obtaining more tweets
- Further cleaning of NLP
 - Including more stopwords
 - Getting rid of Keywords used for gathering tweets
 - Placing language filter
- Create a bi-gram or n-gram datasets
- Create smaller windows





Thank You

Thank for your time and attention to the presentation today.

Github link: https://github.com/Ericusick/Twint_Covid_NYC