Hamin Han

# Machine Learning and COVID-19 Data

Over the last few years, everyone around the globe has been affected by COVID-19. I want to use machine learning to find out if there is any visible pattern related to population and geography. Specifically, I want to use clustering on COVID-19 data and see if it visualize clear clusters formed by the geography and population of the regions. I would assume that the regions with massive populations, especially near populous cities, would have far more COVID cases. But what would that look like on a local scale?

**Related Work:**

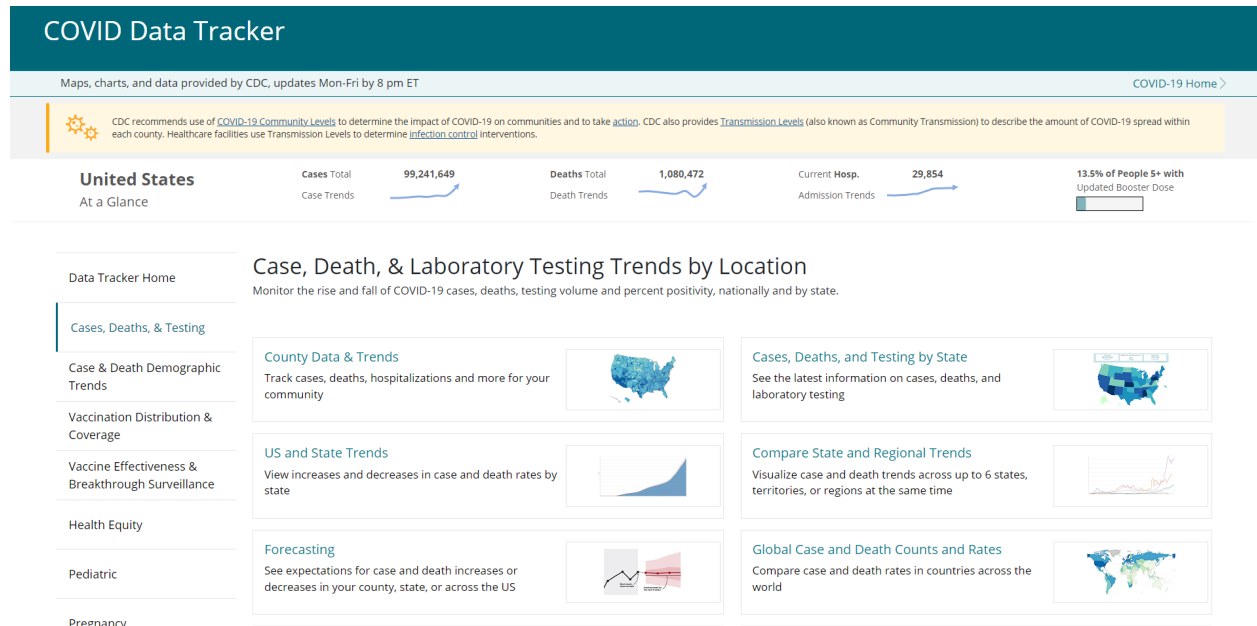https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7998460/

From a published academic journal called *The Geography of the Covid-19 Pandemic: A Data-Driven Approach to Exploring Geographical Driving Forces,* the authors concluded that "the disease initially spread in the densely and heavy populated capital region and over time moved to more distant regions of the country, furthermore single events of large spread and counter measures of these showed a space-time ripple effect in decreasing infection rates". From their investigation of the relationship of geographic factors and COVID-19, they were able to first prove the direct relationship between population size and the spread of the virus. This was also the case for my findings when running clustering on US states and local counties.

**Methodology:**

First and foremost, it is important to have accurate and organized data to use. For my project, there was an abundance of COVID-19 data out there but I used the data provided by CDC, the Centers for Disease Control and Prevention.

LINK: https://covid.cdc.gov/covid-data-tracker/#cases-deaths-testing-trends



From there, I was able to get specific data for recent covid cases in state counties.

(Maryland Counties Covid Cases Data)

And then from there, I do the clustering. First import all the libraries and data

```
# import libraries

import pandas as pd
import numpy as np
import random as rd
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn import metrics
from sklearn.metrics import pairwise_distances
from sklearn.cluster import KMeans
```

```
[ ]  # import the data csv file

data = pd.read_csv('county_level_latest_data_for_maryland.csv', skiprows=2, on_bad_lines='skip')
data.columns = [x.replace("\n", " ") for x in data.columns.to_list()]
data.head()
```

| | FIPS code | State Name | County | Weekly Cases | Weekly Case rate /100k | % Change in weekly Cases from past week | Test positivity rate - last 7 days | Test positivity rate - absolute change | Total diagnostic tests - last 7 days |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 24001 | Maryland | Allegany County | 74 | 105.09 | 21.31 | 5.81 | -1.56 | 1314 |
| 1 | 24003 | Maryland | Anne Arundel County | 541 | 93.40 | 68.01 | 12.64 | 1.56 | 4419 |

Run K-means looking at the Population Density and Weekly Cases columns. But first decide the optimal k-value using the elbow method. (4 in this case)

For reference, "The elbow method is a graphical representation of finding the optimal 'K' in a K-means clustering. It works by finding WCSS (Within-Cluster Sum of Square) i.e. the sum of the square distance between points in a cluster and the cluster centroid." (towardsdatascience)

```
cost =[]
for i in range(1, 11):
    KM = KMeans(n_clusters = i, max_iter = 500)
    KM.fit(X)

    # calculates squared error
    # for the clustered points
    cost.append(KM.inertia_)

# plot the cost against K values
plt.plot(range(1, 11), cost, color ='g', linewidth ='3')
plt.title('The Elbow Method')
plt.xlabel("Value of K")
plt.ylabel("Squared Error (Cost)")
plt.show() # clear the plot
```

Now implement K-means clustering using the Sklearn library.

```python
kmeans = KMeans(4) # From the value taken from the elbow method
kmeans.fit(X)
identified_clusters = kmeans.fit_predict(X)

data_with_clusters = data.copy()
data_with_clusters['Clusters'] = identified_clusters
plt.scatter(data_with_clusters['Population density - county'],data_with_clusters['Weekly Cases'],c=data_with_clusters['Clusters'],cmap='rainbow')
plt.title('COVID-19 Cases in Maryland Counties')
plt.xlabel('Population Density (Persons per Square Mile)')
plt.ylabel('Weekly Cases')
plt.show()
```



Also test the silhouette score and average it to see how good the clusters are.

For reference, "The silhouette value is a measure of how similar an object is to its own cluster (cohesion) compared to other clusters (separation). The silhouette ranges from −1 to +1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters." (Wikipedia)

```python
score = metrics.silhouette_score(X, kmeans.labels_, metric='euclidean')
print('Average Silhouette Score over All Samples: %.3f' % score)

Average Silhouette Score over All Samples: 0.770
```

I did this for the local counties in Maryland and the nearby state Virginia. I also thought to run the same clustering algorithm for the four very populous states, California, Texas, Florida, and New York. As well as doing a similar approach for clustering the 50 states of America.

All of the ipynb files and csv files are included. And kept separate for each clustering.

**Analysis of Results:**

I made a separate powerpoint slides that organized the results and analysis.

But most of the results verified the direct relationship between population density and the number of COVID-19 cases. Additionally, the clusters formed were divided into geographical regions, initially grouping around large populous cities and a second group forming that are near these cities, and the other set of cluseters that are far away. Even if multiple clusters with multiple cities were formed, they would group up near the matching city or regions. The results proved what was being said in the academic journal: "the disease initially spread in the densely and heavy populated capital region and over time moved to more distant regions of the country."

Check the slides for the indepth analysis.

I will also add parts of the slides as reference at the end of the report.

**Limitations:**

This might be due to my lack of experience with the coding, but it's hard to identify which county or state the point on the graph represents. I had to check the values through the csv file and match it up to find out what the data point was and where it's located on the actual map. Also, for the USA states clustering, the data didn't have the population density of the states so I just had to use the number of total cases as a general benchmark for the amount of people there are in the states.

**Potential Follow-up Work:**

One thing is to add more functionality to the code. Be able to see what each data points are and how that might look geographically as well. Maybe fine tune the elbow method and the silhouette score code as well. Just more depth and analysis with the coding since I would consider myself a novice when it comes to Machine Learning and python in general. I'm sure there are a lot of more useful libraries that I could've implemented to make the quality of the project better.

As for a long-term improvement, perhaps going beyond and trying to predict how the new variants will spread geographically using Machine Learning. Just predicting the effects and spread of the virus is impressive but tying it into geographical regions, especially locally in Maryland would be a cool idea to explore in the future.

Hamin Han

EXTRA REFERENCES FOR RESULTS (Check the slides included)

## Maryland Results
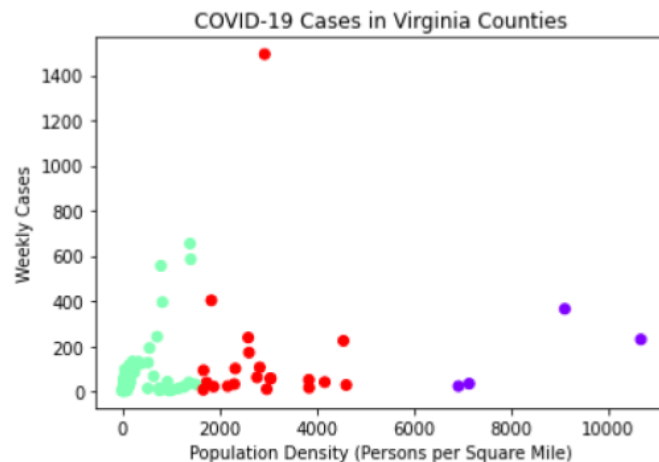
COVID-19 Cases in Maryland Counties



- **Yellow:** Less populated counties that have less cases as expected.
- All the other counties are nearby DC or Baltimore.
- **Purple:** Counties near Baltimore City. The closer the county is to the city, the more cases is has.
- **Red:** Montgomery County and Prince George's County that are right next to DC.
- **Cyan:** Baltimore City. Still a lot of cases but expected more with the population.

Clusters that are easily identifiable and shows the direct relationship between population and COVID cases.

Average Silhouette Score over All Samples: 0.770

## Virginia Results

COVID-19 Cases in Virginia Counties



- **Purple:** Mostly counties near DC.
- **Green:** A lot of the other counties that aren't close to nearby cities.
- **Red:** A mixture of counties near cities such as Virginia Beach, Chesapeake City, and Richmond.

There are clusters being formed and divided into geographical regions but doesn't show a great relationship between population and COVID cases.

This could be due to the data being collected in the recent 7 day period so there might be a lot of outliers.

Average Silhouette Score over All Samples: 0.779

# California Results

### COVID-19 Cases in California Counties



Average Silhouette Score over All Samples: 0.763

- **Yellow:** San Francisco County with the most population density.
- **Cyan:** Los Angeles County with the incomparably high number of COVID cases.
- **Red:** These are mostly counties in the southern parts of California, parts of them are near LA and a handful are near San Francisco.
- **Purple:** The rest of California. Mostly on the northern parts of California or the outer borders of California where it's far from the big cities.

Clusters that are easily identifiable by geographical location, particularly around LA and San Francisco.

Surprisingly, San Francisco County has a small number cases compared to its population density while LA county is on another level with COVID cases

# Texas Results

### COVID-19 Cases in Texas Counties



Average Silhouette Score over All Samples: 0.893

- **Green:** These counties aren't exactly near each other but they all contain big cities such as Dallas, Houston, Austin, and El Paso.
- **Red:** This is the second group of counties that are still nearby one of the big cities.
- **Purple:** All the rest of counties that aren't exactly near any populous cities.

Clusters that are easily identifiable and shows a clear and direct relationship between population and COVID cases.
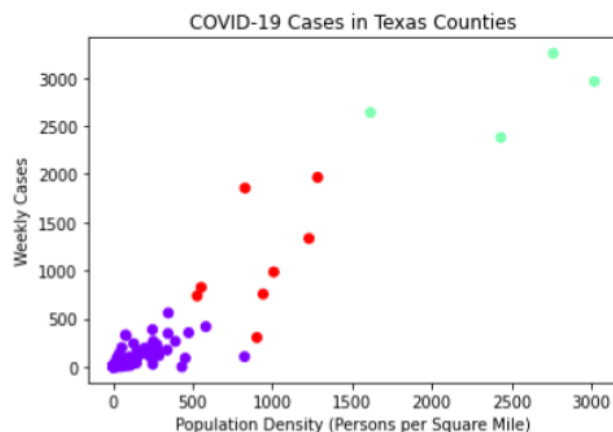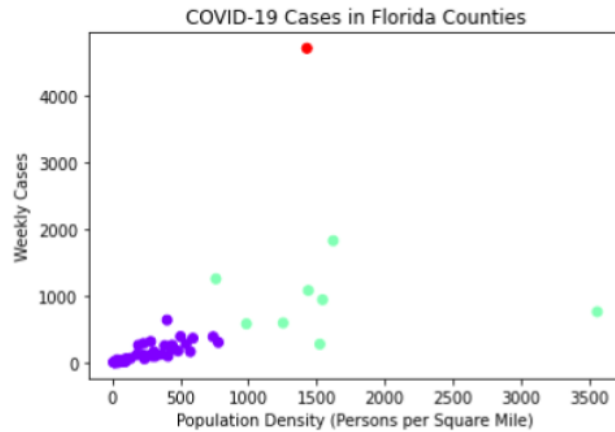
# Florida Results



COVID-19 Cases in Florida Counties
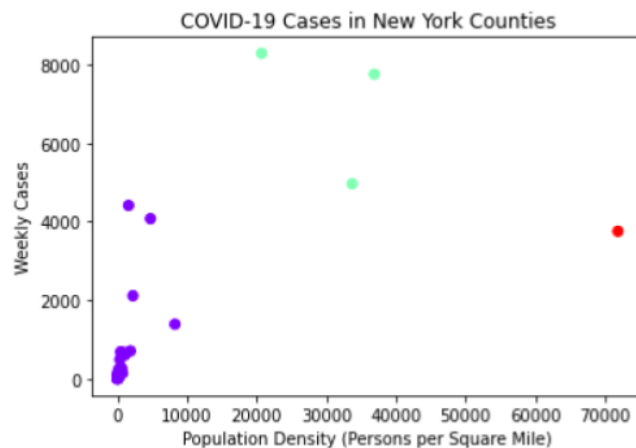
Average Silhouette Score over All Samples: 0.754

- **Red:** Miami-Dade County with the most cases. Most likely due to having a big city like Miami
- **Green:** Counties with or near big cities such as Orlando, Hollywood, and Tampa. These are also the counties with more population density.
- **Purple:** Rest of the counties that are far from big cities and less population densities.

Clusters that are easily identifiable and shows the direct relationship between population and COVID cases. Although there is an outlier with Miami-Dade with the enormous amount of cases while not having the most population density. Most likely due to a big city like Miami being there and spreading COVID easily.

# New York Results



COVID-19 Cases in New York Counties

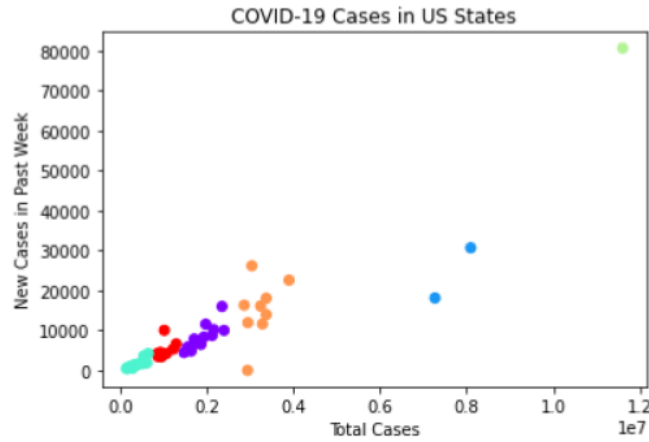Average Silhouette Score over All Samples: 0.932

- **Red and Green** are the counties of NYC
- **Red:** New York County with Manhattan with the most population
- **Green:** Bronx, Kings, Queens counties with Bronx, Brooklyn, and Queens.
- **Purple:** All the other counties.

Clusters that are easily identifiable and shows the direct relationship between population and COVID cases. Large population density with a lot of cases where NYC is. Although New York County with the most population seem to be an outlier, could be because the data is only from the recent 7 days.
For the purple cluster, more cases as the counties get close to NYC towards south of the state

# US States Results



COVID-19 Cases in US States

Average Silhouette Score over All Samples: 0.661

- **Green:** California
- **Blue:** Texas and Florida
- **Orange:** New York, New Jersey, Pennsylvania and other States. Mostly in the Mideast and Midwest.
- **Cyan:** Mountain-Prairie area
- **Red and Purple:** Everything in between. WIth more cases for states with or near large cities

Clusters that are easily identifiable and shows the direct relationship between population and COVID cases.

Although California is skyrocketing with the number of cases.