Matthew Shinder

# **<u>Capstone Two - Project Report</u>**

## *<u>Introduction</u>*

COVID-19 is a contagious disease that was caused by the SARS-CoV-2 coronavirus. While first known to originate in Wuhan in December of 2019, the virus quickly spread around the world creating a global pandemic. The goal of this project is to look at trends of COVID-19 data, specifically between each US county.

## *<u>Datasets</u>*

We will use the US counties COVID 19 dataset (from Kaggle) which contains COVID-19 case numbers dating from as early as February of 2020 and up to May of 2022. After data wrangling and modeling our COVID data, we can look into the US Census Bureau dataset to get demographic information by county. Connecticut data was lost from the main demographic file so an outside file was used to gather the data specifically for counties within that state.

<u>Links:</u>

https://www.kaggle.com/datasets/fireballbyedimyrnmom/us-counties-covid-19-dataset?resource=download

*https://www.ers.usda.gov/data-products/county-level-data-sets/county-level-data-sets-download-data/*

*https://data.census.gov/table?q=P6&t=Race+and+Ethnicity&g=010XX00US$0500000&y=2020&tid=DECENNIALDHC2020.P6*
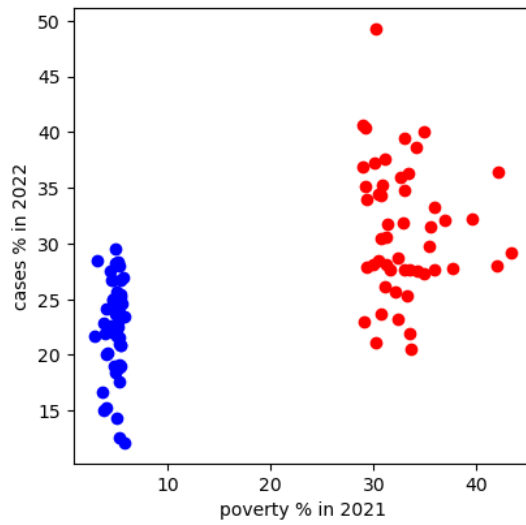
*https://portal.ct.gov/DPH/Health-Information-Systems--Reporting/Population/Annual-Town-and-County-Population-for-Connecticut*
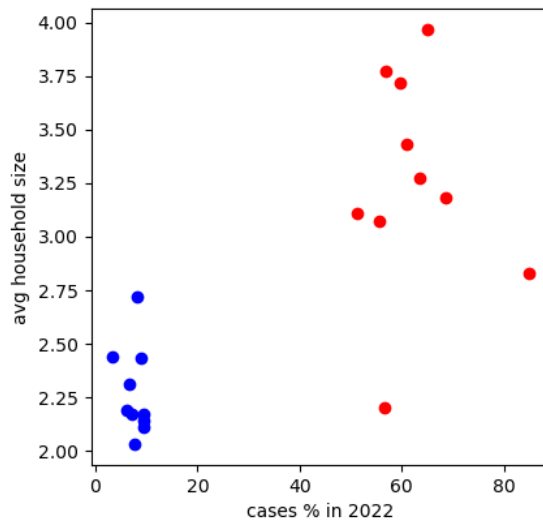
# *Cleaning/Wrangling*

To clean the data, numerous steps were made to remove many values as we had over 2 million data rows of data in our COVID-19 cases dataset. Rows that had invalid data such as null values were removed as well all data associated with US territory locations such as Puerto Rico and the Virgin Islands. When it came to the demographics datasets, the main issue was the naming of columns and which columns from our data source that we would want. The hardest problem in all of these steps was the combination of all data sources. First, we had to make sure all the counties' demographic data sources matched by name and state in all data sources before merging. When doing this step, Connecticut as stated above was missing data, so an outside file was used in order to fill in the missing values as missing the whole state of Connecticut could be detrimental in the findings. The data from the newly formed demographics data frame was then checked with the COVID 19 data frame in order to only include counties that were listed in both. Those that were not, were removed from their respective dataset (very small amount). New York City was combined as the COVID 19 dataset listed the county as a whole vs. the demographic datasets that had them listed by NYC boroughs such as Queens and Brooklyn.
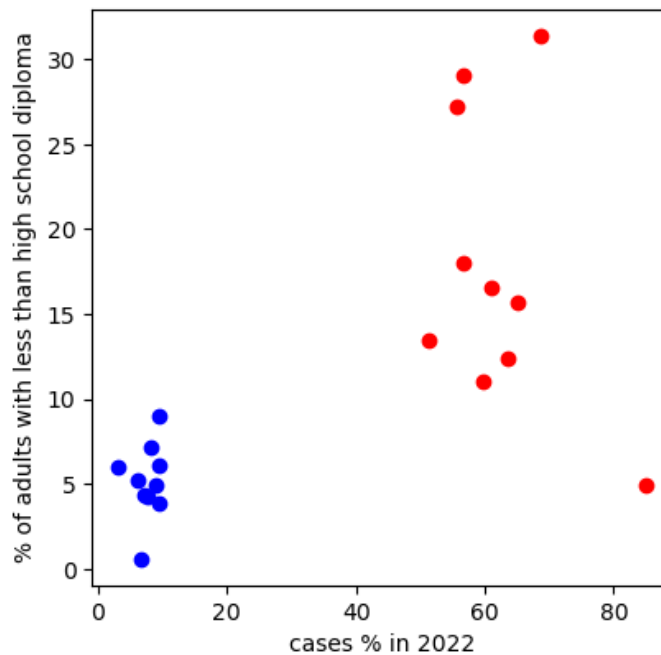
# *EDA*

Two approaches were taken in order to find trends in our data. The first approach was to look at counties that have high covid rates and then compare their demographic stats whereas the second approach would do the same but vice versa. Finding out the most populous COVID 19 case counties proved to just be the most populous counties in the US. To combat this, a new column was created in order to show the COVID 19 case versus the total population of the county, referred to as covid case percentage. Three dates (one per year) were chosen to separate the data based on key events; however, the covid case percentage was an accumulated total by said date. The 2020 date was 11/13/20 as every US County had reported at least one case of COVID 19. For 2021, 6/1/21 was a good pace for COVID 19 as the reported cases had begun to dip down heavily due to the creation and use of the vaccine. Our 2022 was 5/13/22 as that was the end of the dataset given to us. The graphs below demonstrate either approach.

This first graph shows when the counties are sorted by highest and lowest poverty percentage counties, what their covid case percentage was in 2022.
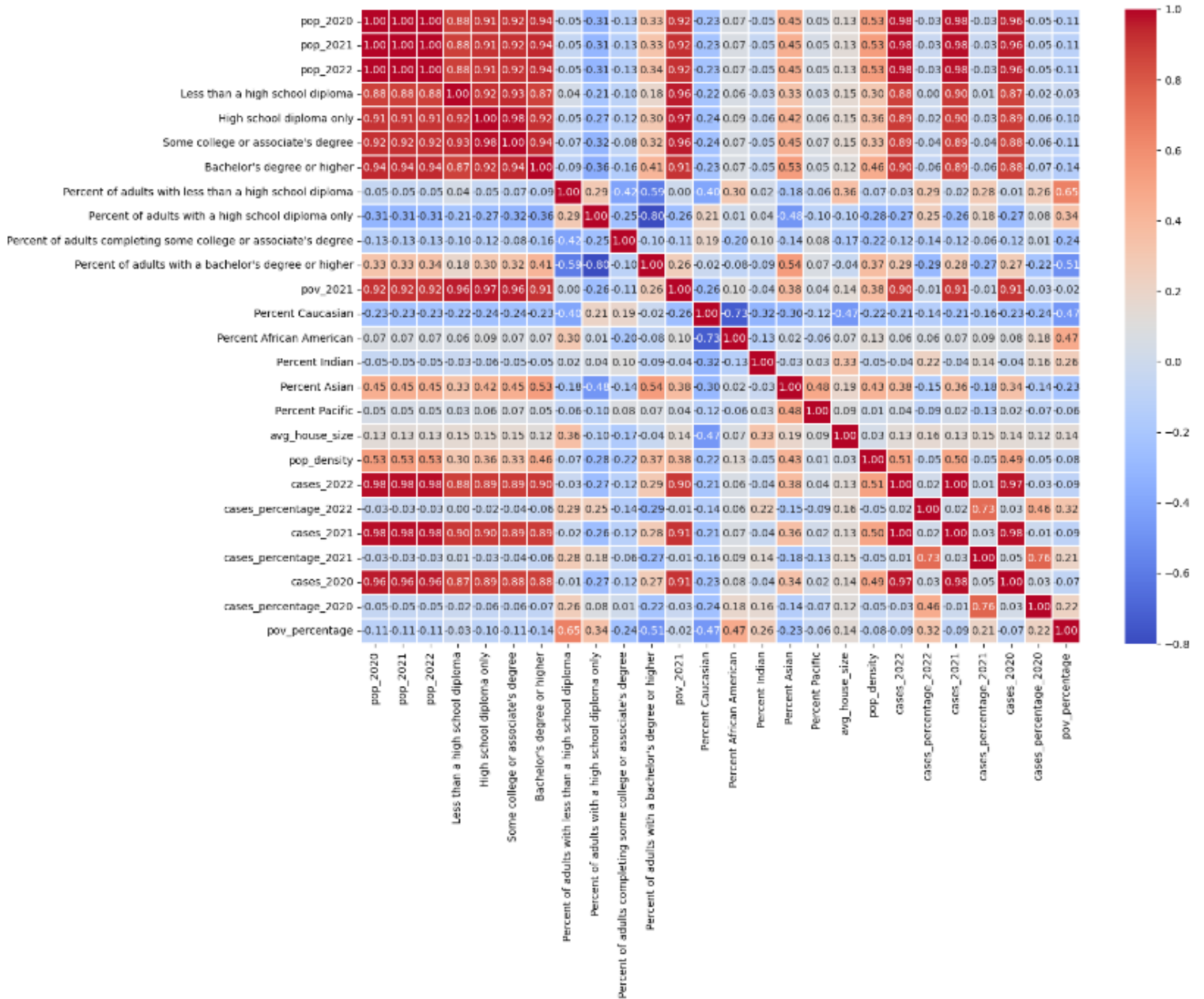


This second graph shows the same trend for the average household size in each county.
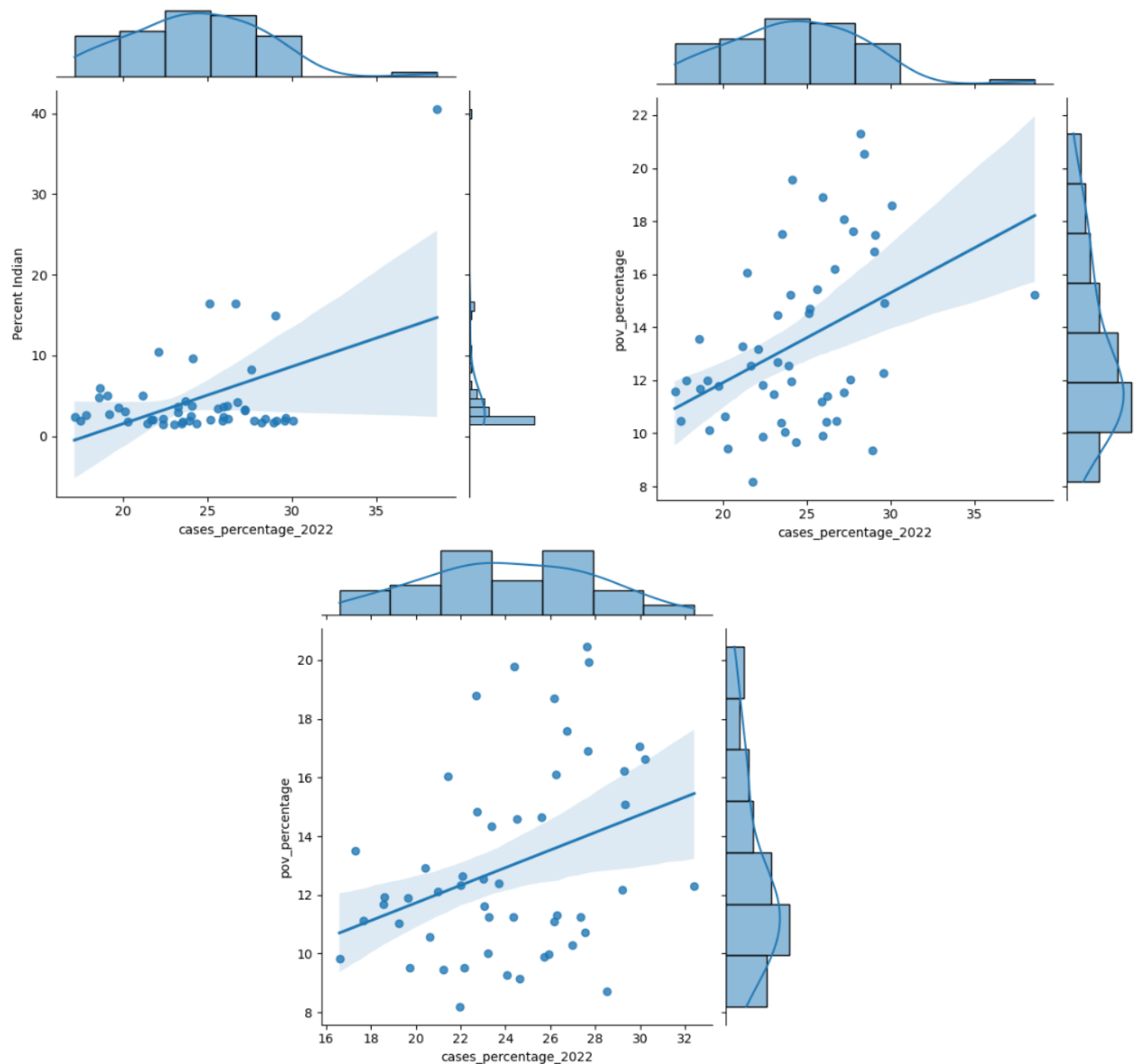
This last graph was sorted by COVID case percentage and then compared with an education stat.

To get a better look at these trends, heatmaps were generated but, to no chagrin, the data did not lead to great results or positive / negative correlations rather regressed to showing no trend. The heatmap on the next page shows most values near 0 where the highly correlated values (dark red or dark blue) are correlated just based on population numbers compared to the percentage columns which is what we are most interested in.

COVID data features correlation

To fix this problem, the counties were grouped by state and then the same heatmaps were generated in hopes to see a stronger correlation as states had different rules for how they handled COVID 19. This isn't the perfect solution as external factors such as super spreader events and state border counties can influence this data but, we can generalize the data in order to see some more noticeable trends.
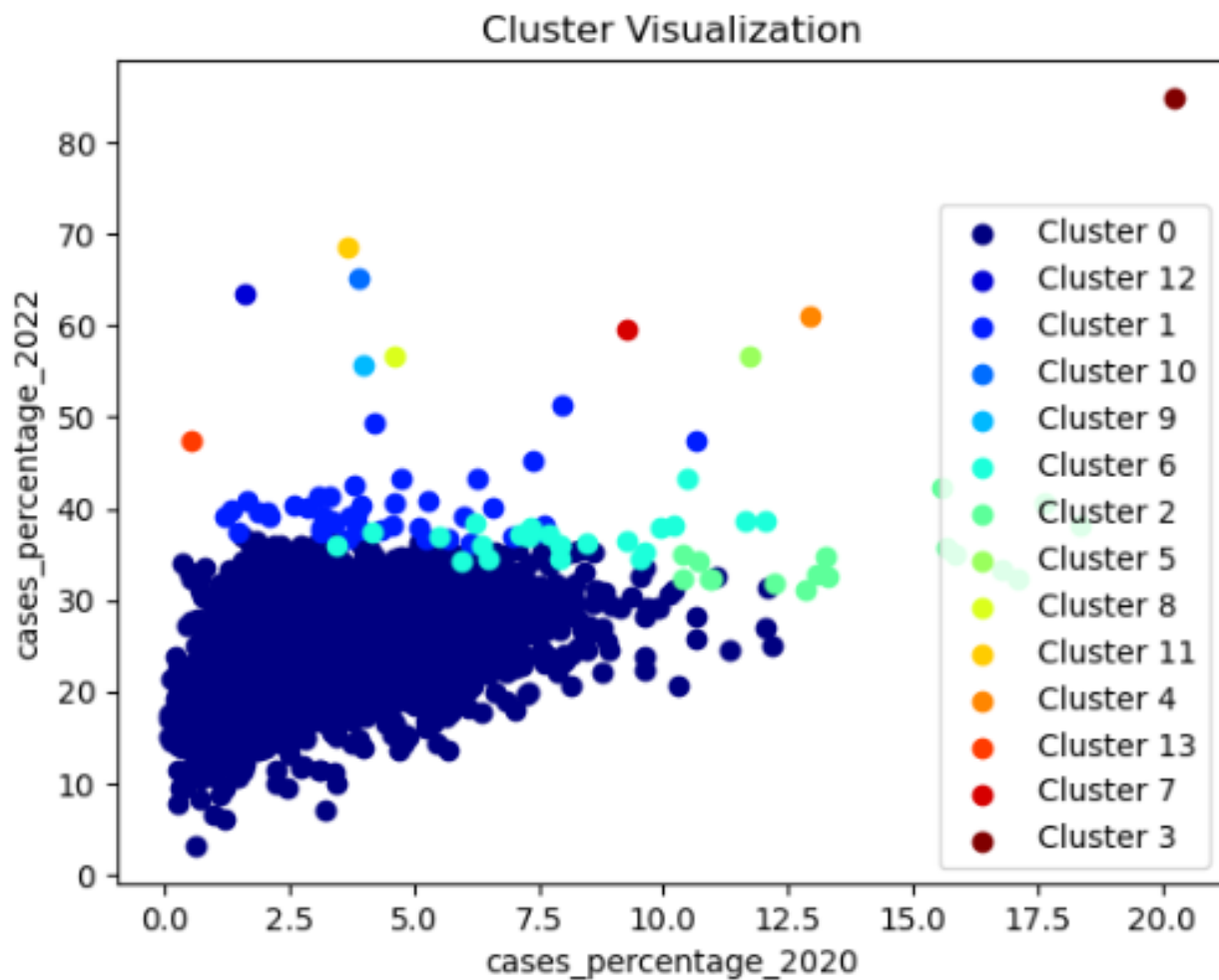


The first two graphs were sorted based on the mean average from each state (represented by a dot) while the bottom graph was sorted based on the median.
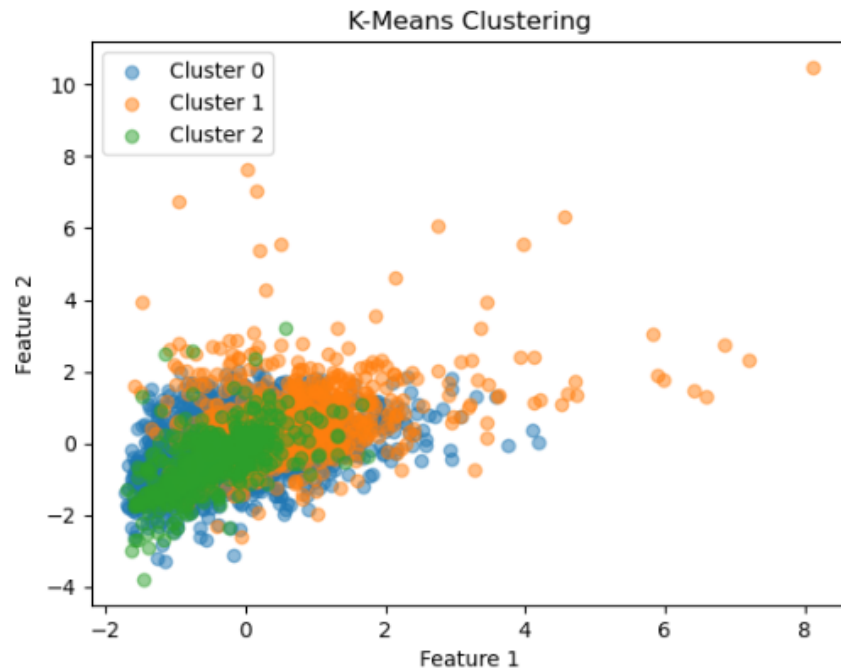
# *Preprocessing and Modeling*

The goals of this step were to set up and model our data for classification models in order to find better trends within our data than what we saw with our EDA step.
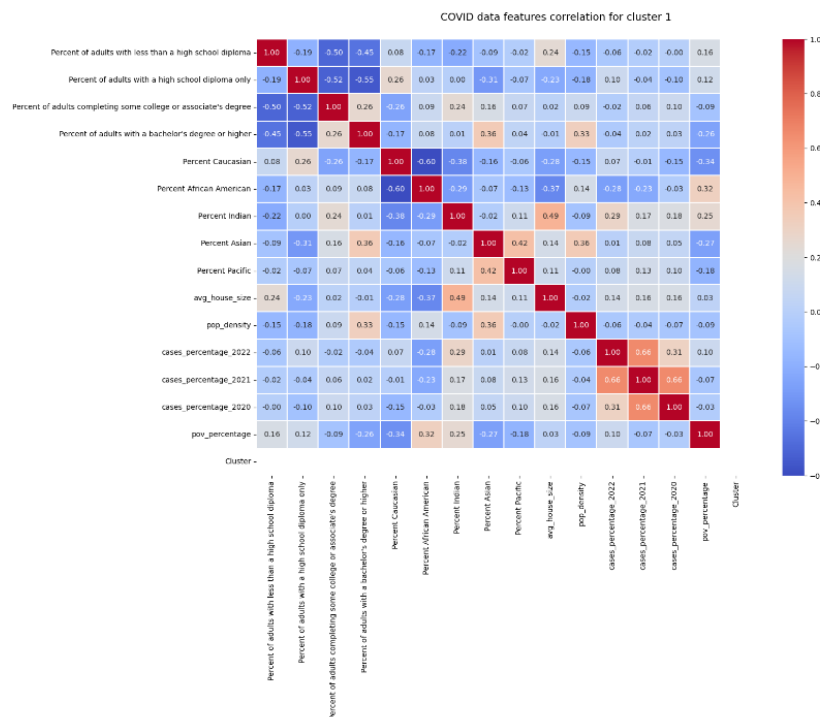
The first models built were with <u>Mean Shift Clustering</u>. Unfortunately, the models built from both approaches proved to be unsuccessful. As shown from the graph below, the models built clusters where one cluster clearly dominated the others when it came to the total number of data points (counties). Cluster 3 only had 1 data point!
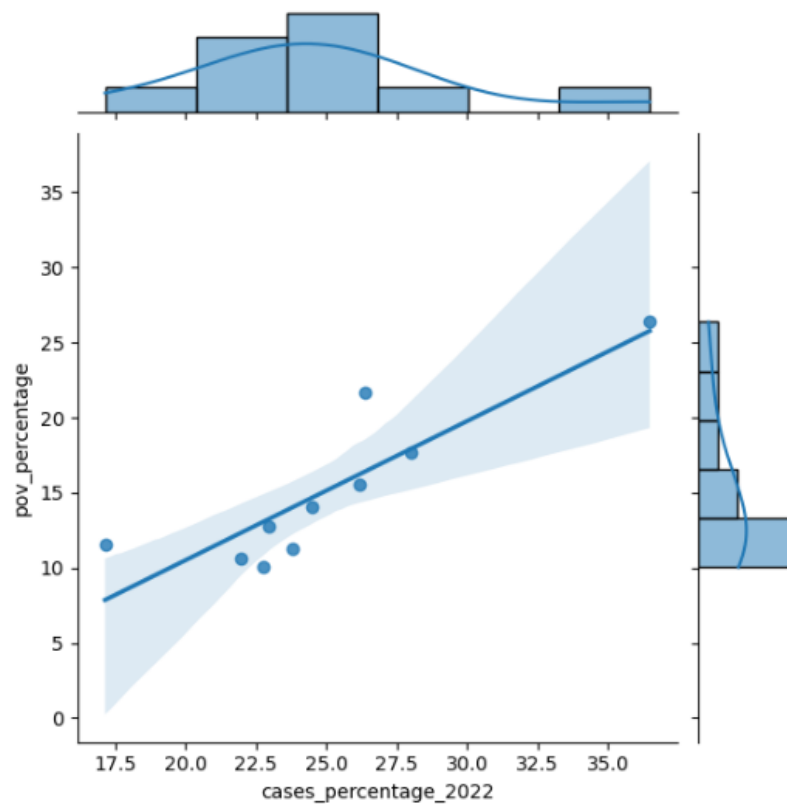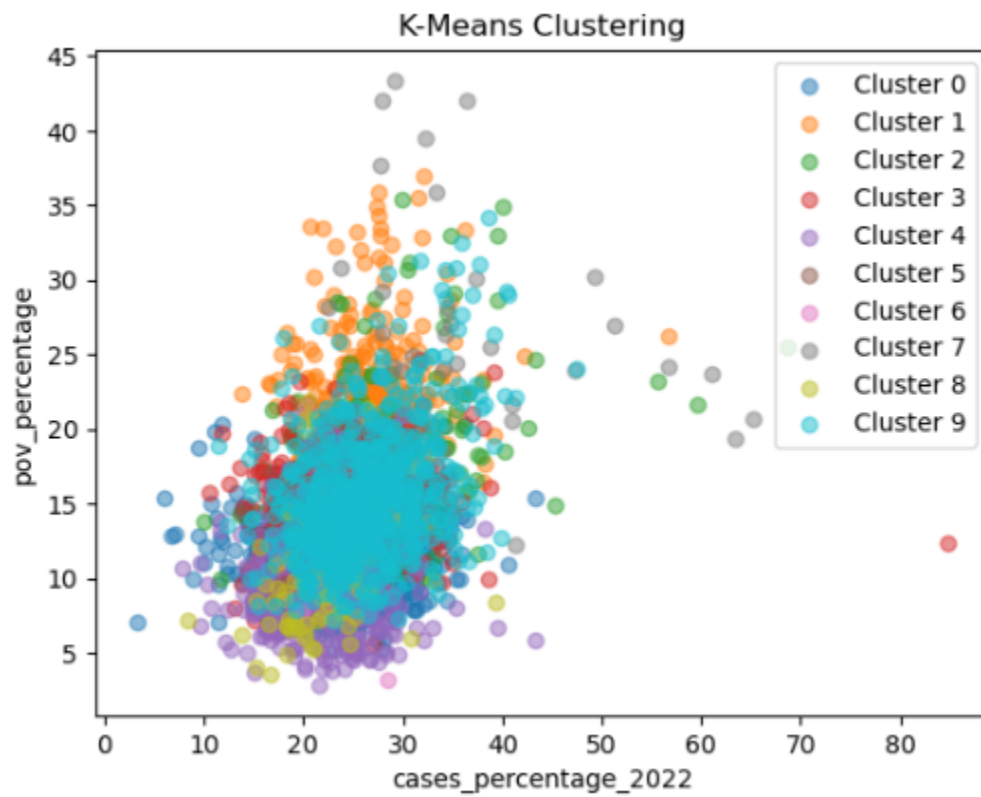
With the combination of <u>PCA</u> and <u>KMeans</u>, we were able to find a better grouping of clusters by reducing the number of clusters. These initial models were clustered by COVID case percentage.



However, when looking at a heatmap of a specific cluster, most if not all of the relations showed no correlation.

When trying from the other approach, it was noticeable that if you were to group the clusters by the mean, a stronger correlation would appear from the models.

# _Analysis_

When looking at the clusters generated by the KMeans modeling algorithm, we can see some noticeable trends between our country clusters. Due to the nature of COVID 19 being a national problem however, these findings may not explain the whole picture due to a plethora of factors that may influence why COVID 19 case totals were high / low in different counties. These factors that can't be listed in statistics definitely played a major role in the outcome of the datasets. An example is whether or not an individual would even get tested in the first place because it was not a mandatory nor daily practice due to costs, so there is no way of knowing how many cases truly happened in the US. However, the biggest trend from our data that we can take away was that the COVID 19 case percentage was high in counties that have high poverty percentages. This might be explained by those who had access to clean and healthy resources during quarantine. While the models may not be perfect, we can look at the results to see what locations might be seen as dangerous when it comes to a new contagious and infectious disease in the future.