# Risk Determination of Exposure to COVID-19

**2020 City of LA and RMDS COVID-19 Computational Challenge**

**Team: Risk Detectives**

**Binlun Feng, Leon Lu, Lei Wang, You Wang, Jennifer Ware, Xiangbo Wen**

**Mentored by:**

**Dr. Yuanjie (Ed) He, Dr. Honggang Wang**

# Introduction

Through literacy research and logical analysis, the risk score for a specific area is determined as the potential increased COVID-19 cases per population density, by the formula:

$$risk\ score\ =\ \frac{Number\ of\ New\ Daily\ Cases*\ 10^6}{Population\ /\ Area\ size}.$$

# Data

## 1. Data sources

Data is collected at the city or community level within Los Angeles County. This includes neighborhoods and communities within the city of Los Angeles and nearby cities inside Los Angeles County. Our data comes from four different sources:

- **Latimes-place-totals.csv**

  LA Times keeps track of COVID-19 in California from multiple health and city agencies. Updated cases are updated on a daily basis. The github repository collects data from LA Times beginning on March 16, 2020 to Present. This is the only source we found that has the information of COVID-19 at city and neighborhood level. From this dataset, we selected these features for our study:

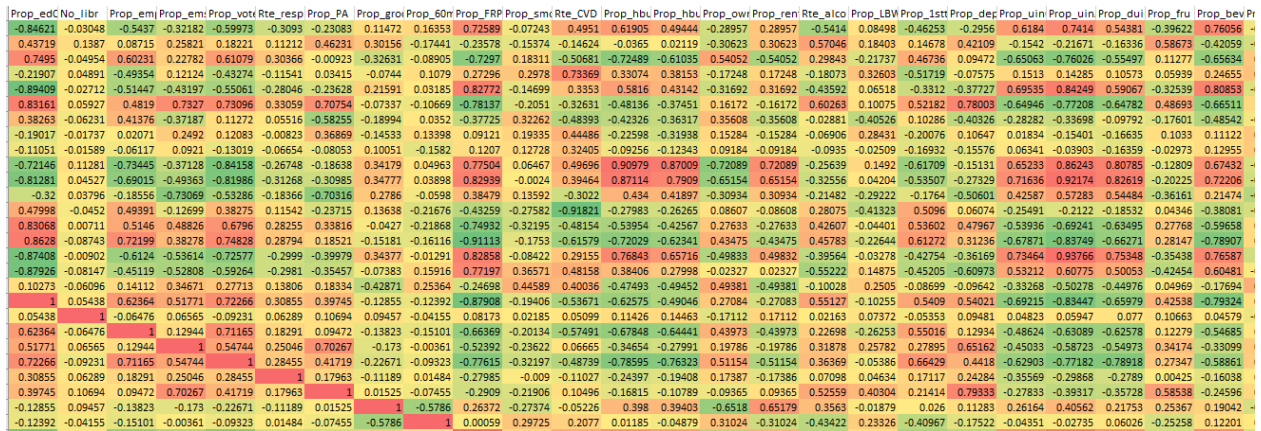| field | type | description |
|---|---|---|
| date | date | The date when the data were retrieved in ISO 8601 format. |
| county | string | The name of the county where the city is located. |
| place | string | The name of the city, neighborhood or other area. |
| confirmed_cases | integer | The cumulative number of confirmed coronavirus case at that time. |
| x | float | The longitude of the place. |
| y | float | The latitude of the place. |

*GitHub link*

  **Data handling:**
  1. Filter out the records which are not in the Los Angeles County
  2. Calculate 'new_cases' by subtracting 'confirmed_cases' on two consecutive days in a certain place
  3. We believe that the risk level in one city is not isolated, meaning that it will be affected by the risk levels in the nearby cities. Therefore, we believe that the nearby new cases are relevant to this study. In order to get the 'nearby_cases,' we summed up 'x' and 'y' which gives us a general idea of the location. We then sorted 'x+y' and summed up five consecutive 'new cases' and assigned this sum value to 'nearby_cases'.

4. We determined to use the next day's new case for our risk score. We created a new column, 'tmr_new_cases' to store this value.

● **Los Angeles County City and Community Health Profiles 2018**

The city and health profile is the next dataset we found accurate to the city and neighborhood level. This dataset is provided by the Los Angeles County Public Health Department and contains demographic, social, and health information. In the beginning, the dataset had over 80 columns. The heatmap listed below shows an overview of the features in this dataset that are highly correlated. In order to reduce redundancy and limitation of the data size, we did not use all the features provided by the dataset. Instead, we manually selected the features that we believed to be more relevant to our study.



*--Segment of Health Profile correlation heatmap*

The features we manually selected for this model are listed below:

| field | type | description |
| --- | --- | --- |
| Prop_18y | float | Proportion of residents ages 17 years and younger |
| Prop_65y+ | float | Proportion of residents ages 65 years and older |
| Prop_Lat | float | Proportion of Latino residents |
| Prop_uinA | float | Proportion of uninsured adults (ages 18 to 64 years) |
| Prop_uinC | float | Proportion of uninsured children (ages 17 years and younger) |
| Prop_rentr | float | Proportion of households that rent |
| Prop_smok | float | Proportion of adults (ages 18 years and older) who smoke cigarettes |
| No_libr | float | Number of public libraries |
| No_farm | float | Number of farmers' markets |
| No_hless | float | Estimated number of homeless individuals |
| Rte_resp | float | Available recreational space (acres per 1,000 population) |
| Prop_groc | float | Proportion of the population living in close proximity to a supermarket or grocery store |
| MHI | float | Median household income |
| Prop_edCG | float | Proportion of adults (ages 25 years and older) with a bachelor's degree or higher |
| Prop_edSC | float | Proportion of adults (ages 25 years and older) with some college education |
| Prop_edLH | float | Proportion of adults (ages 25 years and older) with less than a high school education |
| Prop_3rdg | float | Proportion of public school third graders who are meeting or exceeding California standards for English language arts & literacy |
| Prop_forb | float | Proportion of foreign-born residents |
| Prop_FPL1 | float | Proportion of residents living below 100% Federal Poverty Level |
| Prop_Whi | float | Proportion of white residents |

- **Social Distancing Metrics**
  SafeGraph provides useful foot traffic data, which we used in our study. Specifically, SafeGraph uses panels to detect GPS pings from anonymous mobile devices in different census block groups during the quarantine period. Due to the large size of the data and limited computing capacity, we decided to use the datasets from May 2020. Based on the scope of our study, we filtered out the locations outside LA County as well as removed the columns that display the hourly level.

  The features we selected from the SafeGraph dataset are listed below:

| field | type | description |
|---|---|---|
| origin_census_block_group | string | The unique 12-digit FIPS code for the Census Block Group. Please note that some CBGs have leading zeros. |
| date_range_start | string | Start time for measurement period in ISO 8601 format |
| device_count | integer | Number of devices seen in our panel during the date range whose home is in this census_block_group. Home is defined as the common nighttime location for the device over a 6 week period where nighttime is 6 pm - 7 am. Note that we do not include any census_block_groups where the count <5. |
| completely_home_device_count | integer | Out of the device_count, the number of devices which did not leave the geohash-7 in which their home is located during the time period. |
| part_time_work_behavior_device | integer | Out of the device_count, the number of devices that spent one period of between 3 and 6 hours at one location other than their geohash-7 home during the period of 8 am - 6 pm in local time. This does not include any device that spent 6 or more hours at a location other than home. |
| full_time_work_behavior_devices | integer | Out of the device_count, the number of devices that spent greater than 6 hours at a location other than their home geohash-7 during the period of 8 am - 6 pm in local time. |
| delivery_behavior_devices | integer | Out of the device_count, the number of devices that stopped for < 20 minutes at > 3 locations outside of their geohash-7 home |
| distance_traveled_from_home | integer | Median distance (in meters) traveled from the geohash-7 of the home by the devices included in the device_count during the time period (excluding any distances of 0). We first find the median for each device and then find the median across all of the devices. |
| median_home_dwell_time | integer | Median dwell time at home geohash-7 ("home") in minutes for all devices in the device_count during the time period. For each device, we summed the observed minutes at home across the day (whether or not these were contiguous) to get the total minutes for each device. Then we calculate the median of all these devices. Beginning in v2, we include the portion of any stop within the time range regardless of whether the stop start time was in the time period. |
| median_non_home_dwell_time | integer | Median dwell time at places outside of geohash-7 home in minutes for all devices in the device_count during the time period. For each device, we summed the observed minutes outside of home across the day (whether or not these were contiguous) to get the total minutes for each device. Then we calculate the median of all these devices. |
| median_percentage_time_home | integer | Median percentage of time we observed devices home versus observed at all during the time period. |

*social distancing metric documentation*

**Data handling:**
1. In order to merge this data to the other sources, we used Census Tract Locations (LA) from USC data socrata to decrypt the fips code in the columns of 'origin_census_block_group' to a city and neighborhood's name.
2. Since one city and neighborhood may have multiple fips codes, we aggregate different columns into city and neighborhood levels.
3. Since the data provided by SafeGraph uses data sampling, the real numbers are not very meaningful. Therefore, we calculated the percentage values of 'completely home device count,' 'part time work behavior devics,' 'full time

work behavior devices,' and 'delivery behavior devices' by dividing the 'device count.'
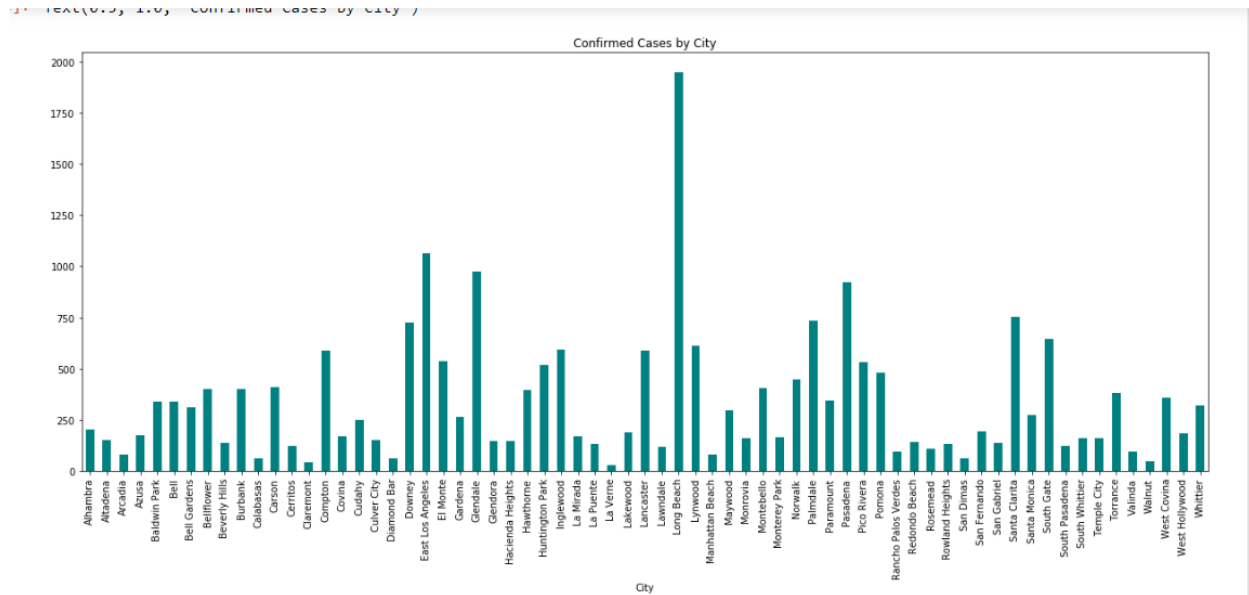
- **Los Angeles Neighborhood Map**

  Based on our exploration, we found the number of confirmed cases are greatly affected by the population size and size of the location. We think the number of cases decided by population density might be a better indicator to show the severity of COVID-19. From the 2018 health profile, we found the total population ('Tot_Pop') for different geographic areas. Although we are missing the area sizes for the different cities and neighborhoods, we were able to find the 'sqmi,' which indicates the square miles of a certain city/neighborhood, from this [website](website) . We used this feature to calculate our risk score.

  After merging the above datasets, our final dataset contained 1626 rows.

## 2. Data Exploratory

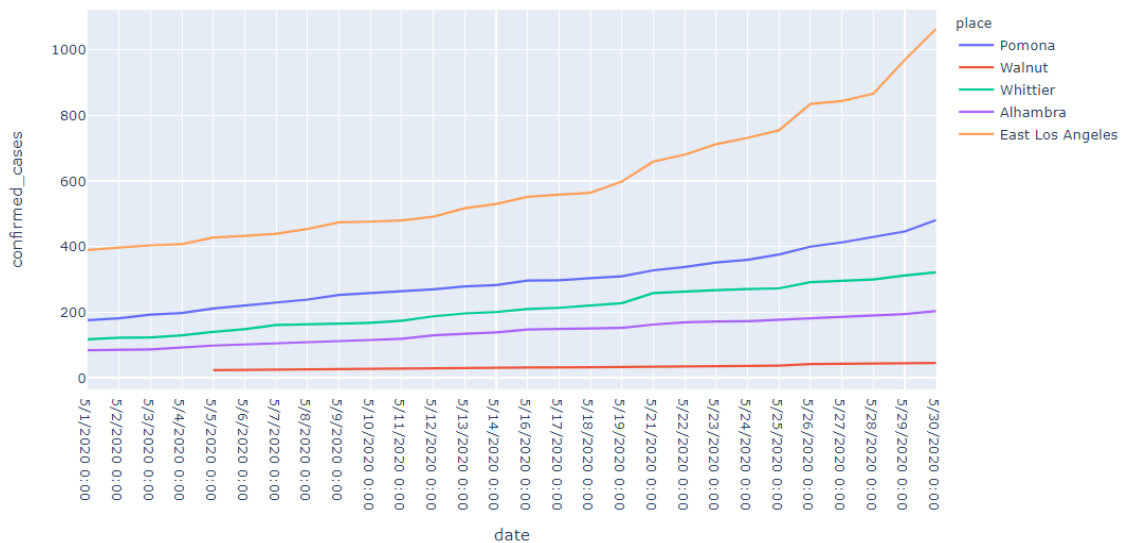**Number of Confirmed Cases by City**



In the beginning of our study, we used the data to see which cities contained the most and the least amount of confirmed cases in order to help us analyze and compute an appropriate risk score. Based on this graph, we found that the city of Long Beach has the most confirmed cases, being approximately 2000 confirmed cases.

**Top Five Cities in Los Angeles by Highest Number of Confirmed Cases**



We believe cities with a higher number of confirmed cases should be considered "high risk" for contracting COVID-19. From the data, we found that the top five cities in LA county, by the number of confirmed cases, are East Los Angeles, Glendale, Long Beach, Pasadena, and San Pedro. Because of the few missing data points in the dataset, the lines do not completely connect together, as seen in the graph above. For this discrepancy, we input the missing data by using mean imputation.

**Trend of Five Selected Cities in Adjacent Proximity to Cities with Highest Case Count**



We selected five cities and neighborhoods surrounding the city of Los Angeles to analyze the trend of total infected cases between 5/1/2020 and 5/30/2020. We wanted to see whether the number of total confirmed cases of a particular city affects the surrounding cities. The data shows that Pomona has the most confirmed cases out of the five selected cities and this number continues to increase. The dataset also demonstrates that the number of confirmed cases in

Walnut did not increase a lot as shown by the somewhat horizontal line in the graph visual. This line demonstrates that Walnut has a very few numbers of confirmed cases. As a result, it can be concluded that the confirmed case count of a city may not significantly affect the case count of nearby surrounding cities.

**Number of Confirmed Cases to Population Total by City**



Here, we analyzed the relationship between confirmed cases and the population count for each city in order to see if the total population number could affect the number of confirmed cases. The graph above shows that there is a relationship between confirmed cases and total population for each city. In other words, cities with a lower population show to have a lower number of confirmed cases while cities with a higher population show to have a higher number of confirmed cases. However, it is important to note that in addition to the total population number, the mobility data for individuals, who travel from other cities, is another important feature to incorporate into the study.
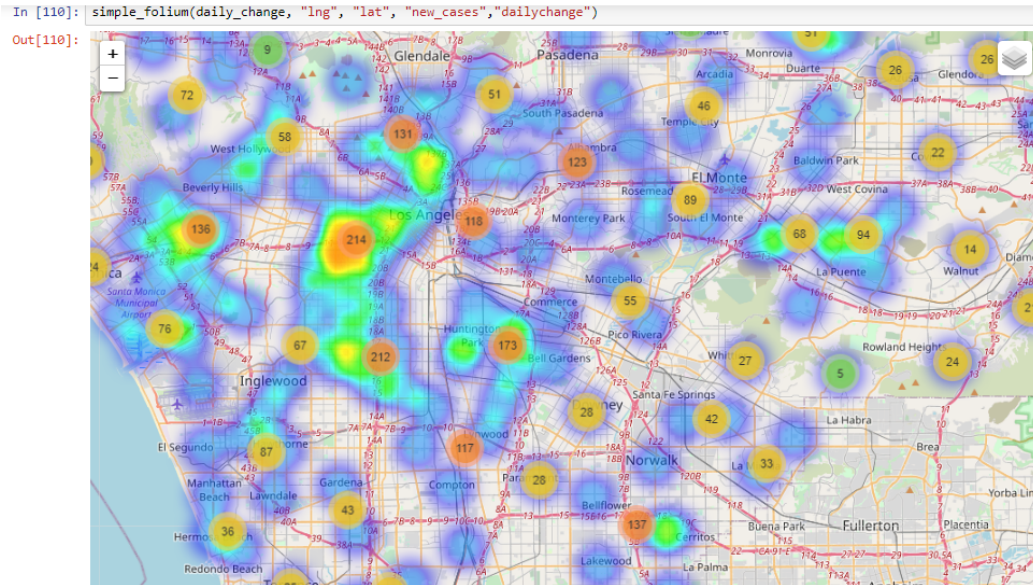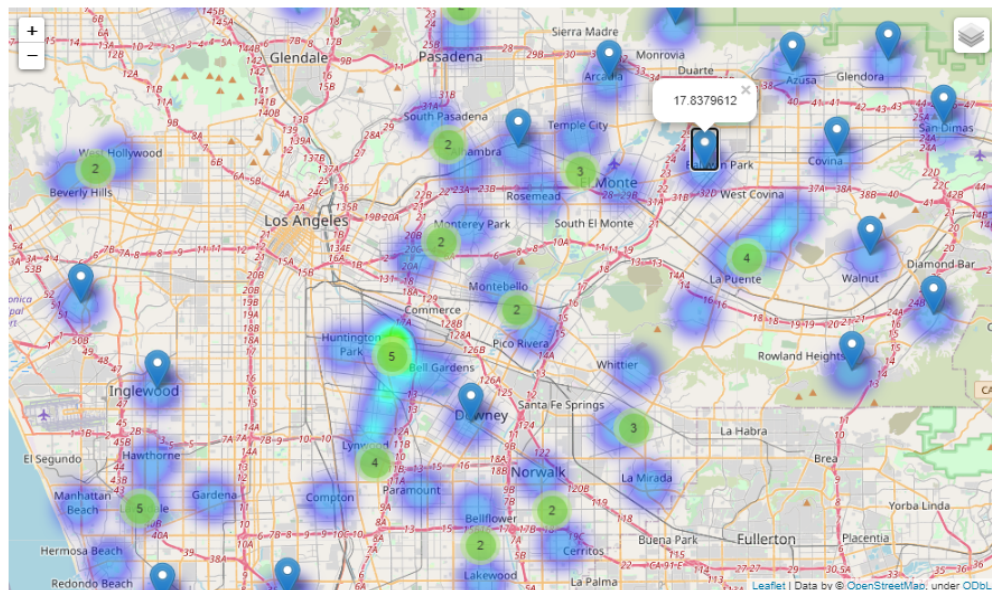
**Number of Confirmed Cases to Device Count by City**

This visual demonstrates the relationship between the number of confirmed cases and the device count for each city. 'Device_count' refers to the number of active technological devices within a given city. While the government recommends that people stay at home, there are still essential workers, including service or delivery workers, who must travel outside of his or her city for work. Because of this, we decided to see whether mobility would affect confirmed case count within a particular city in Los Angeles. As shown in the graph above, the data shows that although there is a relationship between confirmed cases and device count for each city, the relationship is not very strong. By way of illustration, some cities with a lower device count still show very high confirmed cases while cities with a high device count also show to have more confirmed cases.

**Heatmap: Number of Confirmed Cases Throughout Los Angeles**

```
In [110]: simple_folium(daily_change, "lng", "lat", "new_cases","dailychange")
```

Out[110]:



This visual demonstrates a heatmap by number of confirmed cases throughout Los Angeles. For this heatmap, we summed the daily new cases from 5/1/2020 to 5/31/2020 for each city and then incorporated its latitude and longitude. From this heatmap, we can see that the most confirmed cases are in the Downtown Los Angeles area.
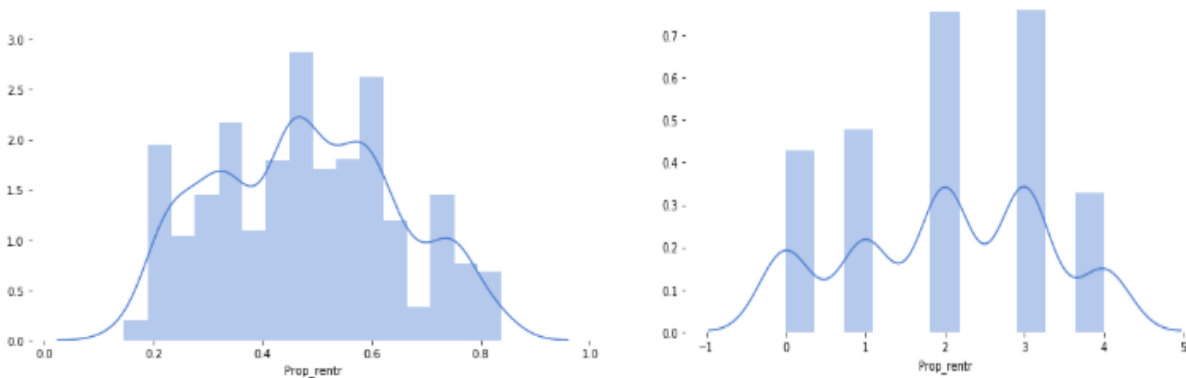
## Heatmap: Risk Score for Each City



This visual demonstrates a heatmap of the risk score for each city throughout Los Angeles, on a scale of 1-5, with 5 being the highest risk. In order to compute the risk score for this heatmap, the latest available data must be used. The risk score prediction for the following day can serve as an important indicator to show what cities people should avoid traveling on that given day.
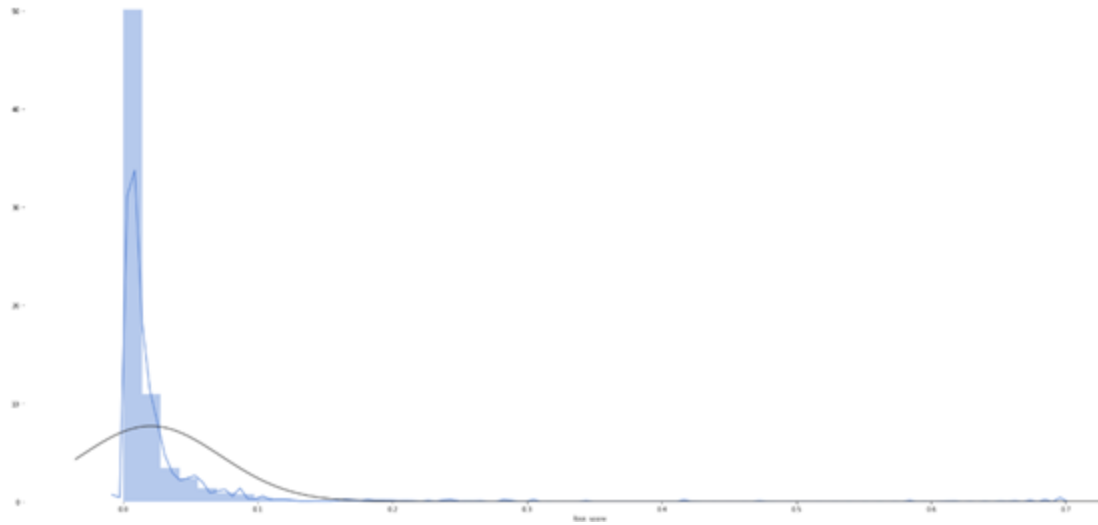
## Methodology

For our predictive models, we incorporated features that we selected from the pool and ran three regression models in order to compare each of their performance. The three regression models we used are Linear Regression Model, Random Forest Regression Model, and XGBoost Regression Model.

Before building the models, we normalized the risk score to scale all the scores between 0 and 1. Also, since the tree-based model works better with discretization data, we use K means based discretization to modify our features into five bins. Figure 1 shows the comparison of before and after discretization.
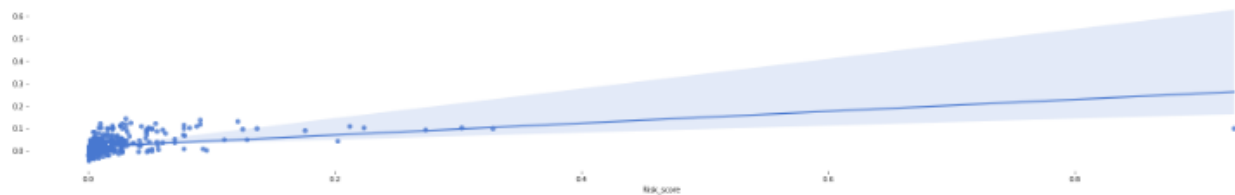
**Figure 1: Before (left) and After (right) Discretization**


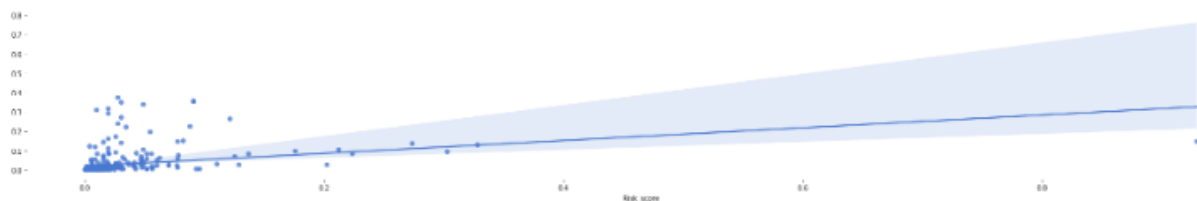
**Figure 2: Risk Score Distribution**
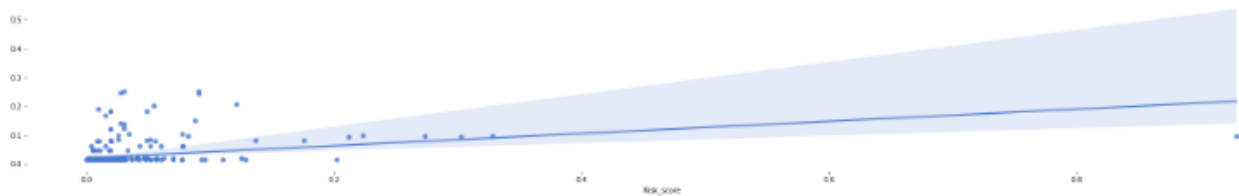
**Figure 3: Linear Regression**



Since the risk score is not normally distributed as the Figure 2 shows, the Linear Regression model does not have a very good predictive ability. However, it still gives the best predictive ability over three models since it gives a better prediction on the risk score between 0 to 0.2, which contains most of the value. As Figure 3 shows the testing set and prediction, the model has Mean Squared Error (MSE) of 0.0026.

**Figure 4: Random Forest Regression**

The Random Forest Regression model has an advantage of dealing with unnormalized distributed value, however, the model can be overfitted easily. So, even though we have an MSE of 0.002 for the training model, the prediction model has an MSE of 0.0039, which means the model is overfitting.

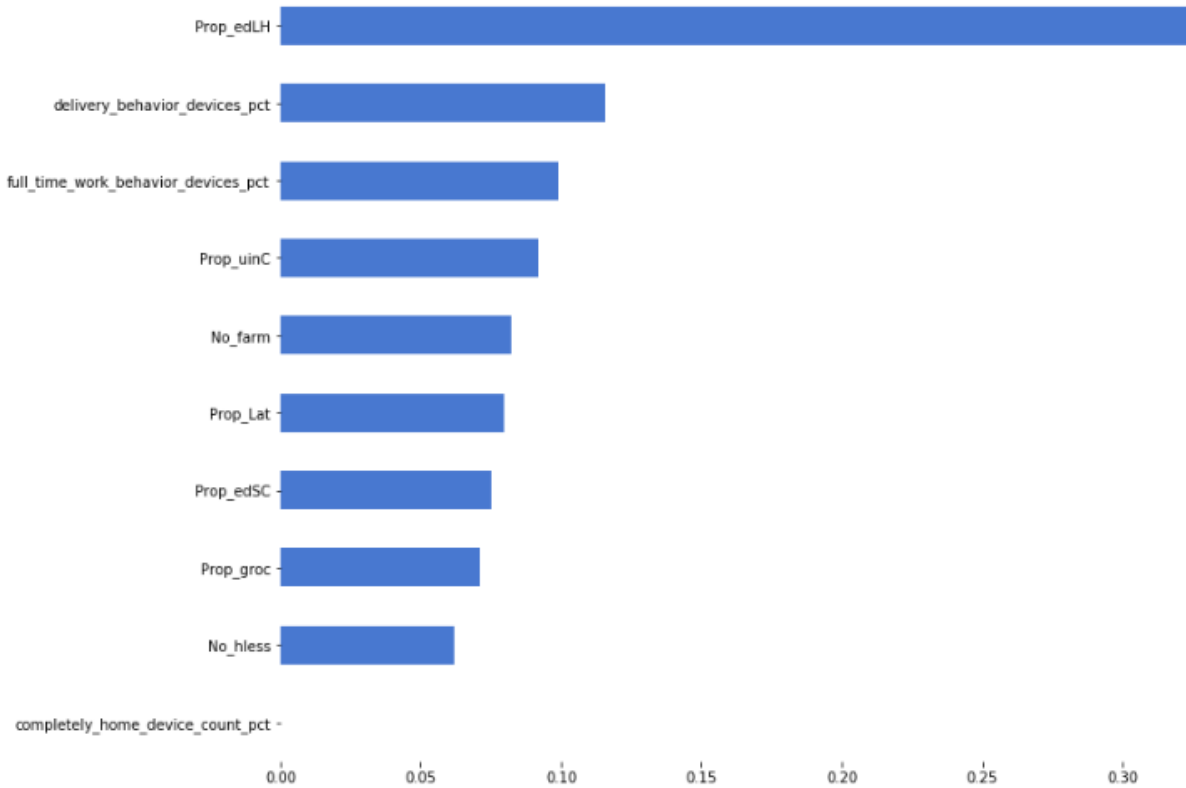**Figure 5: XGBoost Regression**



The XGBoost Regression has the ability to avoid overfitting, high flexibility, and continue on Existing Model. As a result, this model gives a better outcome than the Random Forest Regression model with MSE of 0.0029

Based on the observation by these three models, we can conclude that the Linear Regression Model has a low variance but higher bias and the XGBoost Regression Model has high variance but less bias. Because of this we formed a simple weight averaging ensemble method to further reduce variance, which aggregates 70% of the linear model prediction and 30% XGBoost model prediction for a better outcome of 0.0025 MSE.

## Results

**Figure 6: Feature Importance**

Based on Figure 6, features are ranked by importance given the coefficient associated with the risk score. The feature that influences our model the most is the Prop_edLH (Low education proportion) feature with 0.3 coefficient. The following features and their coefficients are full_time_work_behavior_devices_pct, delivery_behavior_devices_pct, and 'Prop_uniC' (uninsurance child proportion), which has a coefficient between 0.1 to 0.15, as well as, 'No_farm'(Number of farmers market), Prop_Lat (Latino), Prop_edSC (some college education), Prop_groc (living near grocery store), and No_hles (Number of Homeless), which are features that have a coefficient between 0.05 to 0.1.

Based on this graph, cities with a higher number of residents with low education, uninsured children, homeless people, farms market, people living closer to stores, and people that are highly mobile will tend to have a higher risk score.

**Risk Mitigation Recommendations:**

1. One short-term recommendation would be to keep a more watchful eye out for cities with high density and high mobility in order to keep the virus transmission in control, since these cities are more likely to have higher risk scores. This includes enforcing strict government ordinances to deter travel from individuals emerging high-risk areas.

2. Since low socioeconomic areas show to have a high confirmed case count, which can be caused by a plethora of reasons, such as not having the economic resources to stay home or get food delivered, another short-term recommendation would be to enact government programs to keep citizens at home while stimulating local economies. An example of a possible government program to enact would be subsidizing free grocery delivery to low socioeconomic areas to contain the spread of COVID-19.

3. Based on the model and the prominent important features used to compute the risk score, a long-term recommendation would be to focus on informing cities with lower educational levels about the dangers and proper safety protocols regarding COVID-19.

4. Another long-term recommendation would be to make more of an effort to ensure that citizens are not experiencing difficulty accessing healthcare and that those who have access to healthcare, yet choose to remain uninsured, are taking the proper safety protocols especially during times of pandemic.

## Acknowledgement

## References

Los Angeles Neighborhood Map
https://usc.data.socrata.com/dataset/Los-Angeles-Neighborhood-Map/r8qd-yxsr

Los Angeles County City and Community Health Profiles 2018
https://data.lacounty.gov/Health/Los-Angeles-County-City-and-Community-Health-Profi/capb-kusk

Latimes-place-totals.csv
https://github.com/datadesk/california-coronavirus-data#los-angeles-countywide-statistical-areasjson

Social distance metrics

https://docs.safegraph.com/docs/social-distancing-metrics