

Research on Criminal Activities in Seattle from 1996 to 2019

CSE 163 Final Project

Ruofeng Tang, Jin Ning Huang

Summer 2020

University of Washington -- Seattle

Table of Contents

Summary of research questions and short answers	4
Motivation and Background	6
Dataset	6
Methodology	7
Results	11
Challenge Goals	16
Work Plan	17
Testing	18
Collaboration	18
Graphs Appendix	19

Summary of research questions and short answers

On seattle.gov, there are some datasets about crimes and uses of force in Seattle from 1997 to 2019. We want to integrate different datasets to analyze how Seattle changed over the years. Here are the research questions:

1. Are there any relations between certain crimes and certain time periods/locations/finances? How did those relations change over the years? The raw data has various information regarding each crime, and we want to see if there are any patterns by using Python.

Answer: The number of crimes generally increased from 2008 to 2019, other than a drop from 2010 to 2013. It increases from 5 am to 12 pm and drops drastically from 12 pm to 5 am, and crimes rarely happen in February. In terms of location, crime usually occurs in south Seattle and Pioneer Square; however, the crime spread out almost evenly and recently the highest rate occurs in Northeast Seattle in 2019. In terms of finances, we make groups in terms of races and we don't see any unusual pattern as all of the lines are closely relative to each other (beside the native). The only exception is the low-income and poverty that both show the big spread of relatively.

2. Is there a shift in crime types from 1996 to 2019? Where did crimes often happen over the years, and did those locations change over time? We want to use geospatial data to check if high-crime-rate areas shifted over time.

Answer: Larceny-Theft is the major crime type every year and it exceeds other types by a large amount; From 1996 to 2007, crimes happened the most in downtown Seattle

around Pioneer Square, followed by nearby areas like Industrial District, Queen Anne and Capital Hill. After 2008, Capital Hill had a similar percentage of crimes to Pioneer Square. Crime percentage around Northgate was increasing consistently.

3. Does the use of force correlate to crimes in any way? By comparing this dataset to a dataset about the police's use of force from 2014 to 2019, can we check if the use of force is positively related to crime numbers, or are there any other factors causing the use of force?

Answer: Use of force is biased against male by a bit, and racial biases are disappearing over the years; it is generally not biased against locations.

Motivation and Background

The social justice system is always an ongoing topic in the US, and we want to contribute by analyzing crimes and police forces in Seattle. We want to know if Seattle is safe, where in Seattle is not safe, and if the police use force justly. This time of coronavirus is also crucial to our project. We want to conclude the past crime data because it will be incompatible with data after coronavirus, where everything has changed. If we get problematic results, we can learn from them and start anew after coronavirus.

Dataset

All the datasets are from the Seattle government website. All data in 2020 is excluded since 2020 have not ended, and most of our analysis is year-based.

- **“SPD(Seattle Police Department) Crime Data: 2008-Present”**

(<https://data.seattle.gov/Public-Safety/SPD-Crime-Data-2008-Present/tazs-3rd5>)

- Description: This dataset contains every crime reported from 2008 to present in Seattle. It has extremely detailed information, and we will focus on time of crime, location, and type. This dataset is our backbone as it contains the primary information for all of our questions.

- **“Use of Force”**

(<https://data.seattle.gov/Public-Safety/Use-Of-Force/ppi5-g2bj>)

- Description: This dataset records the police’s use of force from 2014 to 2019. It includes the level of force used, the subject’s race and gender, and the location. With this dataset, we can be able to answer the third question.

- **“Seattle Police Department Beats”**

(<https://data.seattle.gov/Public-Safety/Seattle-Police-Department-Beats/nxn-434b>)

- Description: This dataset contains geospatial data of all police beats. The “Use of Force” dataset provides locations as police beat codes, and we will use this dataset to convert the police beat codes as geospatial data.

- **“Seattle Crime Stats by 1990 Census Tract 1996-2007”**

(<https://data.seattle.gov/Public-Safety/Seattle-Crime-Stats-by-1990-Census-Tract-1996-2007/e3zj-s4zh>)

- *Description:* This dataset contains brief information about crimes from 1996 to 2007 in Seattle. It contains the crimes per type each year in each census tract from the 1990 census. This dataset can reinforce the primary dataset of the missing pieces from it. Once again, this information can help us be able to answer all of those questions.

- **“Census Tracts 1990”**

(<https://data.seattle.gov/Land-Base/Census-Tracts-1990/qtbw-f6xb>)

- *Description:* This dataset contains geospatial data of 1990 census tracts. The “Seattle Crime Stats by 1990 Census Tract 1996-2007” dataset provides locations as 1990 census tract numbers, and we will use this dataset to convert the 1990 census tract numbers as geospatial data.

- **“Housing Cost Burden by Race”**

(<https://data.seattle.gov/dataset/Housing-Cost-Burden-by-Race/2cxc-b9bd>)

- *Description:* Displacement risk indicator showing how many households within the specified groups are facing either housing cost burden. We will use this dataset to get the information of the income, housing costs, races, and possible tenure. With that data information and some other crime information, we can be able to answer the first question.

Methodology

Task 1: Focus on “SPD(Seattle Police Department) Crime Data: 2008-Present”

a). When does crime happen over the years? (2008 to 2019)

We will separate crimes by year, and round up the time by hours. For example, a crime that happened at 15:30 is in the category “15:00 - 16:00”. Then we will make a bar chart: x-axis is the time period in a day, y-axis is the number of crimes, and in each time period, there are multiple bars representing each year from 2008 to 2019, where we have detailed data. This graph is a bit complicated to see a trend visually, so I also included four different simple line charts: crime numbers by year/month/year and month/hour.

b). How did crime types change?

We will separate crimes by year and type. Then we will make a line chart: x-axis is the year, y-axis is the number of a type of crime, and there are multiple lines in different colors representing different types. From this graph, we can see what types of crime happened more and how it changed. There are too many types recorded, so we chose the ten crimes with the highest numbers. To see a more clear pattern, we also plotted the line chart with “Year and Month” as x-axis to have more data points.

c). Where did crime usually happen? (**multiple datasets**)

We will use the “Police Beats” dataset to attribute detailed locations to police beats, and plot a map of crime numbers from 2008 to 2019 by police beat. This map tells us where crimes usually happen.

Task 2: Introduce “Seattle Crime Stats by 1990 Census Tract 1996-2007”

a). How did crime types change from 1997 to 2019? (multiple datasets)

We want to combine the dataset from 1997 to 2007 and the dataset from 2008 to 2019. Since two datasets have different notations on crime types, we need to manually decide which types are similar. We will then have a consistent line chart of crime number vs. the year by type. We will perform multiple lines plots to see each individual trend throughout the year.

b). How did crime locations shift from 1997 to 2019? (multiple datasets)

We want five maps together to contrast: crime ratios by 1990 census tract in 2000, crime ratios by 1990 census tract in 2005, crime ratios by 2010 census tract in 2010, crime ratios by 2010 census tract in 2015, crime ratios by 2010 census tract in 2019. (There is a 2000 census, but our data uses 1990 census tracts from 1997 to 2007, so we do not use 2000 census tracts for 2000 and 2005 graphs). We will use the “1990 census tracts” dataset to give 1997 to 2007 crime data locations, and use “2010 census tracts” dataset to give 2008 to 2019 crime data locations, which was performed in task 1c. Different census tracts make the maps harder to compare, but if we only use 1990 census tracts, newer data would be inaccurate as well. If we plot based on more accurate census tracts, we can at least analyze visually. For each map, we want for each census tract the ratio of crimes in that census tract to crimes in that year. If we use crime numbers, there might be a general rise or fall, and then we cannot know if crime locations shifted. From these five maps, we can see how crime locations shifted over the years.

Task 3: Introduce “Use of Force” data

a). Is the police use of force biased against gender or race?

The use of force is recorded in levels, and by “use of force continuum”:

Level 1: Officer Presence, i.e. Officer as an authority

Level 2: Verbal Commands, i.e. Commanding a subject

Level 3: Empty Hand Control, i.e. using empty hand to control a subject

In the dataset, most recorded use of force is level 1, which is not hostile to the subject. A relatively small part is in level 2 and level 3, which can be hostile and potentially biased. We want to set the level 1 use of force as a “population” and check the percentages of each sex or race. Then we check the same percentages for level 2, level 3 use of force. Then we can have two line charts: x-axis is the year, y-axis is the ratio of percentage for level 1 to the percentage of level 2 and level 3 together. The first line chart is about gender, and we have multiple lines, each representing a gender; the second line chart is about race, and each line representing a race. Theoretically, all the ratios should be around 1; if not, there must be some biases. The line chart with races did not work because too many records are “Not Specified”. We changed the formula and calculated the ratio of level 2/3 use of force to level 1 use of force. Then we have a different line chart with refined formula.

b). Is the police use of force biased against locations? (multiple datasets)

Plot the crime maps from task 2c at 2015, 2019 again, and plot use of force maps at 2015, 2019 based on police beats. Use the “Seattle Police Department Beats” dataset to convert the beat codes in use of force dataset to a beat map, then plot use of force maps. We

decided to plot the ratio of level 2, 3 use of force to level 1 for each police beat, but a few police beats did not have level 2, 3 use of force in a given year.

Task 4: Introduce “House Cost Burden” data

a). Selecting the factors of house income, area location, and races

We will filter/select the dataset information that is relating to housing income/area median income, geometry area location (by County), and races. We could also group up to aggregate the average income and/or by races.

b). Are there any relationships between the finance/locations/races and the crime rate?

Plot the regression graphs for finance that would categorize the different levels of income; x-axis would be the year; y-axis would be the income. We might group by 3-5 different levels of income, and we can compare the other plots that display the crime trends. We can also plot another similar setting of the regression except replacing the income for house cost.

Plot the grouped bar graph and multi-line graph for race in terms of income and house cost separately (two graphs). The x-axis would be in year, y-axis would be the unit of USD, and each group per year is categorized uniquely by race. We can compare those plots to the other plots that display the crime trends.

Results

Research Question 1: Are there any relations between certain crimes and certain time periods/locations/finances? How did those relations change over the years? ([Click for graphs](#))

The plan was to have a master graph of (crime numbers) vs. (year and hour), but the product is too complicated to analyze, and the visualization does not help to

understand(1a1.png). So, I have four separated simple line charts instead: crime numbers vs. year/month/hour/year and month.

Crime number vs. year graph(1aYear.png) shows a general increasing number, but there is also a drop around 2010 to 2013. It does not tell much since there are only around ten points on the graph, so we made another with each year and month as a point(1aYear and Month.png). It varies too much to get a convincing conclusion, so we cannot conclude any trend of crimes in Seattle over time.

Though these two graphs are frustrating, the other two give interesting results. Crime number vs. month graph(1aMonth.png) has a low point in February, while the crime number of other months are at the same level. We print out each point, and February is at 58276, while the second to last is 64837; Other than February, the rest ranges from 64837 to 69817. Taking the midpoint of other months, February is only around 86% of their numbers.

Crime number vs. hour in a day(1aHour.png) gives an oscillating graph: the numbers goes from the low point at 5 am to the high points at 6 pm or 0 am, then it drops drastically from 0 am to 5 am. The range is from around 10000 to around 60000, a shockingly large number.

As for the location, the Pioneer Square area had a ratio of 8.7% crime reports, which is the highest than other areas in 1996 (crime_ratio_1996.png). Most crime rates usually occur in the South side of Seattle around that time too. Throughout the year, crime slowly increased until the year 2010 where the map shows a big jump of crime rates in Seattle overall (crime_ratio_2010.png). However, the Pioneer Square area still retained the highest crime ratio which was around 8%. The most recent year, 2019, tells us that crime can still occur at any place and not just Pioneer Square. The map figure (crime_ratio_2019.png) shows that Northeast Seattle had the highest crime ratio, and the next highest was relatively close such as Ballard, Wallingford, Queen Anne, Capital Hill, and etc.

As for the finance factors, it's really hard to analyze and digest the information down because there can be no correlation of income and housing cost that match the pattern with crime reports. For example, the line figure ('income_level_overall.png') shows the highest poverty occurring in 2017, we can't confirm if the majority of the crime would be caused by the poverty group. Similar to the housing cost, people who are in burden for paying the highest housing cost wouldn't mean they are more likely to do crime because again we can't see any clear correlation without needing to assume. The fault is due to the lack of description on the Crime and Police Department datasets; however, we could be able to get a better answer if we have more time on this project.

If you looked at all of the bar graphs and line graphs (race based on finance factor) that consist of race categories, you can see that most data values are all over the place, and don't show any clear correlation. We can be sure that there is no correlation of certain races would be most likely to do the crime with the factors of finance.

Research Question 2: Is there a shift in crime types from 1996 to 2019? Where did crimes often happen over the years, and did those locations change over time? ([Click for graphs](#))

The crime types are dominated by Larceny-Theft. Every year Larceny-Theft has around 25000 cases, while the second largest number is Assault Offenses, less than 10000 cases per year. There are too many types recorded in the dataset, so we only plot the ten most often crime types(1b11.png). The tenth most often type, driving under the influence only has around 1000 cases per year, far less than Larceny-Theft or Assault Offenses. Therefore, ignoring less often cases does not influence the research. The graph does not give clear trends for each line, because the top Larceny-Theft line takes up too much space. So, a new graph without Larceny-Theft is created(1b12.png). The most noticeable feature on this graph is the general increase in 2014. The Fraud Offenses and Motor Vehicle Theft peaked in 2014, while Trespass

of Real Property keeps increasing since 2014 after a fixed number from 2008 to 2014. Other types do not have a clear turn in 2014, but most are increasing around that time. Unfortunately, data can only tell us what is happening, but it cannot give us the reason. To get more data points, we plot the same two graphs with x-axis being “Year and Month” (1b21.png, 1b22.png). The new graphs are messier, but each line keeps oscillating with a period of about 3 or 4 months. Interestingly, all lines have such oscillations, but the periods are slightly different.

If you want to compare individual crime trends, you can see on the graph figure (‘combined_crime_report.png’) assault offenses and non-residential burglary had been ridiculously increasing in post-2007 and are still growing. Interestingly, vehicle theft declined greatly as it could have been with better car safety technology being improved. What’s also interesting is that even though Larceny-Theft is still considered the highest crime type, you can see it declining overall and could be predicted to pass through below of the assault offenses.

Research Question 3: Does the use of force correlate to crimes in any way? By comparing this dataset to a dataset about the police’s use of force from 2014 to 2019, can we check if the use of force is positively related to the crime numbers, or are there any other factors causing the use of force? ([Click for graphs](#))

The Use of Force dataset is hard to process because we do not have two related reference frames: population and crime. We only have the police’s use of force records categorized by subject gender, race, and location, but we do not know the respective overall populations and we do not have crime records to link with use of force records. Therefore, we generally use ratios to seek patterns. Our first attempt to summarize something meaningful is this formula:

$$Ratio\ of\ Bias = \frac{\frac{use\ of\ force\ level\ 2\ and\ level\ 3\ in\ a\ given\ year\ to\ a\ group}{use\ of\ force\ level\ 2\ and\ level\ 3\ in\ a\ given\ year}}{\frac{use\ of\ force\ level\ 1\ in\ a\ given\ year\ to\ a\ group}{use\ of\ force\ level\ 1\ in\ a\ given\ year}}$$

If the police is unbiased, the ratio in numerator and the ratio denominator should be close, and the Ratio of Bias should be near 1. This formula works when we categorize by gender (Subject_Gender.png). The line for male and female both come around one, but the line for “Not Specified” looks strange. The reason is that out of 10740 cases in Use of Force dataset, only 229 cases identify as “Not Specified”. So, if we remove “Not Specified”, we have two lines around one. The line for male is slightly over one, and the line for females is much lower than one. It means the Use of Force is biased against male, since the percentage of male receiving use of force level 2 and 3, Verbal Commands and Empty Hand Control, are higher than the percentage of them receiving use of force level 1, Officer Presence. Similarly, the Use of Force favors females as they receive use of force level 2 and 3 less.

When we plot the ratio by race, the plot looks more abnormal than the one by gender (Subject_Race.png). Many lines exceed one, and the highest line reaches a crazy "70. To make matters worse, the highest line is “Not Specified”, and the number of cases is 1910, the third largest category after “White” and “Black or African American”. When we tried to remove “Not Specified”, the new graph was still uninterpretable (Subject_Race_New.png). Then we scale the ratio differently:

$$\text{New Ratio} = \frac{\text{use of force level 2 and level 3 in a given year to a group}}{\text{use of force level 1 in a given year to a group}}$$

Note that the new ratio is the old ratio times an array of constants:

$$\frac{\text{use of force level 2 and level 3 in a given year}}{\text{use of force level 1 in a given year}}$$

And we plot a new line chart with the four largest races, excluding two races with around 100 cases, and “Not Specified” (3a.png). The four lines change up and down, and they are close, suggesting unbiased use of force by the police. But the more accurate conclusion should be: if the police are involved, in other words a use of force level 1 already happened, then the police are unbiased against races; however, we do not know if police involvement is biased. For

example, if the police thinks a situation needs attention, he or she will not be biased about using verbal commands or empty hand control; however, we cannot conclude if the police is biased when evaluating a situation. In fact, the recorded 4351 cases against white people and 3489 cases against black people may suggest such biases, since white people are 65.7% of Seattle populations while black people are 7%.

(<https://www.seattle.gov/opcd/population-and-demographics/about-seattle#raceethnicity>)

Combining the police beats map, we have 6 different maps to see if use of force is biased against locations. First, there are maps about the percentage of crimes happening in each police beat in 2015 and 2019 (3b2015_0.png, 3b2019_0.png). They are compared side by side to the maps about the new ratios in each police beat in 2015 and 2019 (3b2015_1.png, 3b2019_1.png). The 2015 maps are similar, suggesting level 2, 3 use of force happens the most in downtown Seattle near Pioneer Square, which is reasonable considering the large number of crimes there. However, the 2019 use of force graph shows a noticeable region in West Seattle. To check if this is unique, we plotted 2018 use of force graph, hoping to find a graph similar to the 2019 graph (3b2018_1.png). It turns out that the 2018 graph is vastly different from the 2019 graph. We think this is due to the low number of level 2, 3 use of forces, so we plotted the general use of force graph from 2014 to present (3b.png). This corresponds to the crime maps where only downtown Seattle has a high percentage. Therefore, we conclude that the use of force is not biased against locations.

Challenge Goals

We focused on multiple datasets, and we dealt with messy data from these datasets. For multiple datasets, we have six different datasets, four of which are main datasets we studied thoroughly, and the rest two are spatial data we used to combine the main four datasets. Two of

the four main datasets are kind of independent: Use of Force and House Burden by Race. Their contents are largely different, requiring different methods to get out useful information. They post different challenges to us because we could not expect the problems with each dataset. The other two main datasets are SPD crime data from 2008 to present and crime data from 1996 to 2007. The first is kind of a continuation to the second, so they can be processed with similar methods. But the newer records are an improvement to the older ones, and we had to figure out how to improve the older ones to match them. The two side datasets are used to help with mapping, since the four main datasets do not have any spatial data. Yet we successfully plotted maps with three main datasets using the side datasets.

As for messy data, it comes from coordinating multiple datasets. All our datasets are in .csv form or .shp form, but the contents require lots of pre-processing. One recurring problem is the time. Most of our datasets have information about time, but they come in many different forms. Some are too specific, and we must cut out the correct portions to use them; some are irregular, and we must find a way to consider all possibilities; some are contained in strange object types, and we must find creative ways to get them working. Both SPD Crime data from 2008 to Present and Use of Force datasets come as individual cases, so we need to count them in some ways. We use several dummy variables and count them under different conditions, converting them between different object types.

Work Plan

We are planning to have frequent meetings and split up the work equally at our own time. Our meetings would be appointed by every Monday, Wednesday, and Friday until the end of the summer quarter (or when the project deadline is reached). During meetings, we will virtually call to review our checkpoint on our process, to help troubleshoot any problems that

arise, and to test the edge cases on our codes. We will host meetings right after the CSE 163 class session and will take approximately an hour. We will use GitHub as our Version Control System.

For the rest of the days, we would work on our tasks individually. We would expect to do the project at a minimum of 2-3 hours and no more than 5 hours on our tasks daily. The task would split up evenly to give us a fair amount of work. Those tasks are based on the list of the methodology task section. Ruofeng Tang would focus on doing Task 1 and 3. Jin Ning Huang would do Task 2 and 4. We will write our report together afterward based on the analysis and plots we each made.

Work Plan Evaluation:

Generally, our coordination is fantastic, but the workload is beyond our imagination. Every meeting happened other than the first one. We checked our results and solved many questions, and we constantly communicated through Messages. In our individual times, we worked on our respective parts. Both Ruofeng Tang and Jin Ning Huang seriously underestimated this arduous project as they painstakingly checked documentations for tons of unheard methods. Ruofeng spent around 6 or more hours every other day, and frustratingly, many hours of studying documentations and testing did not go into the final research. Jin spent almost entire days working on his own tasks, more common due to solving the troubleshoot. Jin also looks at documents to determine the datasets for many data cleaning; such as getting rid of NAN values, renaming columns, and filtering out unnecessary data information. Also, they meet a serious real-world problem: there are few obvious patterns in his datasets. So the actual project develops further from the methodology.

Testing

This project is all about visualizing, and graphs themselves can tell if the coding is correct. For example, we can look at the legend to see if the dimensions are correct. During our coding, we usually print out the results after each step and check if everything is right. Other than printing the results themselves, we sometimes print out the length of our results. For example, when we have a Dataframe about crime numbers per month, the length should be 12; when we plot crime numbers by police beats, the length should be 51 since there are 51 police beats. Sometimes the length of the object about police beats is less than 51, then we will find the length to the related objects and explain why.

Collaboration

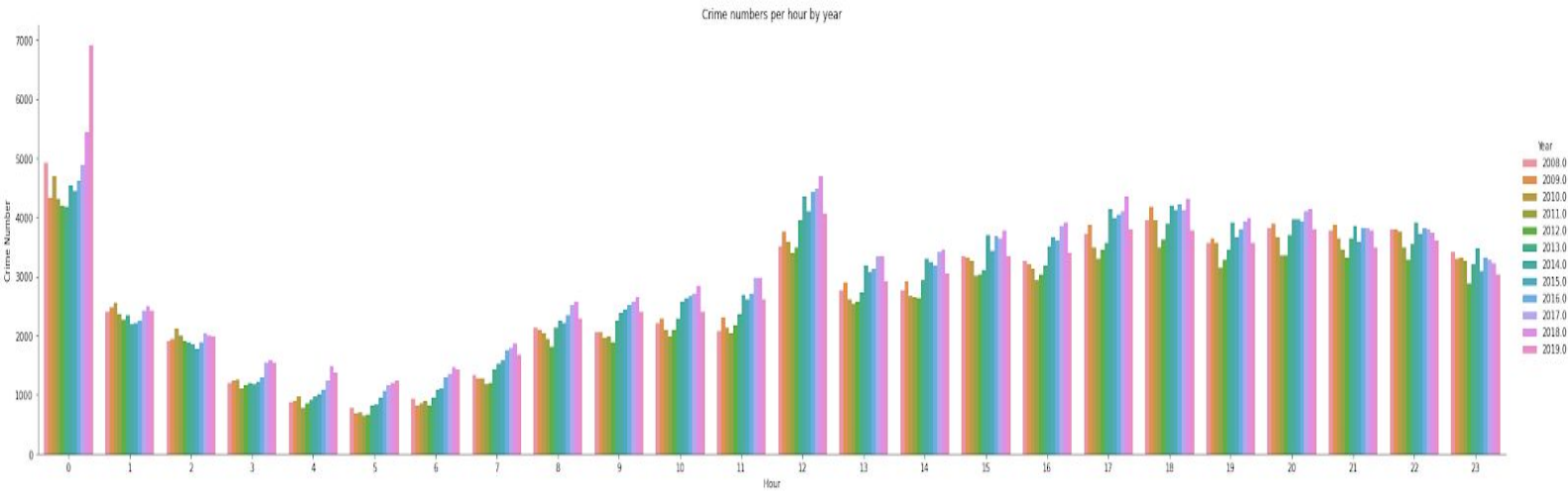
Ruofeng did not receive help from anyone other than groupmate and course staff, but he researched online intensively. Usually, he checks documentations or finds new useful methods on Stack Overflow and then checks their documentations.

Jin also didn't receive any help beside Ruofeng and researched online intensively as well. More often, he checks Python built-in functions and dataset documentations. When he would look at Stack Overflow when he found better methods and check the document if needed.

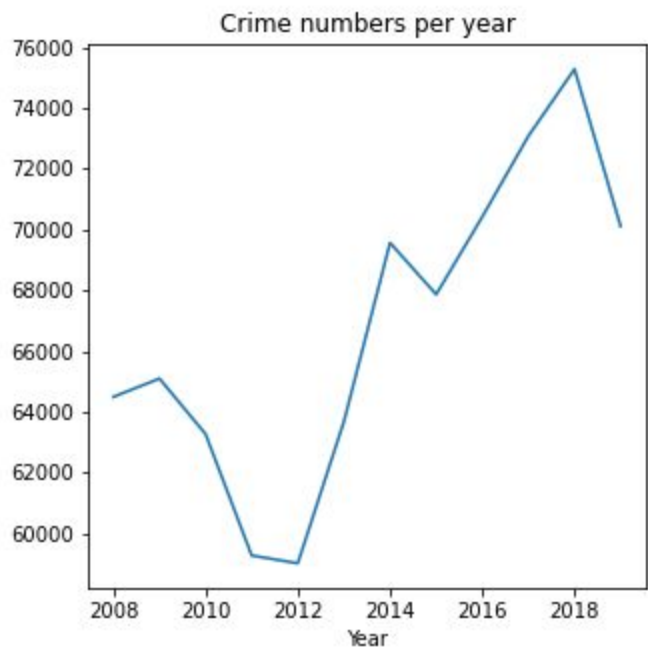
Graphs Appendix

Research Question 1 ([Click to go back](#))

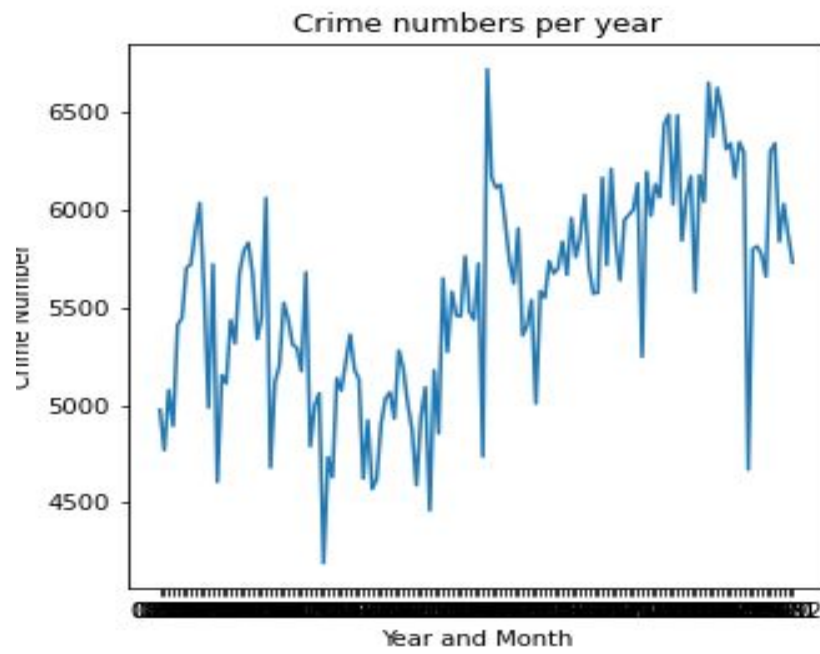
1a1.png



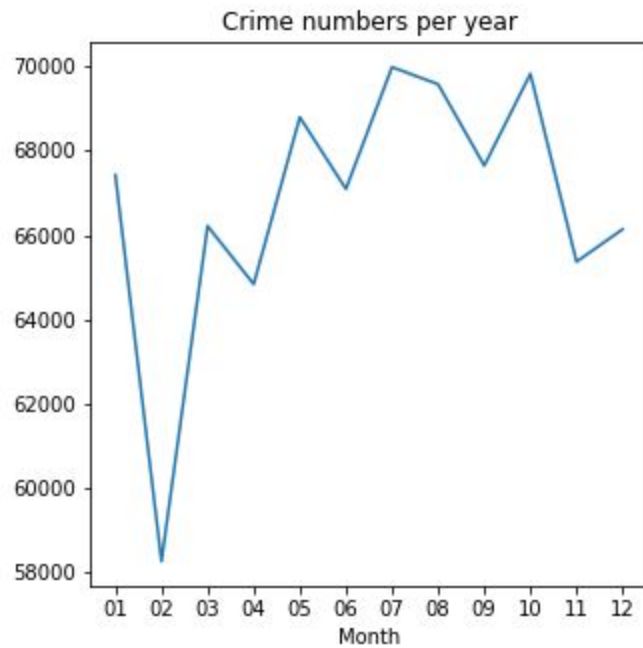
1aYear.png



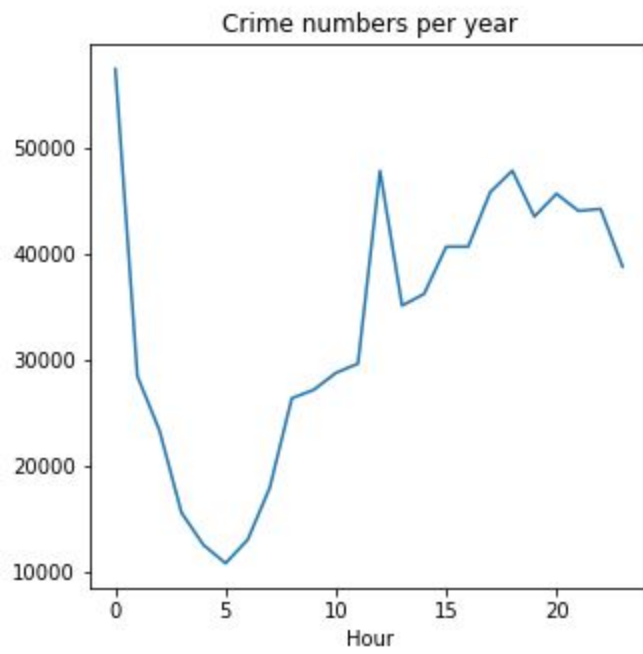
1aYear and Month.png



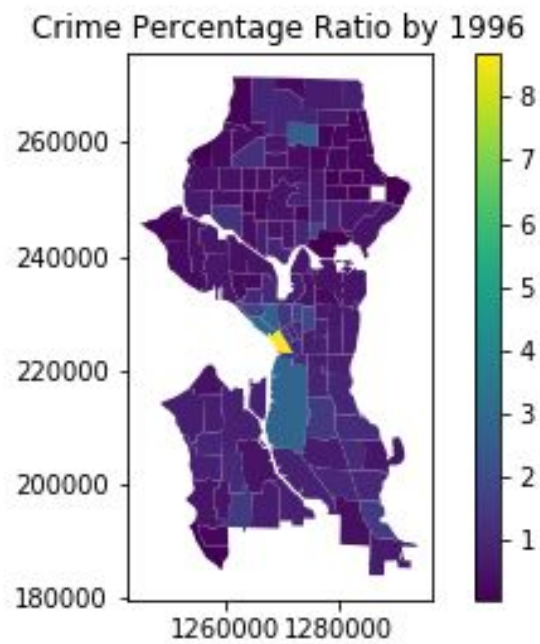
1aMonth.png



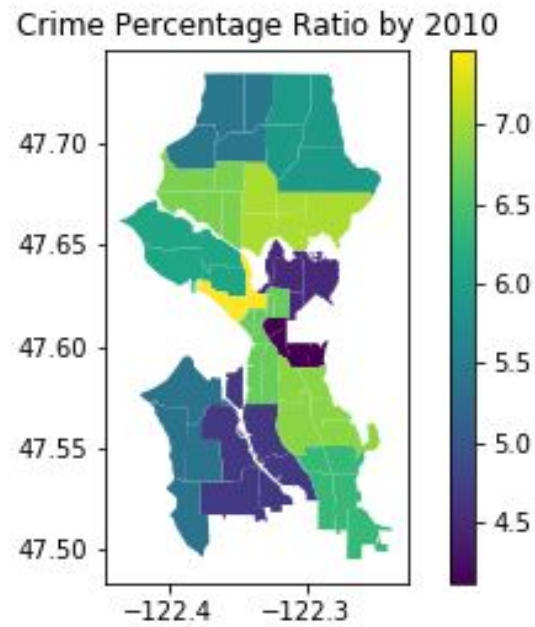
1aHour.png



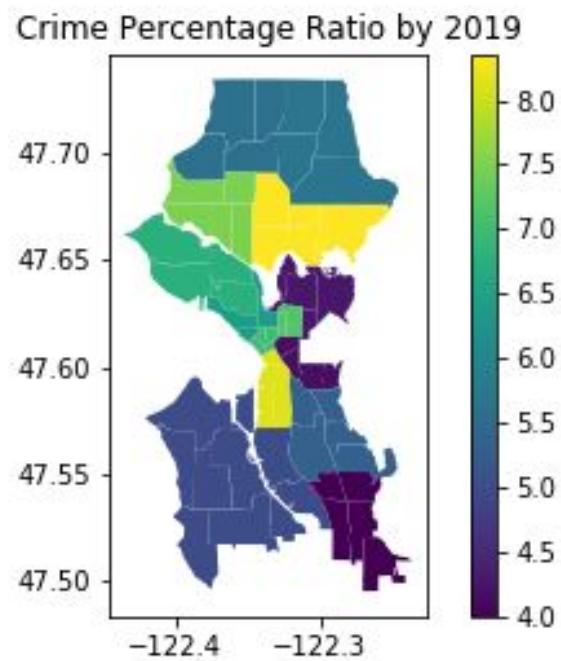
crime_ratio_1996.png



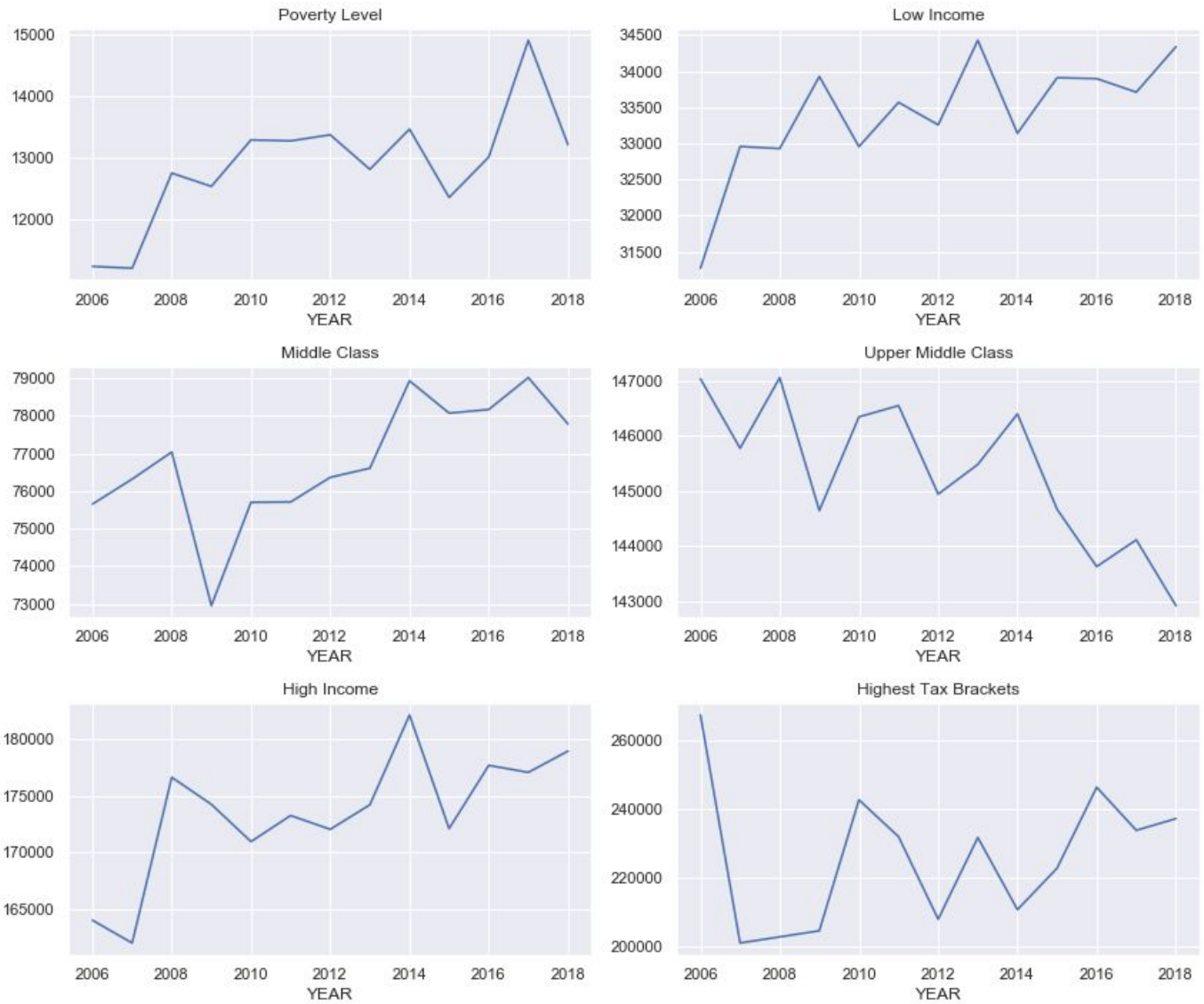
crime_ratio_2010.png



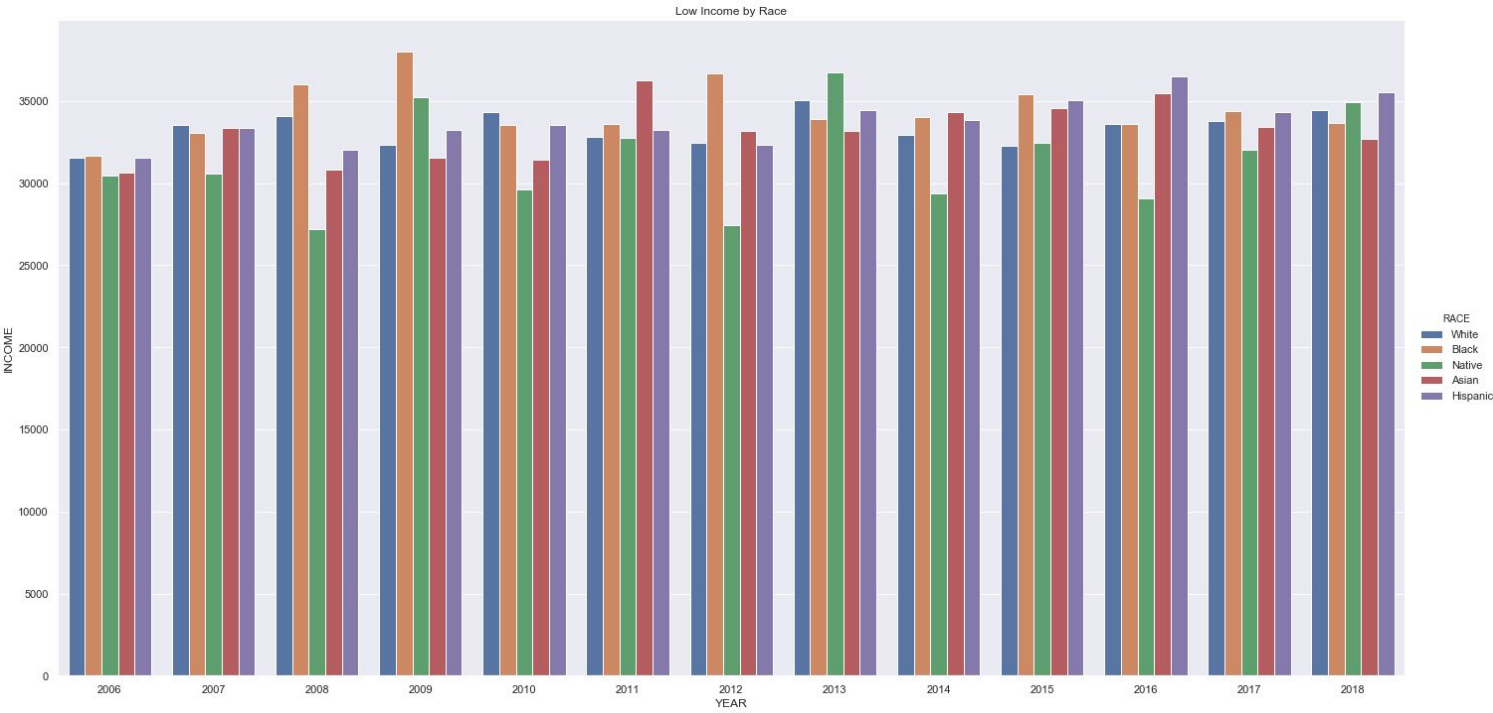
crime_ratio_2019.png

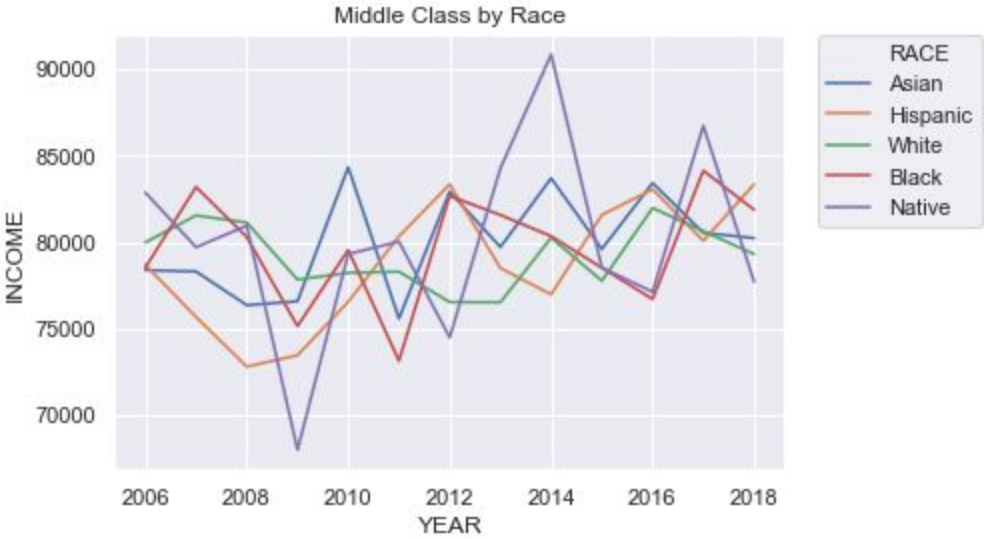
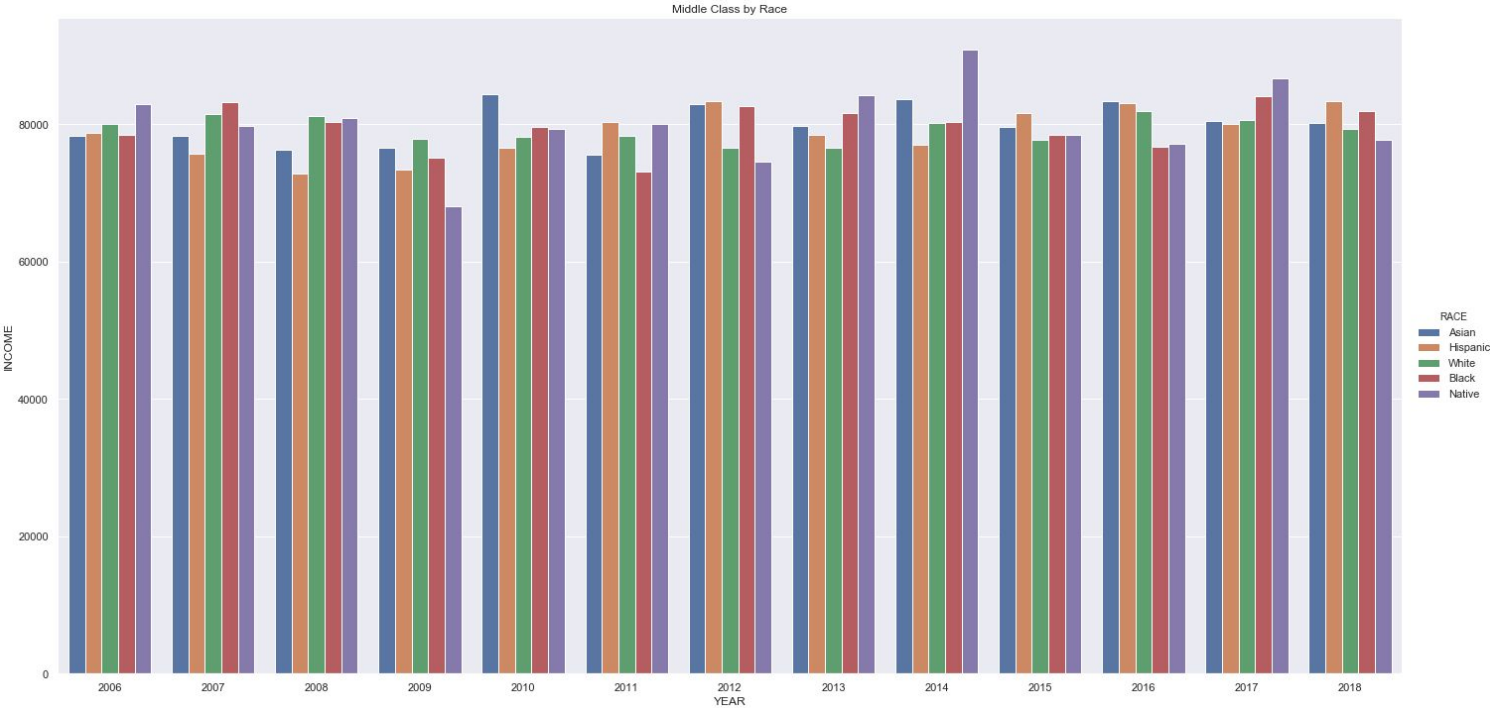


income_level_overall.png



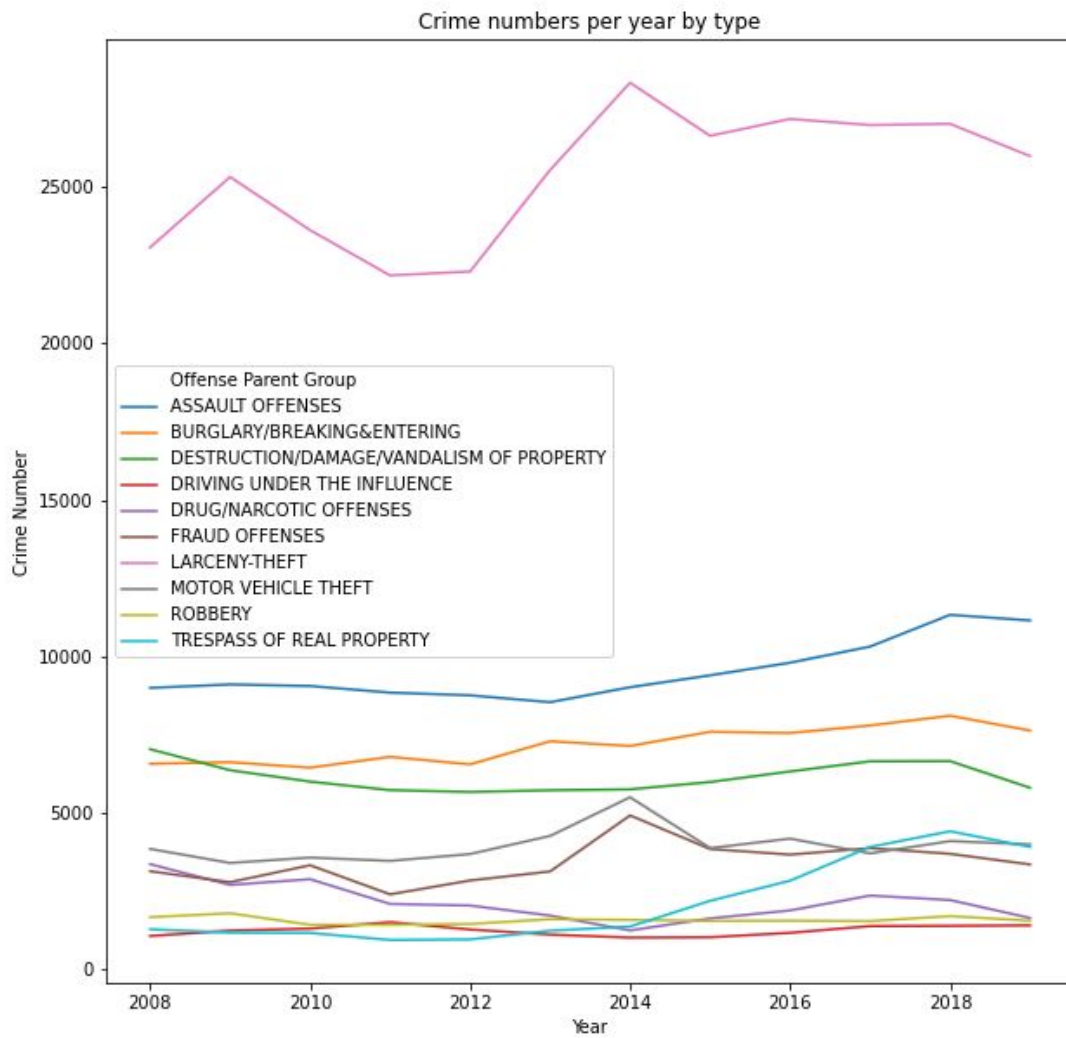
race based on finance factor (will pick only two groups to save space; look more on GitHub)



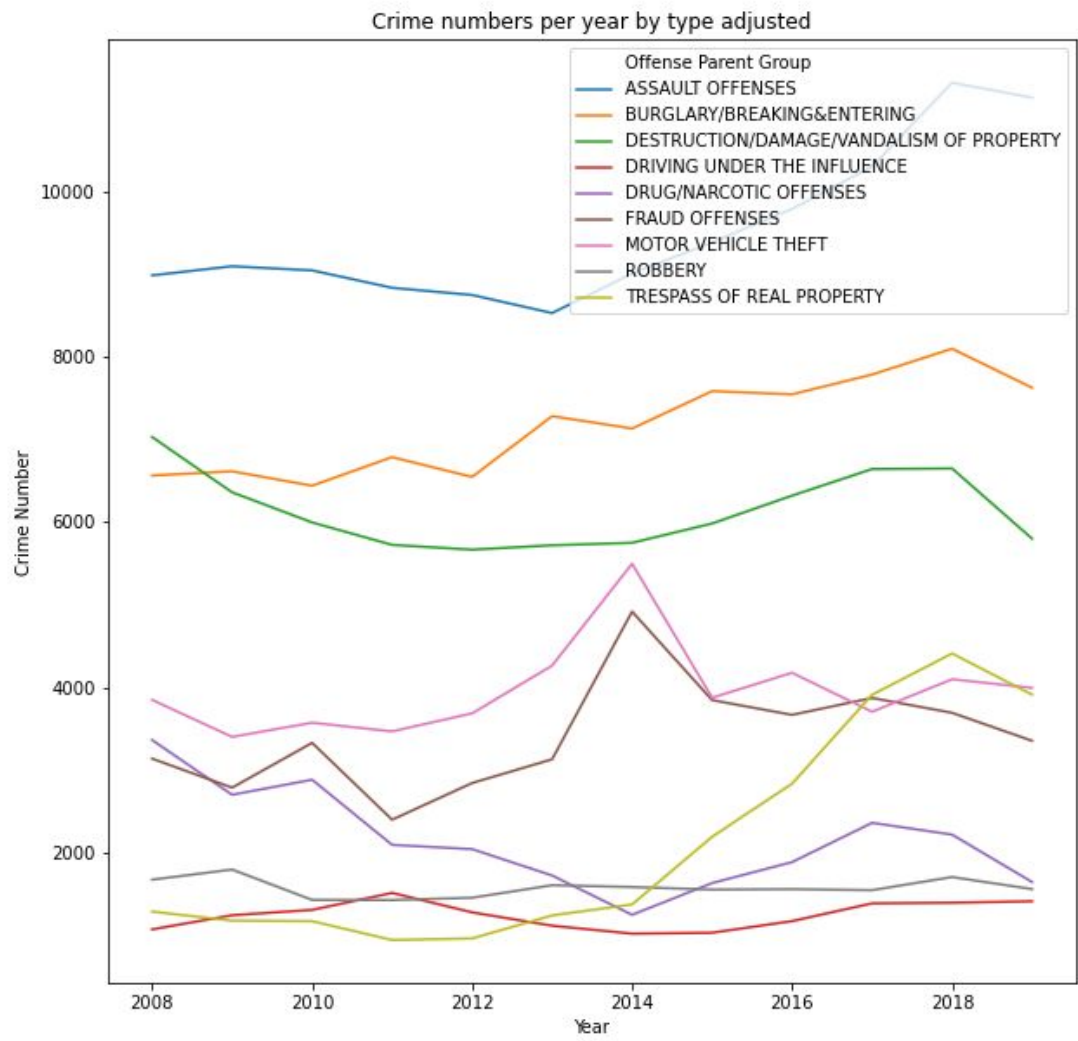


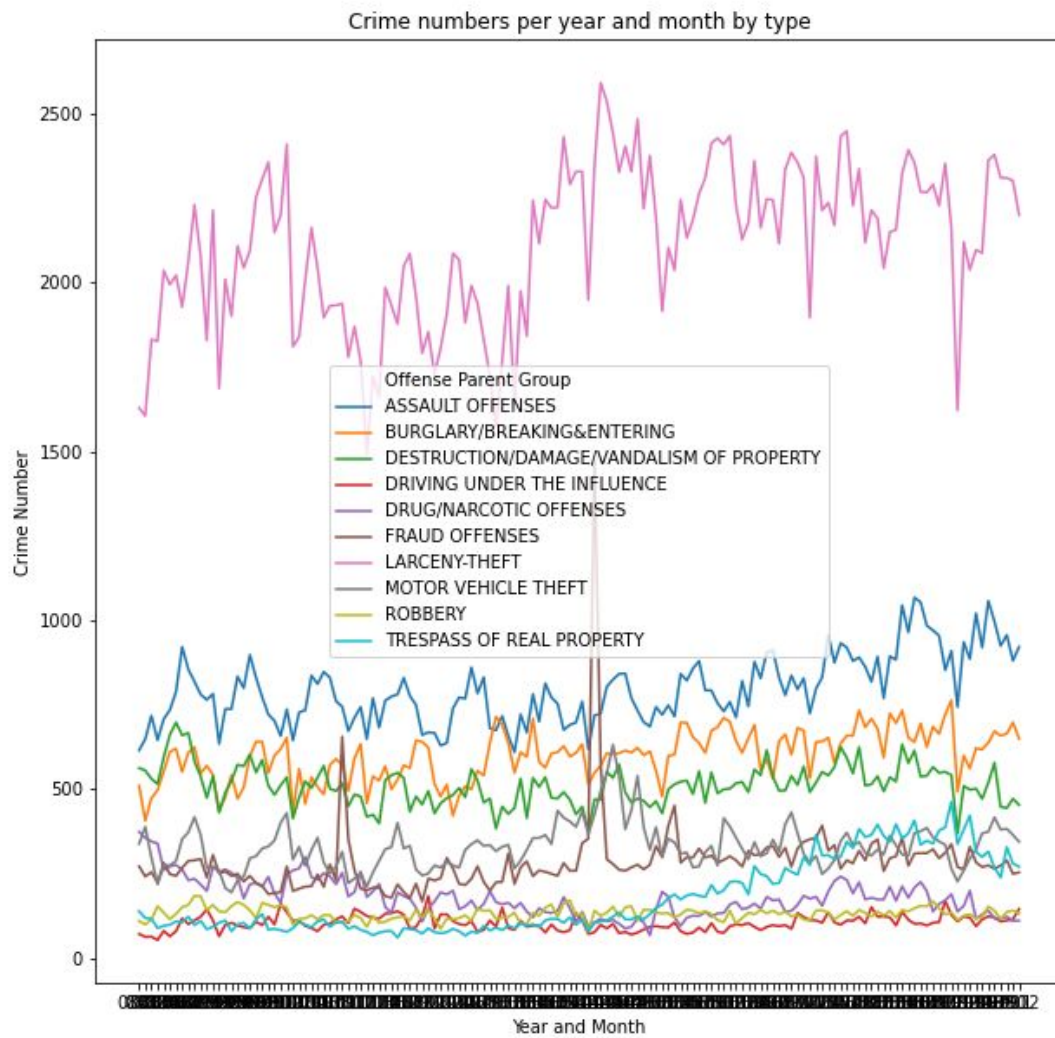
Research Question 2 ([Click to go back](#))

1b11.png

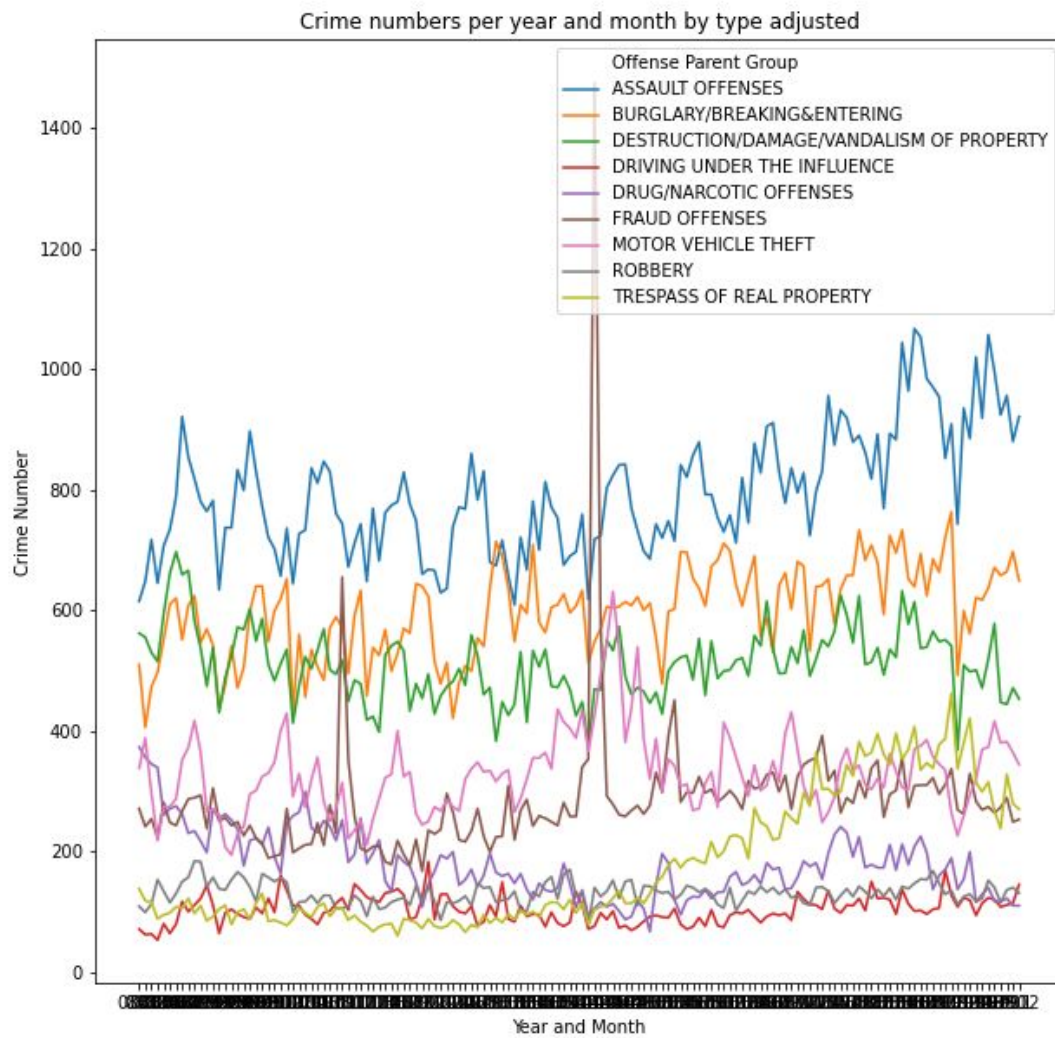


1b12.png

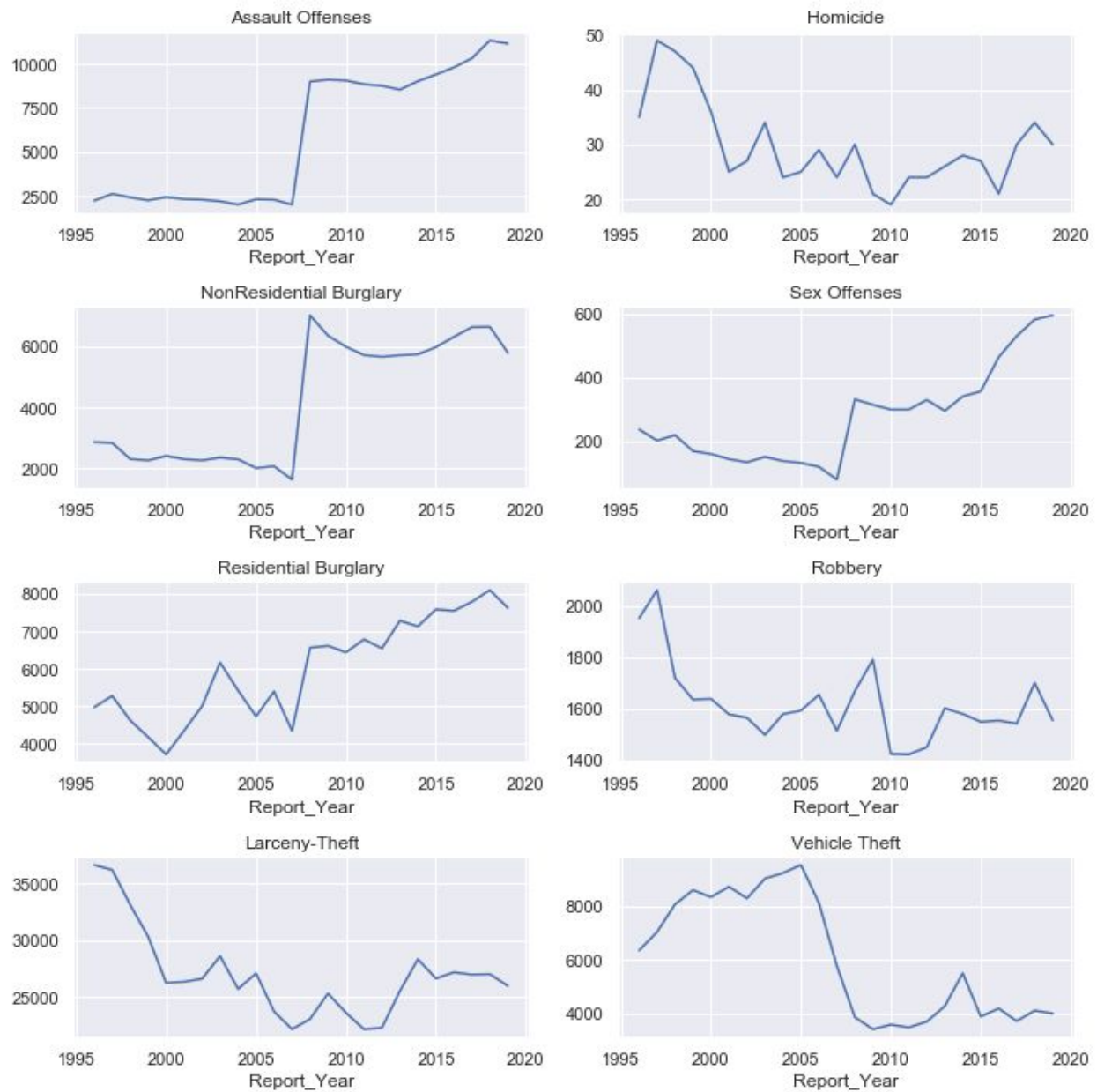


1b21.png

1b22.png

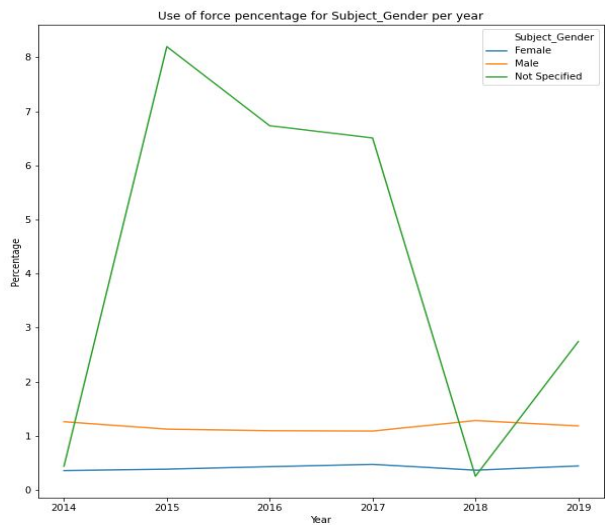


combined_crime_report.png

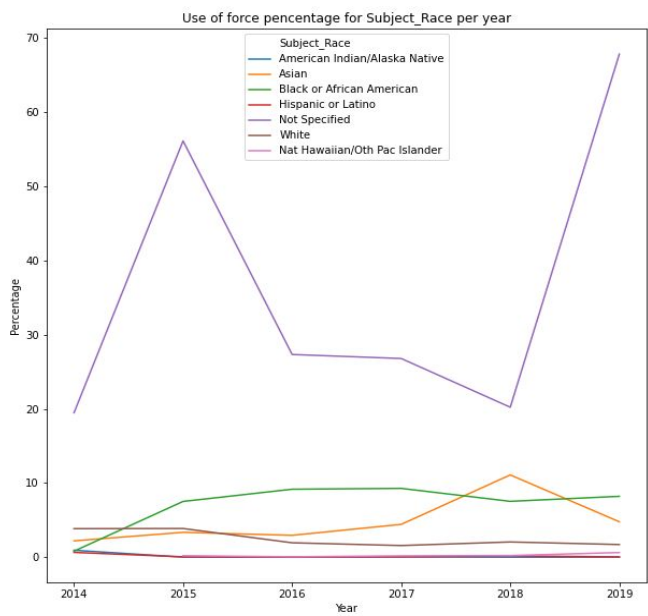


Research Question 3 ([Click to go back](#))

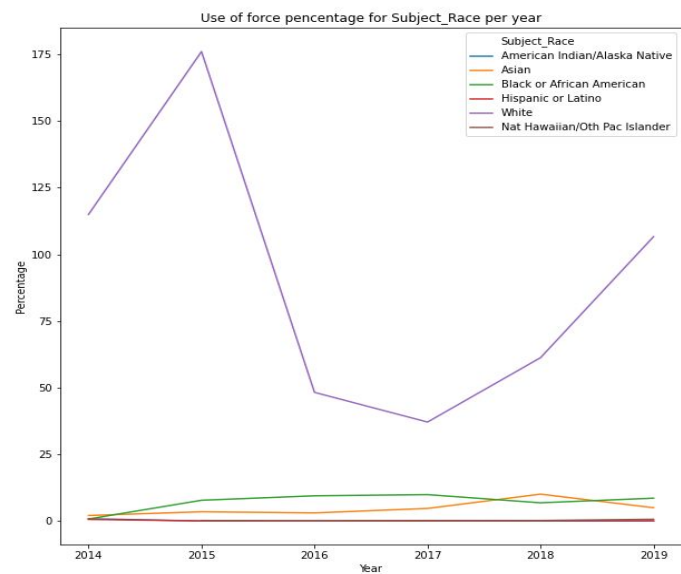
Subject_Gender.png



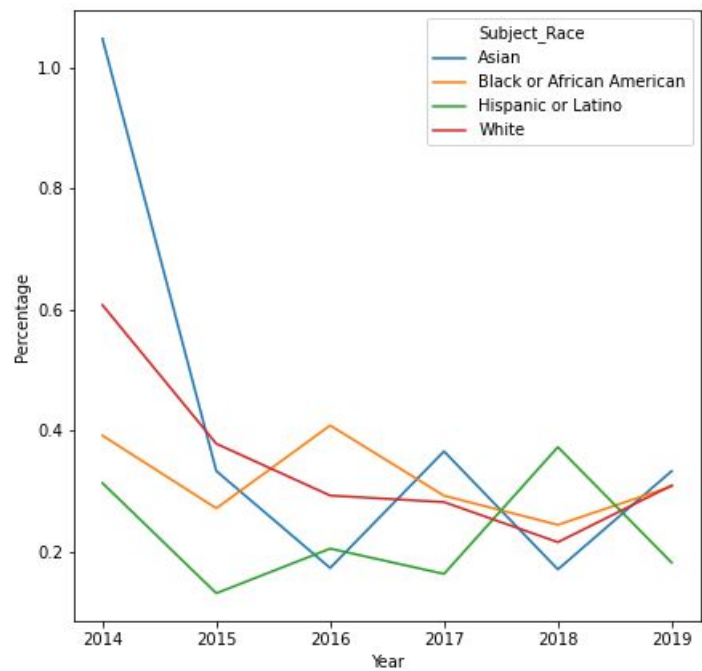
Subject_Race.png



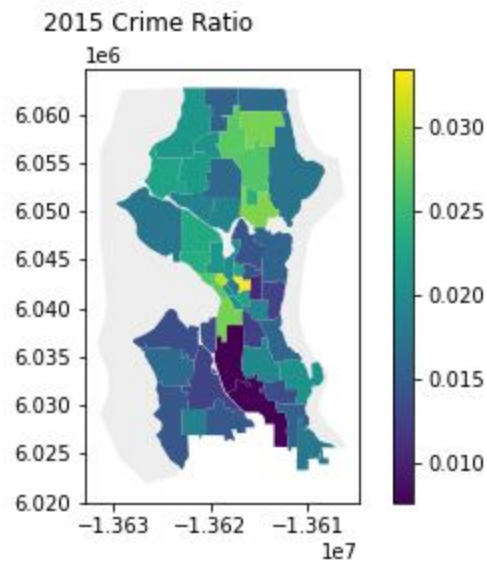
Subject_Race_New.png



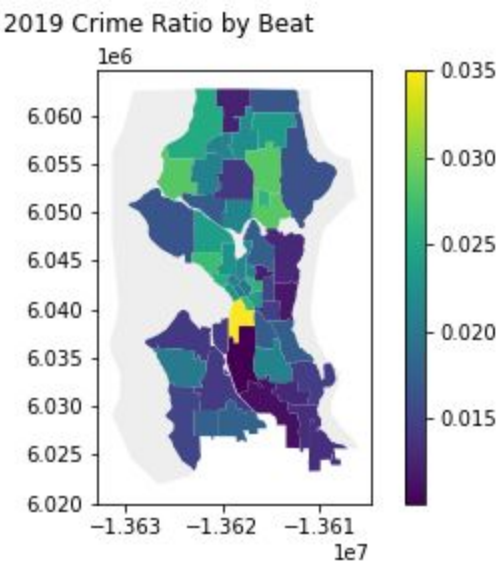
3a.png



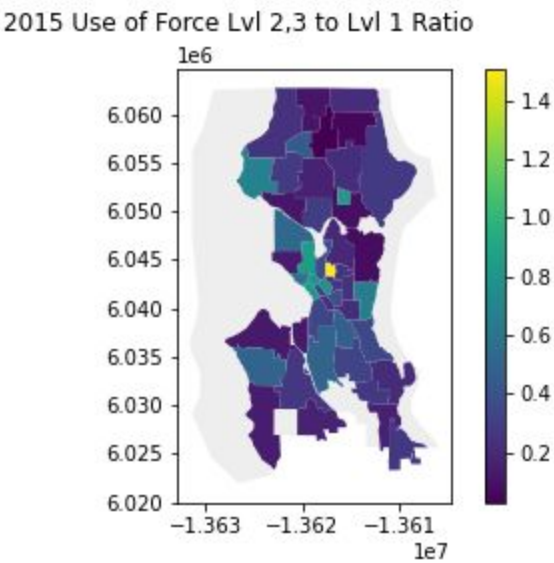
3b2015_0.png



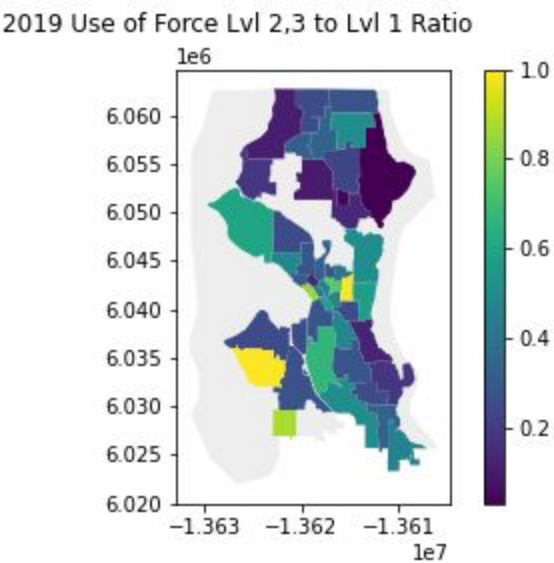
3b2019_0.png



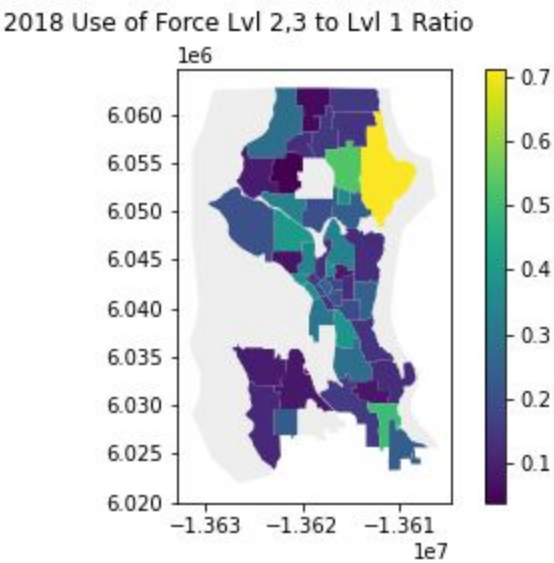
3b2015_1.png



3b2019_1.png



3b2018_1.png



3b.png

