

Exploring NYC Complaints

Dav King

2022-04-22

Note: this version of the document is a clean version, with the code removed and some of its lengthy outputs muted. The full, unmessy version is also uploaded to the Sakai dropbox.

1. Briefly Describe the Dataset

This dataset is an aggregate of cases raised by constituents across every NYC council district office, with complaints about a wide variety of topics. Though different offices gather their data differently, it has been aggregated into a fairly specific set of overarching labels. There are eleven variables in the dataset: a unique identifier for each case (`unique_key`), a listing of which city council district each complaint was raised within (`account`), the date that the case was opened (`opendate`) and the date that it was closed (`closedate`), the type of complaint being raised (`complaint_type`) and a more specific indicator of what was being targeted in the complaint (`descriptor`), and the zip code, borough, city, council district (`council_dist`), and community board (`community_board`) in which the complaint was raised.

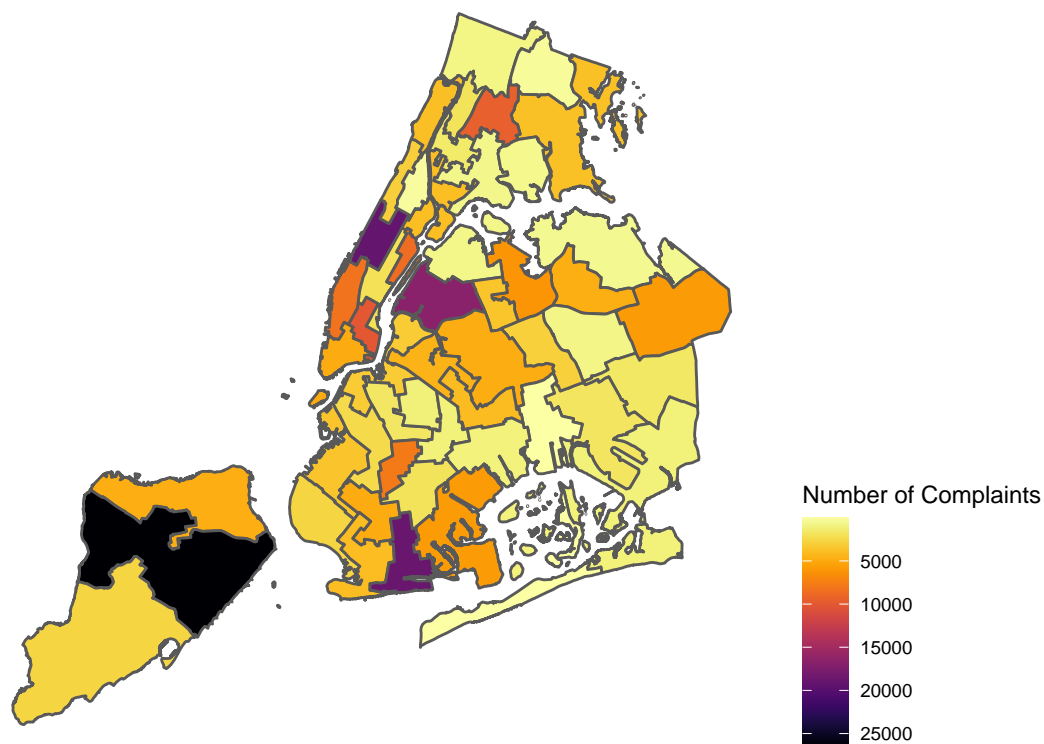
2. Look at each column/variable in more detail

Using the `describe` function from the `Hmisc` package, we get a very extensive description of every variable: the number of observations it has, the number of missing (NA) values that it contains, and the number of distinct observations that it contains. It also shows the lowest and highest values for each variable (which are alphabetical, in character strings), and for numeric variables it also includes several quantile denotations. This is a great way to get a view of a complete dataset with minimal work.

unique_key: there's not much to say here. This variable is an indicator of each separate case being raised to the NYC city council district offices. It starts with the code of the specific `account` with which it was raised, and then appears to be a random, mostly sequential string after that. It is present for every variable and almost entirely unique across all cases, `n = 240027` and `n_distinct = 239958`.

account: this is a descriptor of which of NYC's 51 different districts the complaint was raised within. It begins with the string NYCC followed by a two-digit number from 01 to 51, enumerating the city district. It is present for all variables. In the graph below, I count and visualize the number of complaints raised in each different council district in order to see whether there are any geographic trends to speak of. It does not appear that there is much of note.

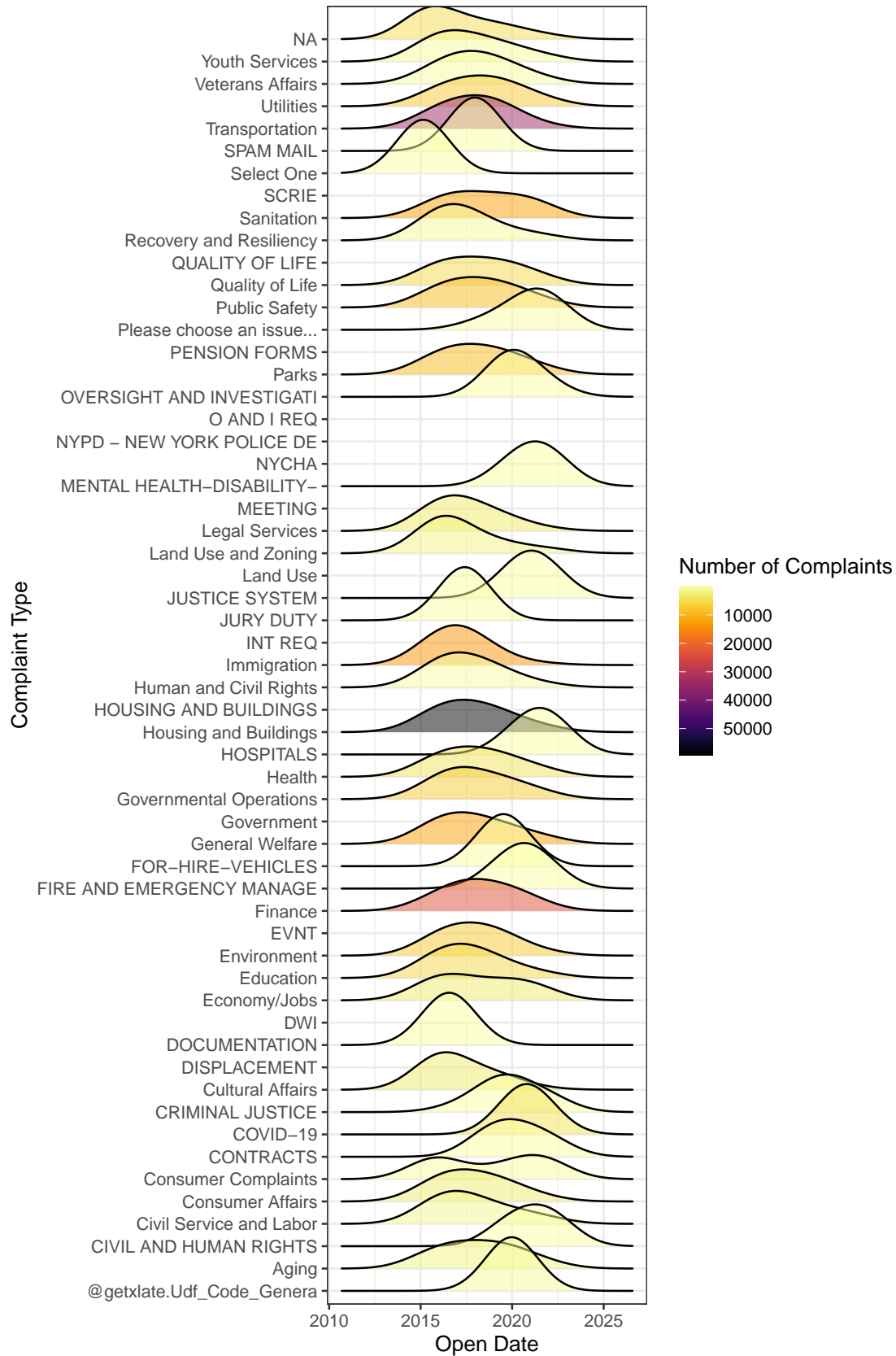
Complaints by Council District



opendate: this variable lists the date that each complaint was opened with an NYC council district office. The first complaint in this dataset was opened on January 1st, 2015, while the most recent complaint was opened on April 7th, 2022. The total number of distinct opening dates (2596) suggests that a complaint was opened almost but not quite every day between those two dates. The mean date for opening a complaint was December 29, 2017, but the median date for opening a complaint was October 2nd, 2017, suggesting a slight left-skew in opening dates for complaints (that is, that people have been opening complaints more frequently recently than they used to).

The graph below shows the distribution of open dates for different complaint types, allowing us to identify any trends in the data structure and any possible outliers. By shading the data according to the total number of complaints in each complaint type, we can identify that the most common complaints had much more broadly distributed peaks in terms of when they were most frequently filed, with medians around the median date of the whole dataset. Though the ridgeline plot smooths out a lot of the detail from this dataset, it still shows that each different complaint type saw a peak at some point in time, with some most common as early as 2015 and others peaking as late as 2021. There is some rhyme and reason to this distribution: for example, complaints about Covid-19, Hospitals, and the Justice System all peaked in the past couple of years, when those issues were notably more salient than they had been in recent times before that.

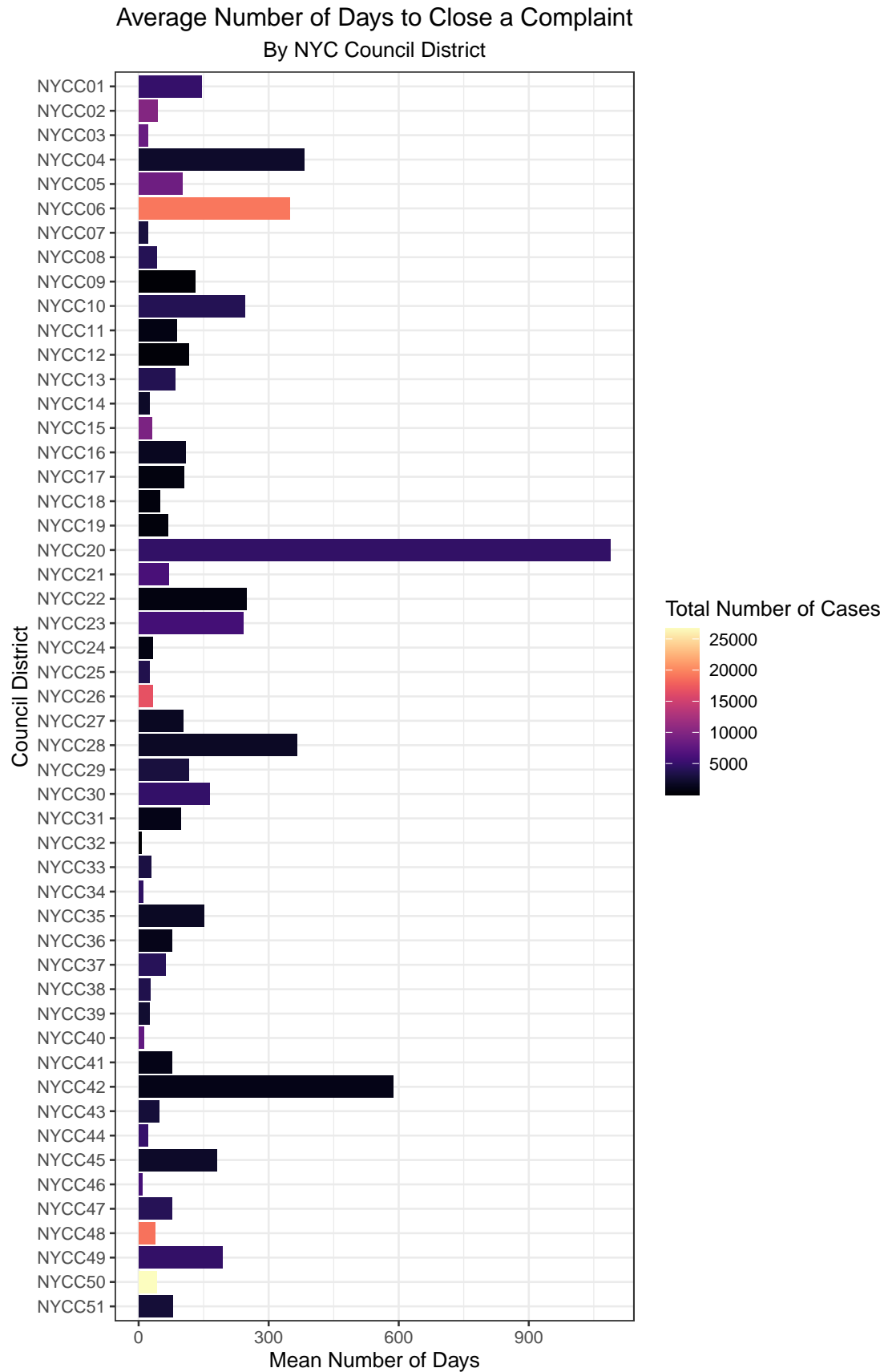
Distribution: Open Date of Different Complaint Types



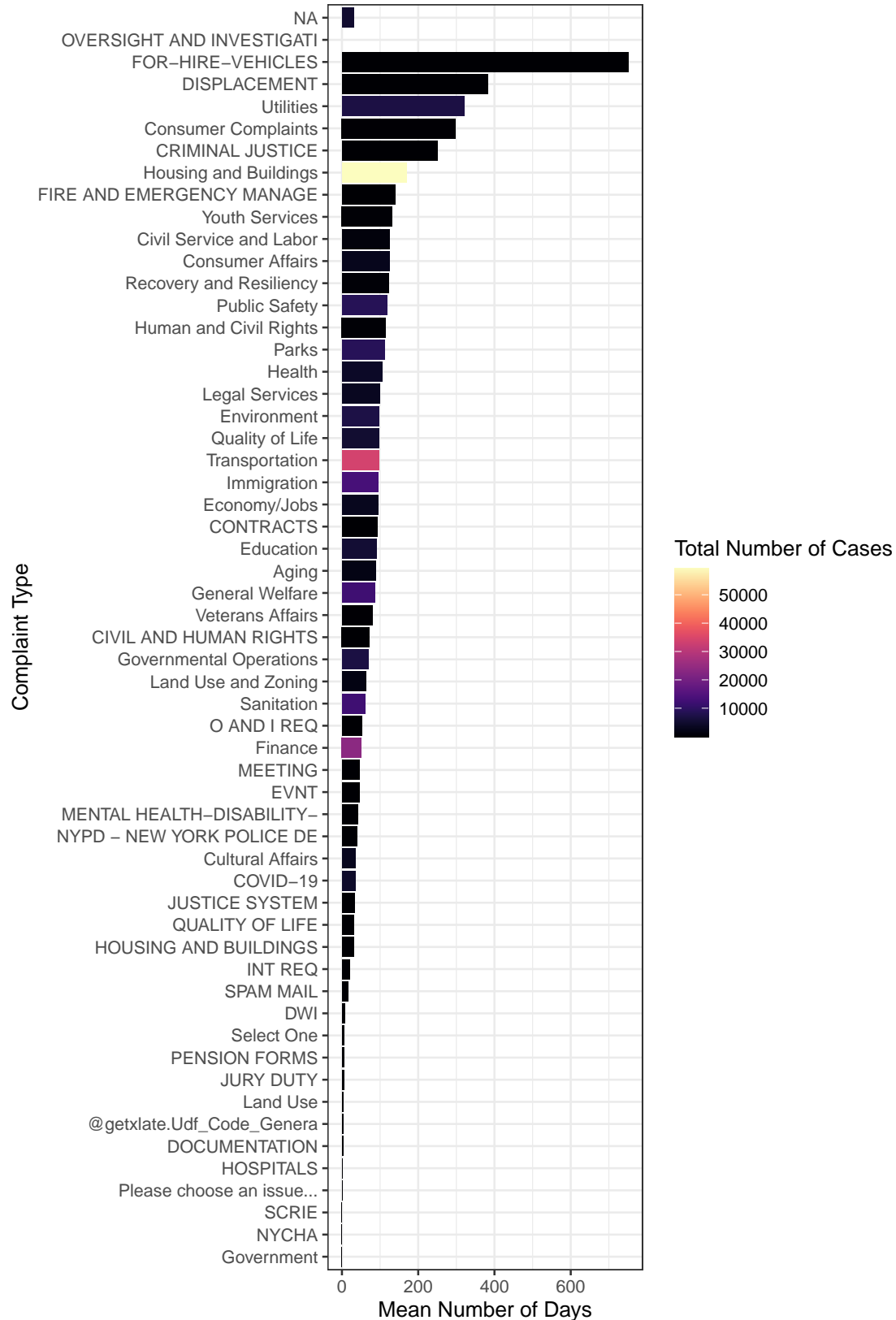
closedate: this variable lists the date that each complaint was closed. Though every complaint had an open date (meaning all of these close dates are for complaints that were opened on or after January 1st, 2015), not all complaints have a close date - it's missing for 16046 of these complaints. Similarly, there are far fewer distinct dates in here, suggesting that on some 500 days between 1/1/2015 and 4/7/2022 there were zero complaints closed citywide. The first date that any complaint was closed was January 2nd, 2015 (all of which happened to be opened on that same date), while the last date that any complaint was closed was April 7th, 2022 - the last date of the dataset. Despite how quickly some complaints were closed on January 2nd, the mean close date was March 23rd, 2018 (a little less than 3 months after the mean open date) and the median close date was January 29th, 2018 (almost four months after the median open date). Not only does this suggest that most complaints take a while to process, it also suggests a right skew in how long they take to close - that is, that most cases are closed relatively quickly, but some take a very long time to be closed.

In the graphs below, we can consider the average number of days to close a complaint grouped by both council district and complaint type, looking for trends and outliers. There is a lot of variability in how long it took council districts to close complaints on average, with some as short as perhaps a couple of weeks and others (very notably the 20th district) taking over 1,000 days to do so on average. This appears to be largely unrelated to the total number of cases that each district office received, though the districts receiving the most cases did tend to resolve them in the shortest amounts of time (probably both because they were receiving much more generic, easier complaints and because they were more prepared to handle them).

When grouping by complaint type, the distribution made much more sense - the complaints that took the longest to resolve on average were *much*, much longer than most other types of complaints, and were very low in terms of raw case counts. There is one notable exception however - the sixth-longest complaint type to resolve on average was also the one that had the outright largest number of cases - housing and buildings. That seems logical, as many people file complaints about their housing but landlord protections are incredibly strong compared to tenant protections. This graph also reveals that no **OVERSIGHT AND INVESTIGATI** complaints were ever resolved, while all **NYCHA**, **SCRIE**, or **Government** complaints were either never resolved or resolved the day they were filed.



Average Number of Days to Close a Complaint
By Complaint Type



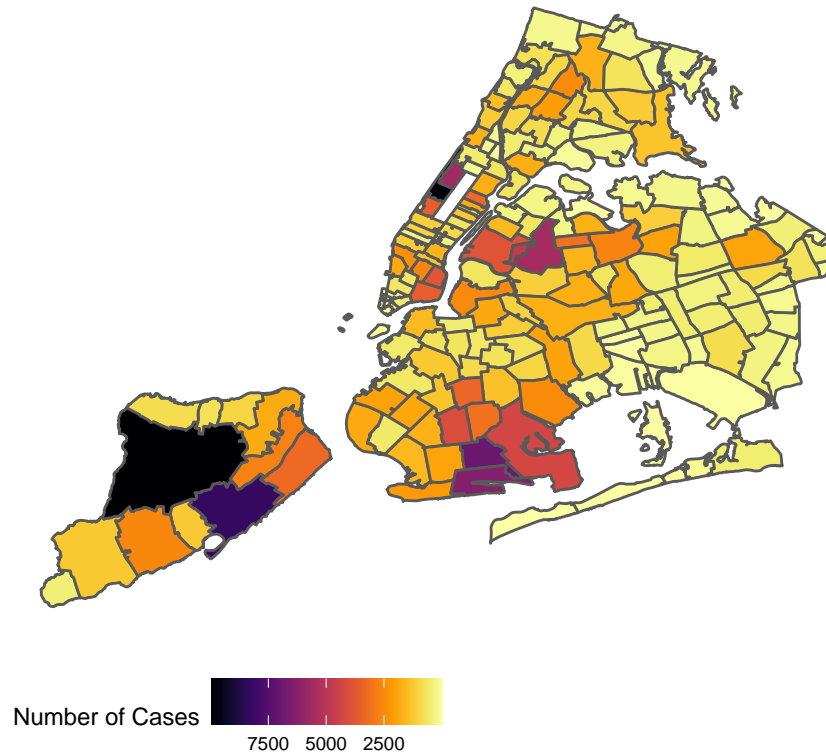
`complaint_type` is the first truly messy variable in this dataset. The code below briefly allows us to inspect it. On initial inspection it contains 56 different types of complaints, but when reducing it to lowercase that number drops to 55. While most of these variables make sense, some are problematic. 5060 are `NA` values, which drops the number of types of complaints to 54. 607 are labeled `@getxlate.udf_code_genera`, which is clearly a tabulation error - as are the 143 complaints labeled `please choose an issue` and the 5 labeled `select one`. Additionally, some of the variables look like they should arguably be merged together, and many of the complaints at very low frequencies are riddled with misspellings or are oddly specific enough that they seem functionally useless within this framework.

`descriptor` is a series of more specific explanations for what the complaint being filed was about. With just over 1,000 distinct responses, it's impossible to try to view and clean all of them - but it's very clear that a lot of these data points are serious issues, as some are just numbers and others simply say `--Select--`. Additionally, while not a large proportion in the grand scheme of things, over 13,000 of these complaints are missing a specific descriptor from this column. The code below allows us to view the most common descriptors (anything seen over 1,000 times in the dataset) along with the complaint type that they were housed under in order to identify trends in the meaningful chunk of the data. Looking at the list, we can see that most of the common descriptors came in the `Housing and Buildings` complaint type, with `Finance` and `Transportation` also rather popular complaint types. As the number goes down, the complaint types diversity (which is to be expected).

`zip` gives the zipcode of the address where the person filing the complaint resided. Fortunately, every observation has a zipcode attached to it, and the total number of distinct zipcodes (1046) is far less than the total number of zipcodes found within New York City. Unfortunately, not all of these zipcodes are within NYC. Though the lowest zipcodes in this dataset are 00000, 00001 and 00002, the lowest zipcode in the country is 00001 and it is found in rural Alaska (as is 00002). However, we can match this with the zip codes found in the `nycdogs` package, and by using a convenient `inner_join` determine exactly which zip codes are actually within NYC.

The map below counts the number of cases opened in each zip code and displays them in a map. This shows the relatively low number of cases opened in Queens, especially as compared to areas in South Brooklyn or Manhattan. Overall, there is not a super large trend in geographic distribution visible on this map.

Cases Opened by Zip Code



`borough` gives the borough in which the person filing the complaint resided. Though some 16 thousand observations are missing this variable, all of the rest do fall into the five boroughs of NYC. Brooklyn residents filed the highest proportion of complaints (28.5%), followed closely by Manhattan residents (27.2%), then Queens residents (21.3%) and Staten Island residents (14.8%) with Bronx residents bringing up the rear at a measly 8.3%. To some extent, this likely corresponds to the relative population of each borough.

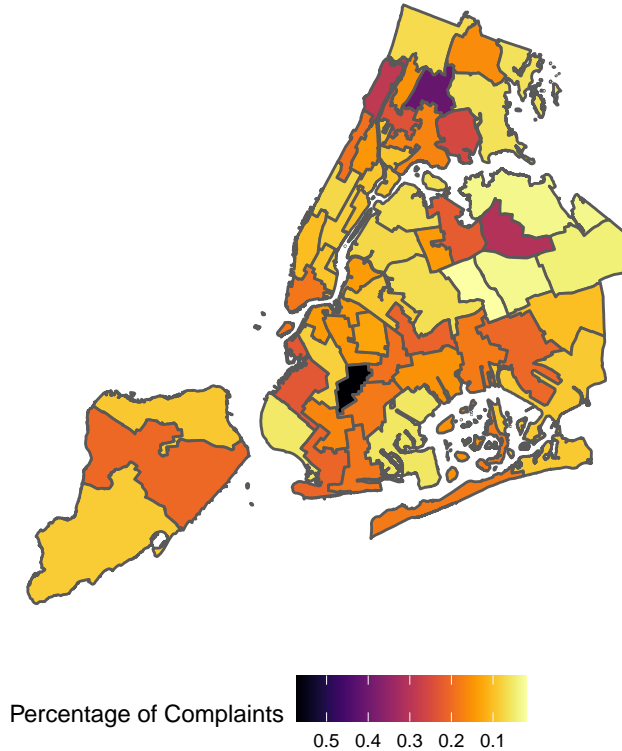
`city` corresponds to the city in which the person filing the complaint resided. This data is missing for far fewer observations (just under a thousand), and has nearly a thousand distinct labels. Once again, some of these are not real cities (e.g., 0, 1 Montague Street, Your City), and others are duplicates (Yonkers, New York and Yonkers, Ny). These might be harder to verify than the zip codes, but they are still a valuable source of data to the extent that they can be trusted.

`council_dist` corresponds to the council district in which the person filing the complaint resided, which is not necessarily the same as the `account` variable corresponding to the council district in which the complaint was opened. This is a messy variable, as council district numbers have anywhere between 0 and 2 zeroes preceding them and thus there are serious overlap issues. NYC has 51 districts, but there are 154 distinct values within this variable. This variable can be cleaned for reasonable use without too much effort, though a substantial number of cases are lost when doing so (as some numbered higher than 51). There are nearly 10,000 missing observations, which is a substantial amount.

The code below mutates the `council_dist` variable into a usable format, identifies the proportion of cases in `account` that do not match up with `council_dist`, and displays this output in a map in order to identify which districts had a lot of complaints filed by out-of-district residents. Looking at the map below, we can see that a surprising number of complaints filed in the Bronx were filed by people who resided in a different council district. Far and away the district with the most complaints originating from residents of a different district, however, was District 40 - with 57.2% of its complaints being filed by residents of a different district (compared to the next highest, District 15, at 41.2%, and followed by district 20 at 30.7%). I can honestly find no reason for this to occur. District 40 is a majority-black district in Brooklyn, but there is only one

thing of note about it - it is close to the geographical center of Brooklyn. District 15 is similarly located at the geographic center of the Bronx, but it makes more sense, as upstart Congressman Ritchie Torres' former district and the home of both the Bronx Zoo and the New York Botanical Garden.

Complaints Filed by Resident of Different District



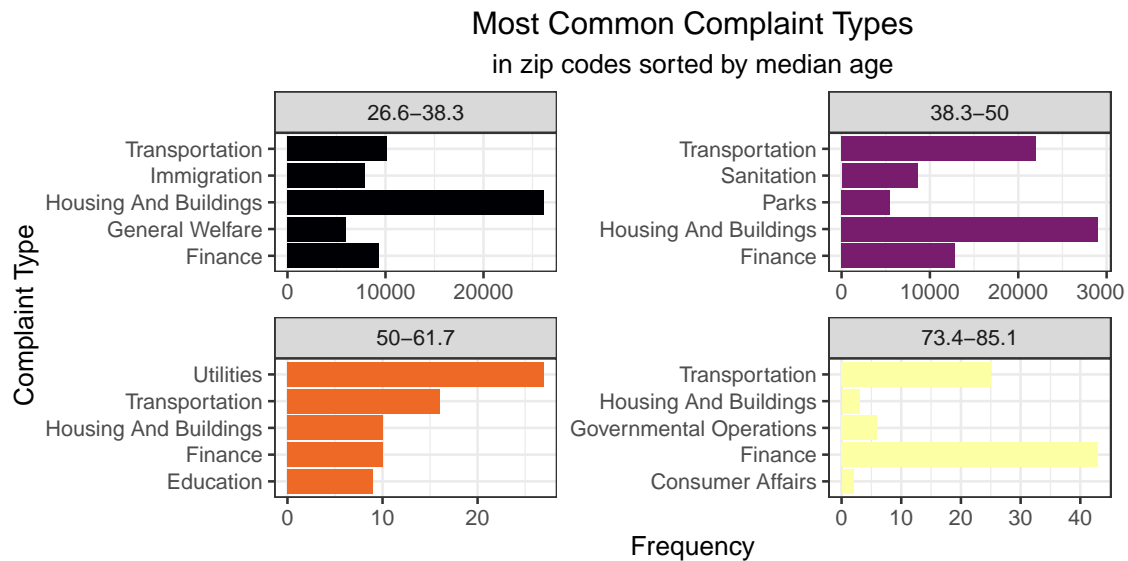
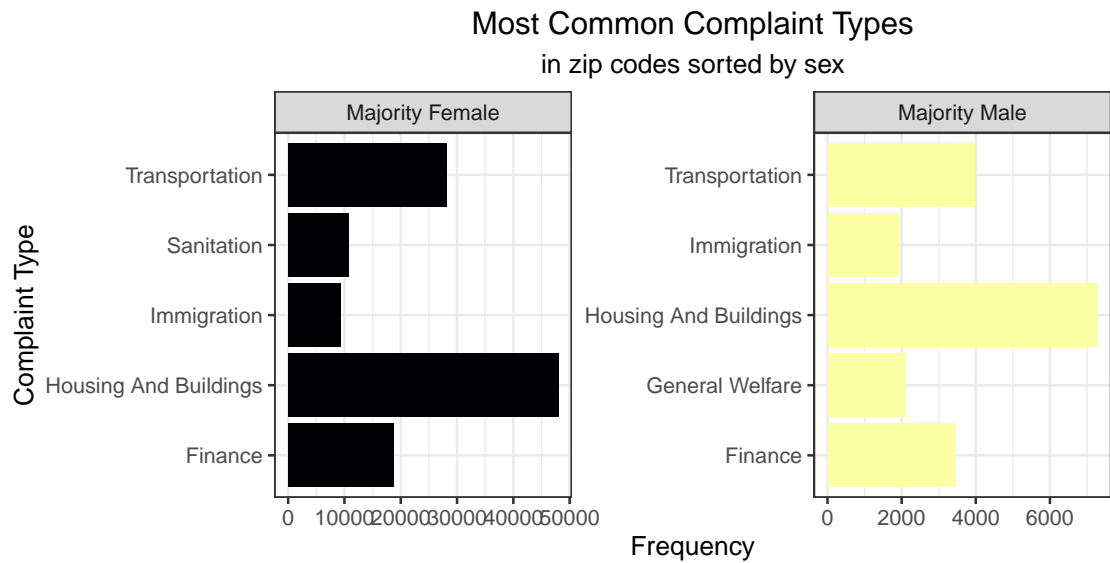
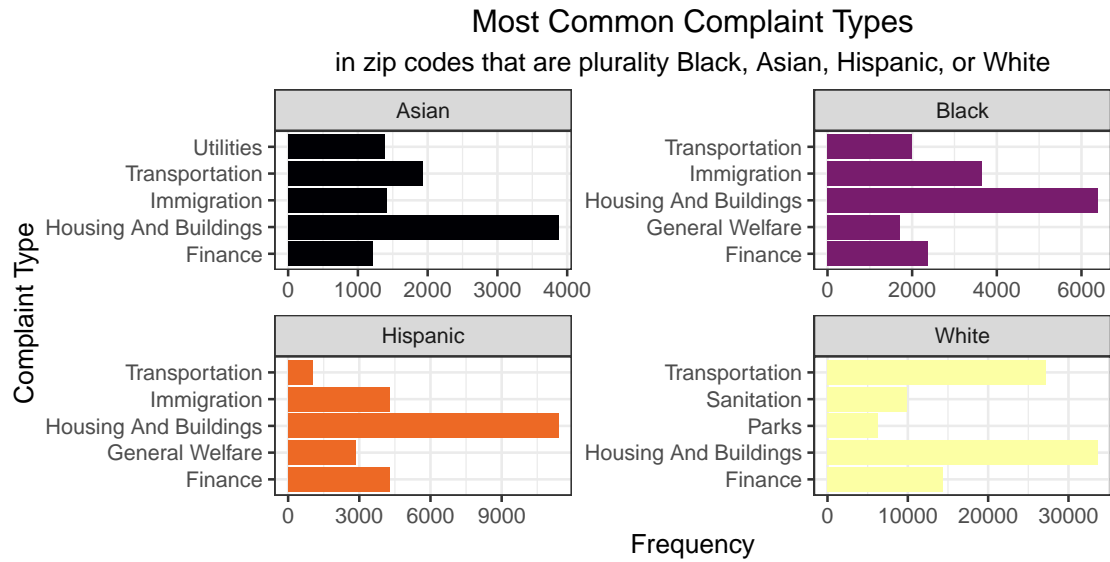
Finally, `community_board` designates the specific community board of the constituent. Though missing for just over 16,000 observations, there are only 94 distinct variables - thus, we can view all of them. Most of them appear to correspond to real community boards, but some 30 or so of these boards listed are just numbers or otherwise meaningless information (primarily corresponding to single-digit numbers of responses).

3. Explore in more depth: at least three polished plots or maps

- Show your work.
- Provide a motivation for and discussion of each plot.
- You can merge in external data if you like.

Plot 1: Complaints by Demographics

I was curious to see whether there were differences in complaints filed by the demographics of the filer. However, this information is not included in the dataset. While there is no way to fill in that information, there is a somewhat crude way to do it - by identifying the most common complaint types filed by residents of zip codes that are constituted of a majority of people belonging to various demographic groups. Drawing from the 2020 five-year ACS, the variables I chose to select were race, median age, and sex. Note that while no math is done to put the complaints in per-100,000 format, setting `scales = "free"` enables us to compare all of these data points on equal dimensions - essentially serving the same purpose.



Across all races, housing and buildings was still far and away the most common complaint type. Finance and transportation were also seen in the top five categories across all race groups. Perhaps unsurprisingly, one of the biggest sources of complaints was immigration issues for zip codes that were predominately every race *except* for White - particularly for Black and Hispanic plurality zip codes. Parks was a complaint seen most frequently in plurality White zip codes, likely because parks are often attached to gentrification in areas such as NYC. General Welfare, which includes many descriptors often related to poverty (such as food stamps, homeless encampments, and social security disability) was in the top 5 for Black and Hispanic plurality zip codes, while Utilities was common in Asian-plurality zip codes and Sanitation was in White plurality zip codes. While this obviously is still a crude measure of race in complaint types, a lot of these findings are rather sensical. This graph also gives some vague insight into the relative proportions of NYC residents who are Asian, Black, Hispanic, and White, though this is where putting the data in per-100K format would have served a useful purpose.

The graph dividing zip codes according to sex is probably functionally useless, as most zip codes are slightly majority female (reflecting the US sex divide as a whole) but none are greater than 55% in either direction, suggesting that this data breakdown is a poor proxy for the complaints filed by members of each sex. Still, there is one difference of note - one of the top five issues in female-dominated zip codes was Sanitation, while this was replaced by General Welfare in male-dominated zip codes. Other than that, these data show almost identical trends for both groups.

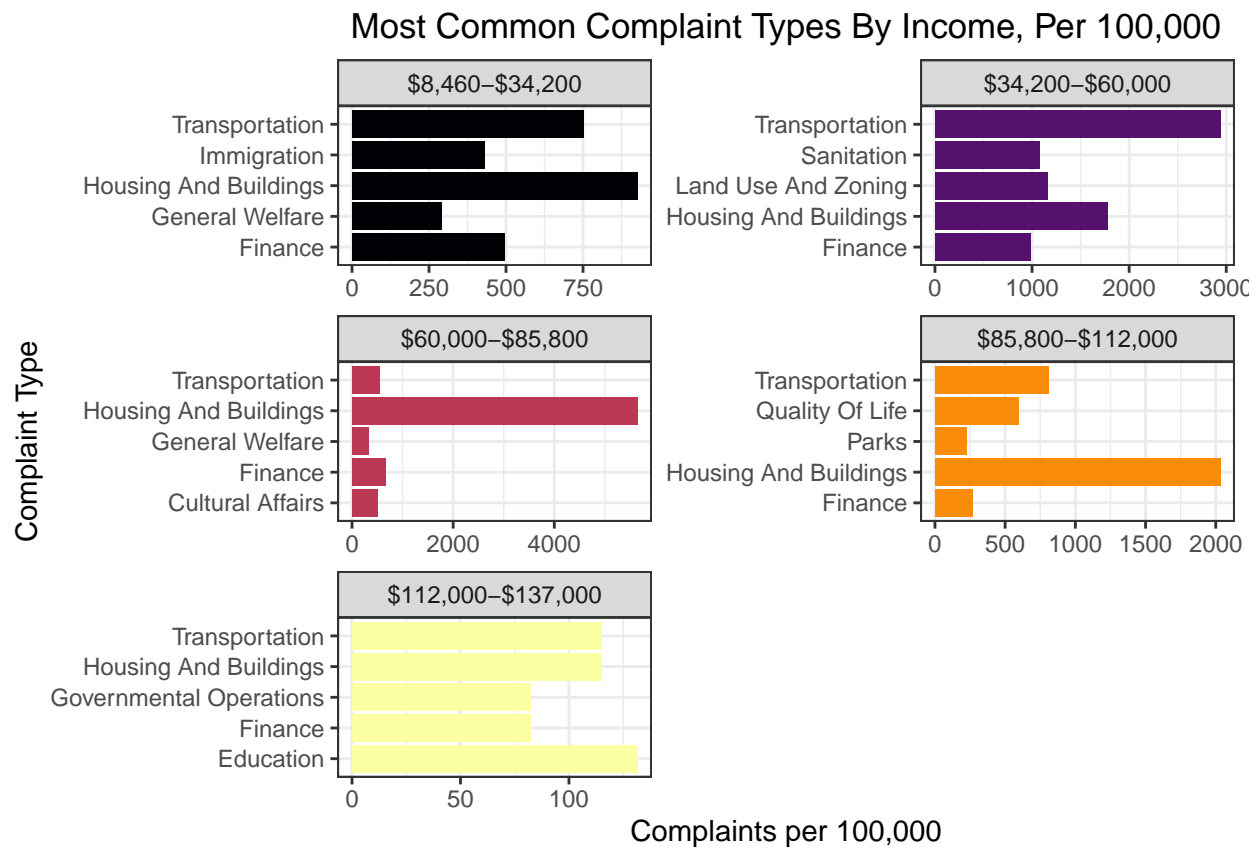
The graph dividing zip codes into different categories of median age is striking. To be clear, most zip codes have low median ages - though the data were split into five different intervals, zero zip codes fell into the 61.7-73.4 range, and almost every zip code had a median age below 50. However, there were seriously different trends in the different age groups. The first two groups again saw housing and buildings, transportation, and finance receive the largest number of complaints. However, zip codes with a median age of 26.6-38.3 saw immigration and general welfare round out their top five most common complaint types, while zip codes with a median age of 38.3-50 instead saw sanitation and parks fill that role. This is remarkably similar to the Black/Hispanic vs White split in the first of these graphs, and suggests that the prevalence of those specific complaint types may well be dictated largely by income. Meanwhile, zip codes with a median age from 50-61.7 were quite different - utilities were far and away the most common complaint type, followed by the usual transportation, housing and buildings (but not the most common this time!), and finance, with education rounding out the rear. Meanwhile, for zip codes with a median age of 73.4-85.1 (so, essentially, retirement homes), complaints about finance were far and away the most common (as people who are retired often have such issues), with a reasonable number of complaints about transportation (again, logical, as many old people slowly lose their ability and right to drive) and almost no other complaint types being filed regularly (with the most common remaining group being governmental operations, again related to retirement benefits).

While none of these plots can be said to truly display the actual trends by these demographics, as the proxy making up entire zip codes blunts a lot of the data, many of these trends indeed make a lot of sense. In the next plot, we will look more closely at the impact of income that was hypothesized as a result of some of these plots.

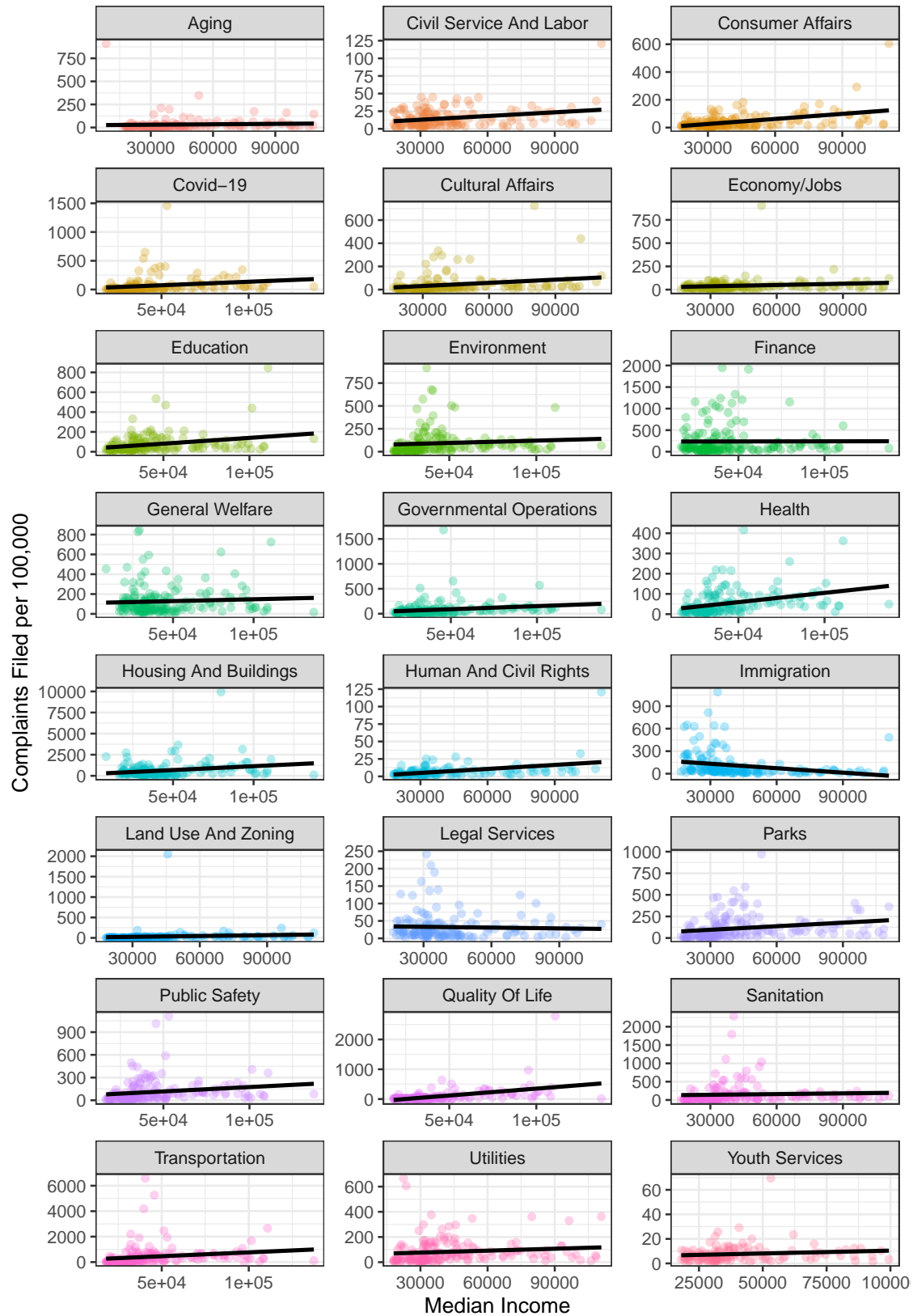
Plot 2: Complaints by Income

In order to consider how income affects these different complaint types, I drew two different plots (call it one and humor me, they're the same general idea). In the first, I split the data into five equal intervals of median income, selected the five most common complaint types per 100,000 in each of them, and displayed the output in a bar chart. This enables us to see what the most common complaint types are among people of different social classes in the same way as the previous demographic plots - not on a case-by-case basis, but still on an aggregate level. In the second, I calculated the per-100,000 complaint rate of each type of complaint according to zip code, filtered out any complaint type that showed up 100 or fewer times, and plotted it against median income. Though most of these trends are clearly not linear, the linear regression

lines can still reveal something about the nature of these relationships - namely, whether richer zip codes tend to file each complaint type more or less frequently than poorer zip codes do.



Relationship between Median Income and Complaints Filed In NYC Zip Codes



In the first graph, a couple of notable trends show up. First of all, even on a per-100K basis, almost no complaints are filed in zip codes falling under the richest bracket - probably logical, when considering the rate at which rich people actually use government services (instead of just paying for the problems themselves). Rates are also somewhat lower in the lowest bracket of zip codes. As was somewhat predicted by the takeaways from the previous set of plots, immigration was only an issue in the lowest bracket of incomes, while general welfare was only in that group and the \$60K-\$85.8K group. Like usual, housing and buildings, transportation, and finance were in the top five issues across all categories. There are some other notable trends in this plot, such as education and governmental operations only being issues in the richest zip codes or quality of life and parks only being an issue for the next-richest set.

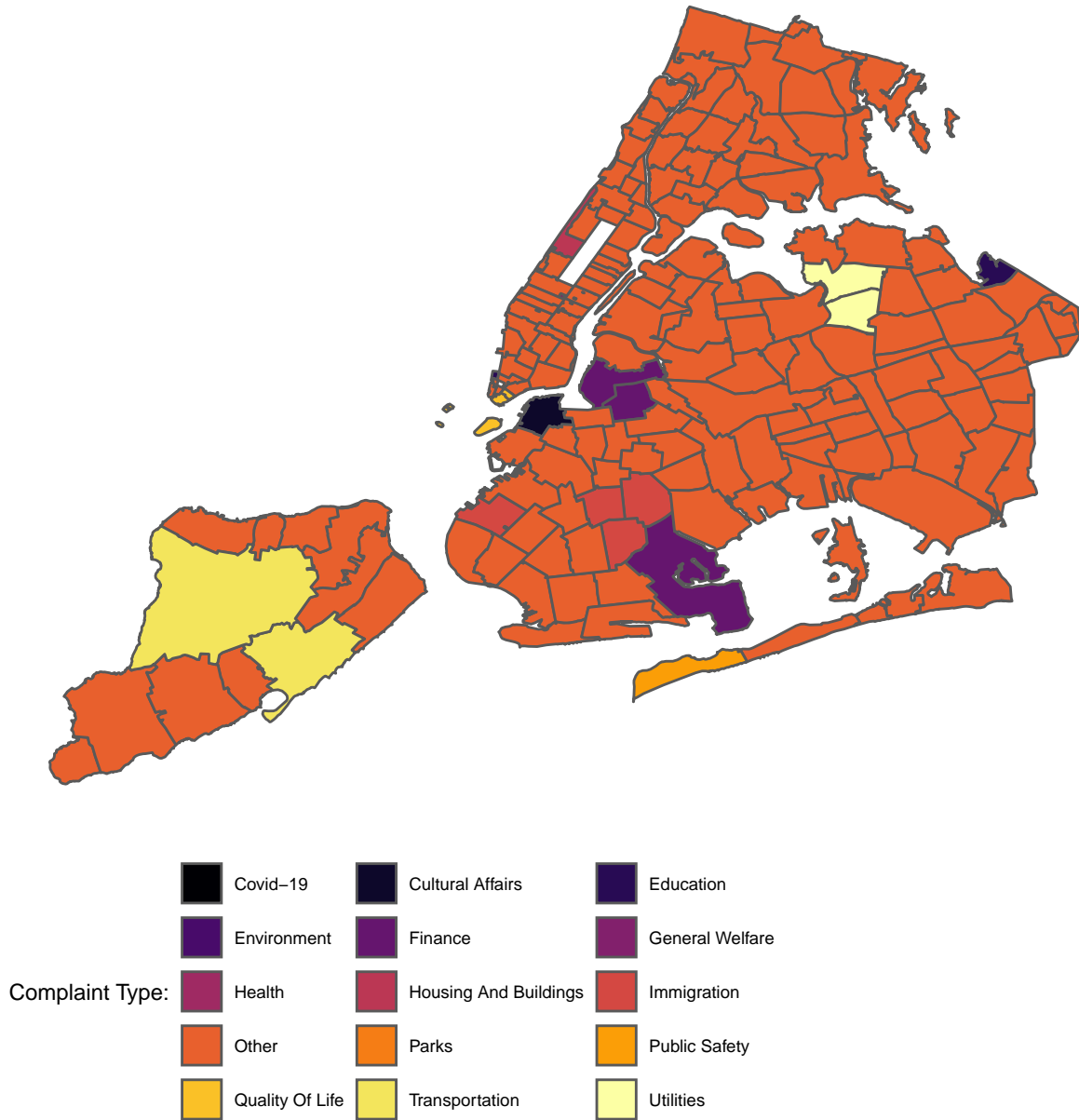
The second graph suggests that a) almost all relationships between income and per-100K complaints are incredibly weak and b) that almost all of those relationships are positive - that is, that people at higher incomes file complaints at slightly higher rates than do people at lower incomes (once again, through the proxy of aggregate zip code median incomes). Notably, there are two complaint types for which there is a negative relationship - immigration and legal services. That is, similarly, a logical trend, given the nature of these complaint types in NYC.

Plot 3: Distribution of Complaints

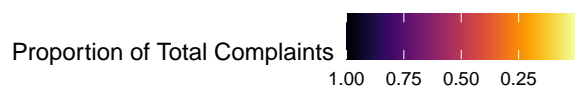
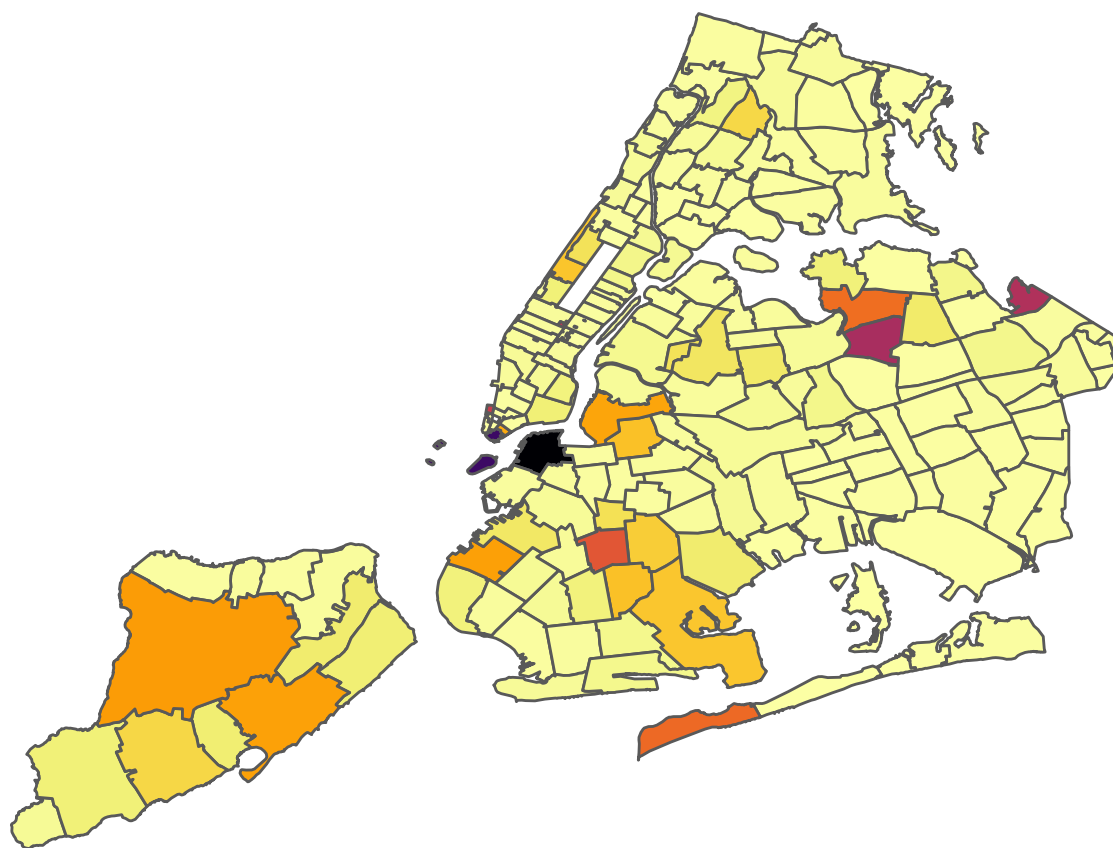
In the last set of plots, I wanted to look at the geographic distribution of these different complaint types. Thus, using the geography data from the `nycdogs nyczip`s package, I drew a series of four maps of the city. The first is a map of the most common complaint types filed in each NYC zip code, marked as “Other” when no complaint type was over 10% of the total complaints in that zip code (which was the case for most of them). This enables us to see whether certain complaints are more common in different geographic areas of the city. The second map shows what percentage of the total number of complaints filed in each zip code were the most common complaint type, an addendum to the first map (originally, these were supposed to be one combined map, but there was no interpretable way to include this much information in the one plot). The third is simply a map of median incomes, allowing us to see geographic trends, and the fourth is an image of how many complaints per 100K were filed in each zip code.

Most Common Complaint Types

By NYC Zip Code

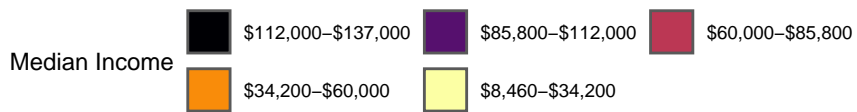
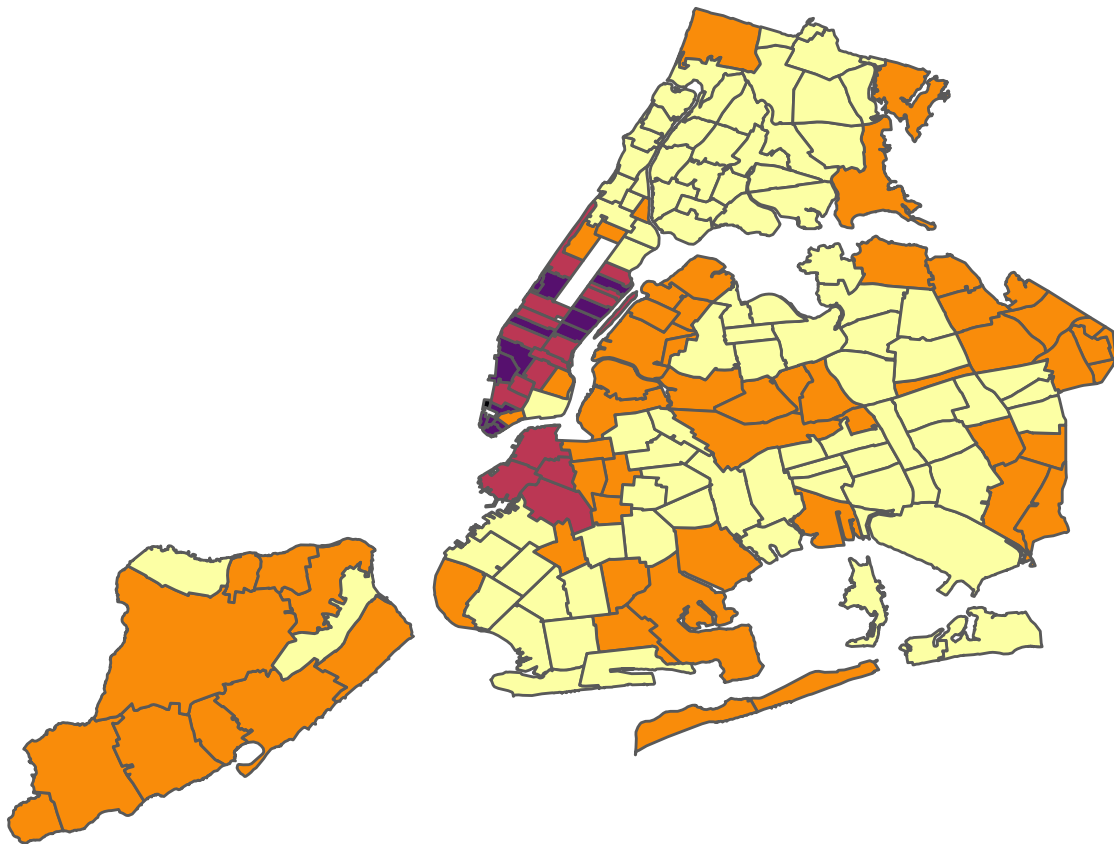


Most Common Complaint Types – Proportion of Total Complaints By NYC Zip Code

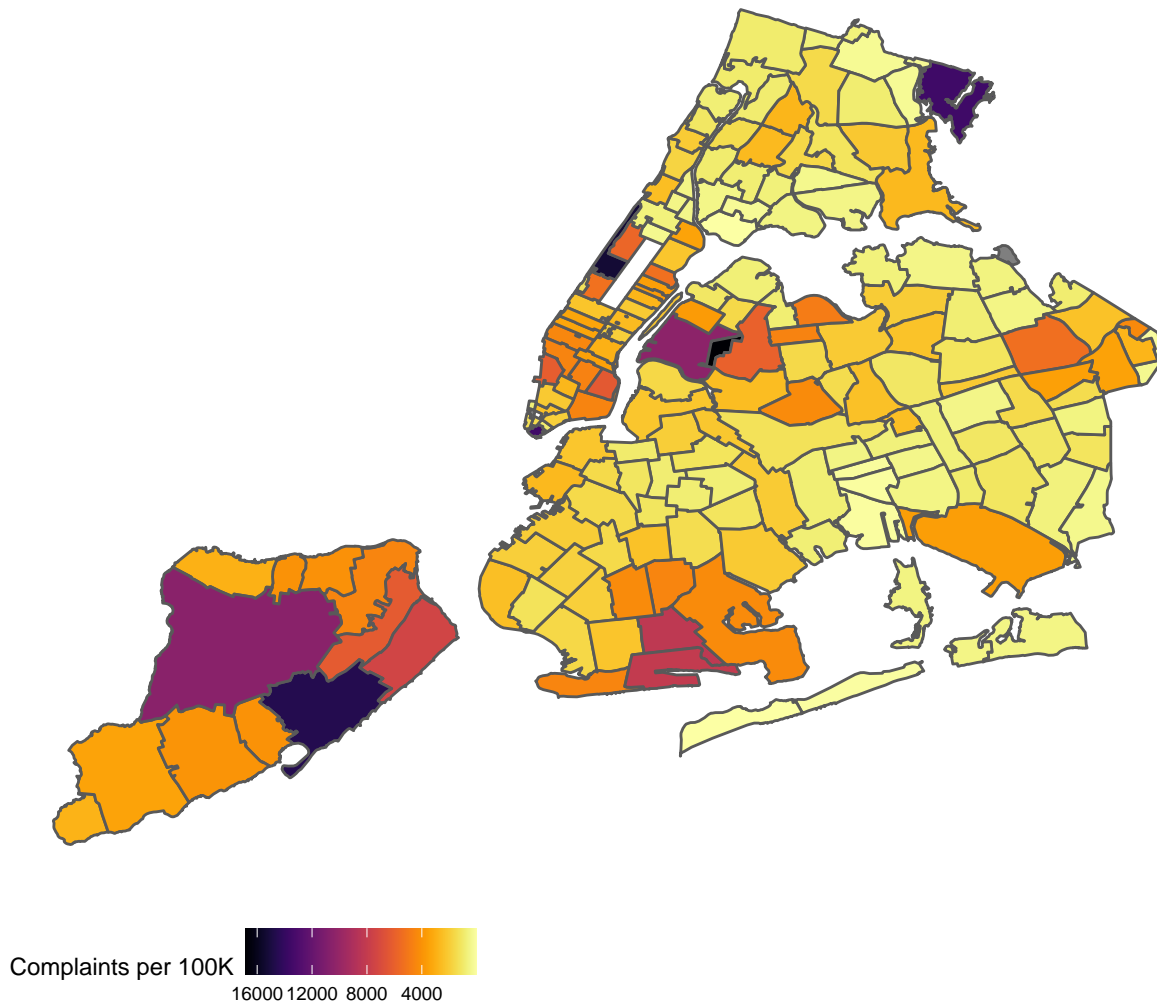


Median Inome

By NYC Zip Code



Complaints Filed per 100,000 By NYC Zip Code



The first map shows the most common complaint types in each zip code. Strikingly, almost no zip codes were dominated by a single complaint type. Even housing and buildings, which ran away with the title of most common complaint type, was strikingly missing from the map in which almost no zip codes were over 10% comprised of any one complaint type. The second map makes this very clear - it is visible that almost every zip code had an incredibly low proportion of complaints filed from their leading complaint type.

The third map, showing median income, is almost a damning tale of income inequality in NYC. All of the wealth is clustered in Manhattan (and a chunk of Brooklyn just across from it), while no zip code reaching even the third category of income is found outside of that area. Strikingly, however, this does not seem to match up very well with any other map considered yet - Manhattan did not stand out in the vast majority of complaint variables yet considered. This more or less continues in the fourth map, where there is a slightly higher proportion of per-100K complaints filed in Manhattan zip codes than the rest of the city but not really anything to speak of. The biggest anomaly here, actually, is Staten Island - which seems quite high in terms of per-100K complaints filed. In all fairness, if I lived in Staten Island, I too think I would file a substantial number of complaints.

4. Short final discussion: Scope and limits of this data

There are certainly some limits to these data. While these data are very thorough in terms of the location that the complaint originated from, what it pertained to, and when it was filed (all excluding NA values, of course), they are seriously lacking in some important other information, including the demographics of who filed the complaint. That would be important information for any government agency to include, as they should strive to see exactly who it is that uses their system and adapt their system to meet the people who are being underserved by it. Unfortunately, they are unable to do so.

That also serves as a major drawback in these plots. Instead of being able to aggregate complaints across specific demographics, we are reduced to looking at them in zip codes (or other geographic areas) composed of a *plurality* of each demographic group. Great amounts of information are lost in doing this, and a lot of assumptions are made. Though the plots generally seem to make sense in the context of these demographics, the specific information would bolster them greatly.