

STAT 133 Final Project: Visualizing Berkeley Police Data

Amy Zhu, Mengyu Li, Rebecca Reus, and Shangjun Jiang

9 August 2016

Abstract

For the STAT 133 final project, our group was interested in finding some ways of visualizing the relationship between the Berkeley Police Department and different groups of people, especially within the context of different populations and neighborhoods within Berkeley. To obtain our results, we accessed, cleaned and explored the four Berkeley Police Department data sets available to the public on the Berkeley Open Data website using R, the statistical programming language. These data sets were of police stops, calls for service, arrests, and jail bookings. There were 16,255 stops (including stops for traffic, suspicious vehicle, bicycle, and pedestrian) from a 16-month period that began on January 26, 2015 (in compliance with General Order B-4, which ordered this data set to be collected for the public) and ending on April 30th, 2016. The calls of service included 4,911 calls, which were collected over the most recent a 6-month period, from January 20, 2016 to July 13, 2016. The arrest data included 232 arrests from the previous 1-month period, from June 20, 2016 to July 19, 2016. The jail bookings data set included 248 jail bookings over the same 1-month period. In addition to the police data, we accessed the Berkeley city and the Alameda County 2010 census data in order to obtain demographic information about Berkeley as a whole and about the demographics of the neighborhoods inside of Berkeley. The three census data sets that we used are also available through Berkeley's open data site, including the Berkeley 2000-2010 census data, which gave us the total population counts and percentages of each group of people in Berkeley for the 2010 census, the Alameda County 2010 census tract population data, which gave us the demographic information for 33 different neighborhoods (or "tracts") inside Berkeley, and the Berkeley 2010 census tract polygons shapefile (.shp), which provided the information we needed to visually display the 33 neighborhoods. Rather than limiting the project to one set of data, incorporating multiple data sets gave us a more holistic understanding of subject and expanded our ability to understand and visualize the data. Our findings suggest the hypothesis that (1) black people in Berkeley are more likely than members of another group of people to be stopped by Berkeley police, and that (2) black people in Berkeley are arrested and jailed at disproportionate rates. We recommend future statistical analysis be done to test our hypothesis.

Introduction

First we will detail our process of cleaning the police and census data and obtaining spatial information for the police data. Next, we will provide our analysis, which includes visualizations of the police and census data and some preliminary findings related to these visualizations. Afterwards, we will discuss problems with the data and recommendations for future study on this subject.

Cleaning Data

Berkeley Stop Data Cleaning Process

The original stop data set included 16,255 stops recorded by the Berkeley police over a 16-month period. For each stop, the police recorded the date and time, location (as a character string, usually in the form of a street block or a street intersection), incident type (referring to whether the stop was traffic, pedestrian, bike, or suspicious vehicle), and a column called disposition. The disposition column was a character string typically including 6 characters per individual as follows: the first character was race, the second gender, the third age range, the fourth reason (for the stop), the fifth enforcement (the result of the stop), and the sixth

car search (yes or no). Each individual that was involved with the stop was recorded using comma separation in the same row and column followed by an additional 6 characters. Additional optional dispositions were sometimes included, again using comma-separation from the 6-character strings (though not necessarily in an order), ranging from one to three characters. These additional disposition abbreviations represented the following options: Primary Case Report, MDT Narrative Only, Arrest Report Only (No Case Report Submitted), Incident Report, Field Card, Collision Investigation Report, Emergency Psychiatric Evaluation, Impounded Vehicle, and Officer Made a Stop of More Than 5 Persons. We found it reasonable to assume that these additional dispositions were related to police paperwork procedures, and since they were used so sparingly, we decided to not consider most of them. However, in spite of only being used less than 20 times, we found the Emergency Psychiatric Evaluation disposition to be of interest, and so we chose to investigate this one in addition to the main 6 character set of dispositions, though we chose not to focus on it during the limited scope of this paper.

Cleaning the stop data was relatively time-consuming compared to the other data sets. First, the dispositions column in the original data set contained many columns and sometimes rows worth of information, though the comma-separated pieces were without order (for example, sometimes the optional disposition abbreviations came before the 6-character string instead). We handled this problem using data tidying techniques, such as separating disposition information into separate row entries for each individual assessed, in the case of multiple persons stopped, and splitting the column by isolating the six character dispositions into different columns. The optional disposition character strings were moved and separated into more columns.

Additional cleaning for the stop data included changing the stop date and time variable to the `lubridate` date format. We also created an hour column and day column in order to examine the stop occurrences by hour over a 24-hour period and across a 7-day week.

We used the stop data's character string locations into location coordinates using a Google Maps geolocation service provided by the library `ggmap` with the function `geocode`. This task proved difficult because of common data entry mistakes. For example, there were common misspellings of street names and the use of several different abbreviations to mean the same thing. Additionally, we found that the geolocation provided by Google does not accept forward slash characters in the column to indicate street intersections, so these forward slash characters were changed to `and` for the location processing, which was successful. Google also does not accept the term "block" to indicate a range of addresses, and so the word "block" was removed from the location column for location processing. Most of these types of errors were dealt with on an automated basis before geolocation, using `str_detect` and `str_replace` from the `stringr` package. The number of problems with the location data that were left over afterwards were small enough to handle on a case-by-case basis. Further, because Google limits geolocation to 2500 queries per 24-hour period, it took approximately a week to process all the locations into coordinates from the stop data (this includes the time it took to fix spelling errors).

The reason why the Incident Number column was not removed was because two teams from our group worked on cleaning different parts of the stop data at the same time; one team worked on separating the disposition column into variables, while the other team worked on creating the columns for longitude and latitude. After both teams had finished their cleaning task, the resulting data frames were then merged by Incident Number.

At the end of the cleaning process, the stop data had been whittled down from 16,255 observations to the 14,291 police stops (due to missing disposition values) that we will use in the following analysis. Please see Figure 1 for a glimpse at the clean police stops data frame.

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |
|--------------|------------------|-------------|----------|--------------|-------|------------|--------------|-------|--------|---------|---------------|-------------|-----------|------------|------------|--------------|-------------|------|-----|
| Incident.Nun | Call.Date | Tir | Location | Incident.Typ | Other | Individual | Dispositions | Race | Gender | AgeRang | Reason | Enforcement | CarSearch | lat | long | Arrest | Emergency.P | Hour | Day |
| 1 | 2015-000048##### | 2000 | BLOCK | 1194 | NA | 1 | BM4ICN | Black | Male | 40+ | Investigation | Citation | No Search | 37.8726801 | -122.27074 | Not Arrested | No | 7 | 2 |
| 2 | 2015-000048##### | 1700 | BLOCK | 1194 | NA | 1 | BM4ICN | Black | Male | 40+ | Investigation | Citation | No Search | 37.8731719 | -122.2938 | Not Arrested | No | 7 | 2 |
| 3 | 2015-000048##### | M L KING JR | T | NA | | 1 | OF4TCN | Other | Female | 40+ | Traffic | Citation | No Search | 37.8716087 | -122.27303 | Not Arrested | No | 9 | 2 |
| 4 | 2015-000048##### | M L KING JR | T | NA | | 1 | OM4TCN | Other | Male | 40+ | Traffic | Citation | No Search | 37.8716087 | -122.27303 | Not Arrested | No | 10 | 2 |
| 5 | 2015-000048##### | UNIVERSITY | T | NA | | 1 | OF2TCN | Other | Female | 18-29 | Traffic | Citation | No Search | 37.8716087 | -122.27303 | Not Arrested | No | 10 | 2 |

Figure 1:

Berkeley Arrest and Jail Bookings Cleaning Process

The number of observations available for the arrest data set and the jail data set were significantly smaller than the number of observations for the stop data set. The arrest data had 232 observations and the jail bookings data had 248. These observations were obtained by the Berkeley Police department over the same month-long time period between June and July of 2016. They both contained similar variables, including a case/arrest/booking number, date and time, type, and subject information (name, race, sex, D.O.B., age, height, weight, hair, eyes, and occupation) and statute information (type, description, agency, and disposition). Cleaning this data set also required the dates and times to be put into `lubridate` format, and we also created an hour and day column. See Figure 2 for the cleaned arrest data frame and Figure 3 for the cleaned jailings data frame.

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q |
|---|-------------|--------------|---------------------|-----------------------|------------|------------|---------------|-------------|--------|--------|------|------|--------------------|--------------|------------------------------|---------------|---|
| 1 | Arrest Numb | Date and Tim | Arrest Type | Subject | Race | Sex | Date of Birth | Age | Height | Weight | Hair | Eyes | Statute | Statute Type | Statute Desc | Case Number | |
| 2 | 16347 | 06/20/2016 | COURT FILEC | Frank Moore Black | Male | 08/25/1982 | 32 | | | | | | Warr - Out (P PC; | | Outside Warrant Misdemeanor; | | |
| 3 | 16431 | 06/30/2016 | ON-VIEW BY | Michael Harc Black | Male | 05/07/1992 | 24 | | | | | | Warr - Out (P PC; | | Outside Warrant Misdemeanor; | | |
| 4 | 16349 | 06/21/2016 | COURT FILEC | Julia Elizabeth White | Female | 02/26/1998 | 18 | 5 Ft. 0 In. | | 90 | BRO | HAZ | Warr - Out (P PC; | | Outside Warrant Misdemeanor; | | |
| 5 | 16345 | 06/20/2016 | ON-VIEW BY | Robert Lee C Black | Male | 06/02/1952 | 63 | 5 Ft. 9 In. | | 130 | BLK | BRO | 1203.2 - F; 4; PC; | | Probation Vi | 2016-00036449 | |
| 6 | 16346 | 06/20/2016 | ON-VIEW BY | LAUREN LOU White | Female | 06/05/1985 | 30 | 5 Ft. 5 In. | | 138 | BRO | GRN | 243 (E)(1); | PC; | Battery: spo | 2016-00036562 | |
| 7 | 16594 ##### | SUSP. OF FEL | Jalisa Nicole Black | Female | 03/08/1990 | 26 | 5 Ft. 2 In. | | 140 | BLK | BRO | 459; | PC; | Burglary: | 2016-00040113 | | |
| 8 | 16582 | 07/17/2016 | ON-VIEW BY | Ismael Valen Hispanic | Male | 05/26/1990 | 26 | 5 Ft. 7 In. | | 145 | BLK | BRO | 23152 (A) - N VC; | | DUI: Alcohol | 2016-00042026 | |
| 9 | 16353 | 06/21/2016 | ON-VIEW BY | Cheryl Denis Black | Female | 05/31/1983 | 33 | 5 Ft. 6 In. | | 130 | BLK | BRO | 1203.2 - M; € PC; | | Probation Vi | 2016-00036799 | |

Figure 2:

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |
|---|-------------|-------------|-----------------------|--------|------------|---------------|-------------|--------|--------|------|------|-------------------|-------------------|--------------|---------------|-----------------------|-----------------------|-----------------|-----------------|---|
| 1 | Booking Nun | Booking Dat | Subject | Race | Sex | Date Of Birth | Age | Height | Weight | Hair | Eyes | Occupatio | Statute | Statute T | Statute Desc | Arrest Date | € Case Numbe | Booking Age | Disposition | |
| 2 | 2016-00001 | 06/20/2016 | LAUREN LOU White | Female | 06/05/1985 | 31 | 5 Ft. 5 In. | | 138 | BRO | GRN | 243 (E)(1); | PC; | Battery: spo | 06/20/2016 | 2016-000365 CA0010300 | | BAILED | | |
| 3 | 2016-00001 | 06/21/2016 | ALONDRECK Black | Male | 01/14/1993 | 23 | 6 Ft. 2 In. | | 175 | BLK | BRO | Warr - Out (P PC; | | Outside War | 06/21/2016 | 2016-000014 CA0019700 | | CITE-JAIL | | |
| 4 | 2016-00001 | 06/21/2016 | Julia Elizabeth White | Female | 02/26/1998 | 18 | 5 Ft. 0 In. | | 90 | BRO | HAZ | UNEMPLO | 11364 (A); | HS; | Possess narc | 06/21/2016 | 2016-000366 CA0010300 | | SANTA RITA JAIL | |
| 5 | 2016-00001 | 06/22/2016 | JORDAN ELIJ Black | Male | 04/05/1997 | 19 | 6 Ft. 1 In. | | 145 | BRO | BRO | | 647 (H); | PC; | Disorderly cc | 06/22/2016 | 2016-000366 CA0010300 | | CITE-JAIL | |
| 6 | 2016-00001 | 06/22/2016 | Aneicia John Black | Female | 08/31/1996 | 19 | 5 Ft. 4 In. | | 115 | BLK | BRO | Warr - Out (P PC; | | Outside War | 06/22/2016 | 2016-000366 CA0010300 | | SANTA RITA JAIL | | |
| 7 | 2016-00001 | 06/21/2016 | Davin Williar White | Male | 04/01/1971 | 45 | 6 Ft. 3 In. | | 190 | BLU | BLU | | 23152 (A) - N VC; | | DUI: Alcohol | 06/21/2016 | 2016-000367 CA0010300 | | CITE-JAIL | |

Figure 3:

Berkeley Census 2010 Data Cleaning Process

The Berkeley census 2010 data was obtained from three separate data sets, as previously mentioned. We first considered the summary information provided by the city website, which was easy to clean and use in our analysis, since it provided a summary of the demographic information for each group in Berkeley as a whole. However, in order to understand the demographics within Berkeley in a way that could provide a spatial visualization, we needed to use the census 2010 Alameda County population tract data, which provided the population and race data by tract number for the entire county. We manipulated this data into a data frame containing the 33 census tracts inside Berkeley, each row representing a single tract and the columns providing demographic information for that tract. To be able to map these 33 tracts, we also needed the census 2010 Berkeley tract polygons shapefile. After converting the shapefile to a mapping data frame and changing the coordinate system to the one used by Google (this was done using the `readOGR` function from the library `rgdal`), we were able to `left_join` this new polygonal mapping data frame to the data frame of the 33 census tracts by the unique census 2010 tract number. This new data frame allowed us to map population information for 33 different areas. The goal was to visualize the population density of different groups of people within the city by neighborhood, which provided us with a context for the sptatially visualized police data.

| | A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|-----------|------------|----------|-------|-------|-------|-------|--------------|-------------|--------------|---------------|---------------|------------------|
| 1 | Id | Population | Hispanic | White | Black | Asian | Other | Percent.Blat | Percent.Oth | Percent.Whit | Percent.Asiat | Percent.Hispt | Percent.Berkelet |
| 2 | 4238 | 2925 | 128 | 2549 | 49 | 182 | 29 | 0.01675214 | 0.00991453 | 0.87145299 | 0.06222222 | 0.04376068 | 0.02598152 |
| 3 | 4222 | 3144 | 358 | 1882 | 402 | 484 | 155 | 0.1278626 | 0.04930025 | 0.59860051 | 0.15394402 | 0.11386768 | 0.02792681 |
| 4 | 4224 | 4196 | 343 | 2324 | 191 | 1360 | 110 | 0.04551954 | 0.02621544 | 0.55386082 | 0.32411821 | 0.08174452 | 0.03727127 |
| 5 | 4225 | 4658 | 323 | 2889 | 104 | 1290 | 97 | 0.02232718 | 0.02082439 | 0.62022327 | 0.27694289 | 0.06934307 | 0.04137502 |
| 6 | 4223 | 3387 | 311 | 2224 | 196 | 651 | 106 | 0.05786832 | 0.03129613 | 0.65662828 | 0.19220549 | 0.09182167 | 0.03008527 |
| 7 | 4218 | 2007 | 100 | 1587 | 59 | 219 | 25 | 0.02939711 | 0.0124564 | 0.79073244 | 0.10911809 | 0.04982561 | 0.01782732 |
| 8 | 4236.01 | 2642 | 203 | 1934 | 92 | 385 | 90 | 0.0348221 | 0.0340651 | 0.7320212 | 0.14572294 | 0.07683573 | 0.02346776 |
| 9 | 4216 | 3558 | 188 | 2872 | 59 | 368 | 66 | 0.01658235 | 0.01854975 | 0.80719505 | 0.10342889 | 0.05283867 | 0.03160419 |

Figure 4:

Data Analysis

Stop Data Analysis

Figure 5 below provides a map of the location information of the Berkeley police stop data mapped by location. We have excluded some observations on the outskirts of Berkeley in order to zoom in to the street level. Note that the `alpha` variable in `ggplot` is set to a small number in order to show the incidents that occur multiple times in the same place. In other words, the darker the point, the more police stops have occurred at that location.

Berkeley Police Stops, 2015–2016

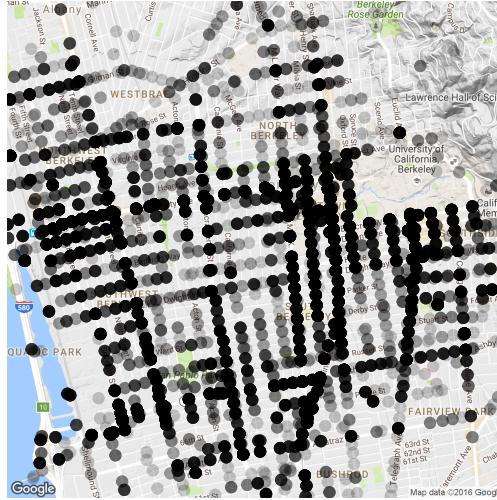
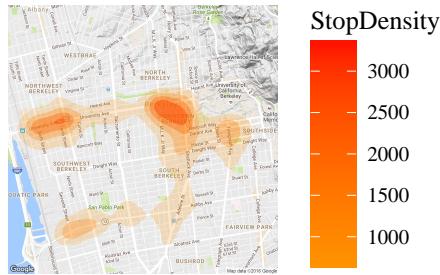


Figure 5: Berkeley Police Stops

Another visualization of the stop data mapped in Figure 5 is the two-dimensional density map of the police stop data in the left pane of Figure 6. One of the most interesting aspects of the stop data in this image is that most police stops seem to be clustered west of UC Berkeley Campus in the Downtown area. This was contrary to what some of us might have thought before, which was that the area with the higher police stops would be the area with the higher population. To visualize the relationship between police stops and population density, consider the map in the right pane of Figure 6, which visualizes the percent population of Berkeley for each of the 33 census 2010 tracts. As you can see, the least densely populated area is the census 2010 tract containing the University of California. This makes sense because the area of this tract includes a campus and some student housing areas, where permanent residents are less likely to live. One of the most population-dense areas is Southside, the portion of Berkeley directly south of the UC campus. Note that even though the tract including the Southside area is the most densely populated of the 33 census 2010 tracts in Berkeley, it is not the area where most police stops occur.

Berkeley Police Stop Density Map, 2015–2016



2010 Population Density

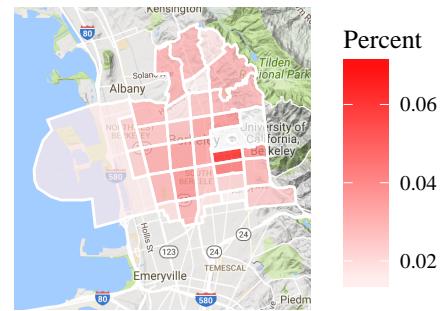


Figure 6: 2015–2016 Police Stop Density (left) and 2010 Census Tract Population Density (right)

In order to break down the stop data further, we mapped the data by each of the disposition variables. We found race to be one of the most interesting variables. Consider the stop data density map faceted by race in Figure 7.

Berkeley Police Stop Density Map by Race

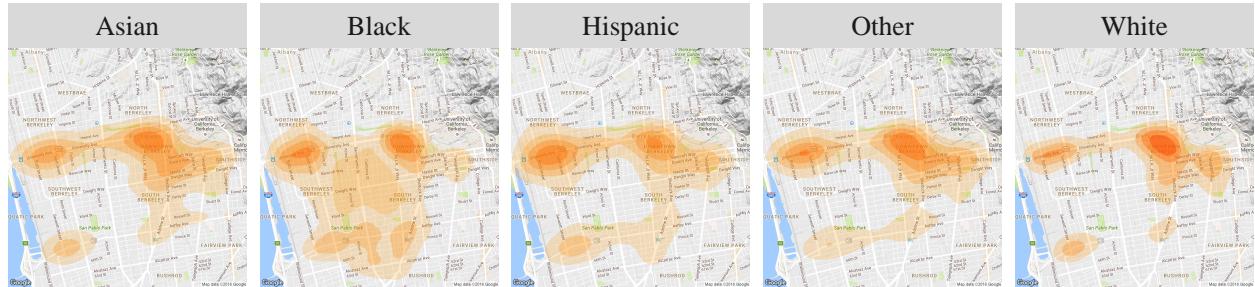


Figure 7: Police Stop Density Map by Race

Figure 7 seems to suggest that the group of people stopped the most in Berkeley are African Americans. For instance, in Central Berkeley, the stops for the races defined by the Berkeley Police Department as Asian, Hispanic, White, and Other seem comparatively lower than for Blacks in the same area. In order to consider the context of these stops, we considered the population distribution for these areas, shown in Figure 8. Note that the population density visualized by Figure 8 is conveyed by the range between white (0%) and red (100%). As you can see, White people make up a striking majority in all census 2010 tract areas, which seems to suggest the hypothesis that the population of Black people in Berkeley is being disproportionately affected by police stops.



Figure 8: Berkeley Population Density Map by Race

Now, with the population percentages of each group in hand, we were able to compare the results of our stop data analysis to the percent population of each group. In Figure 9 (left), the **Percent Stopped** chart conveys

the percent of each group stopped compared to the total number of people stopped. In Figure 9 (right), the Percent Population Stopped chart conveys the proportion of people in each group stopped relative to the population of that group living in Berkeley. Even though the people stopped by police are not necessarily Berkeley residents, the Percent Population Stopped chart suggests that some groups of people in Berkeley might be more likely to be affected by police stops, especially Black people.

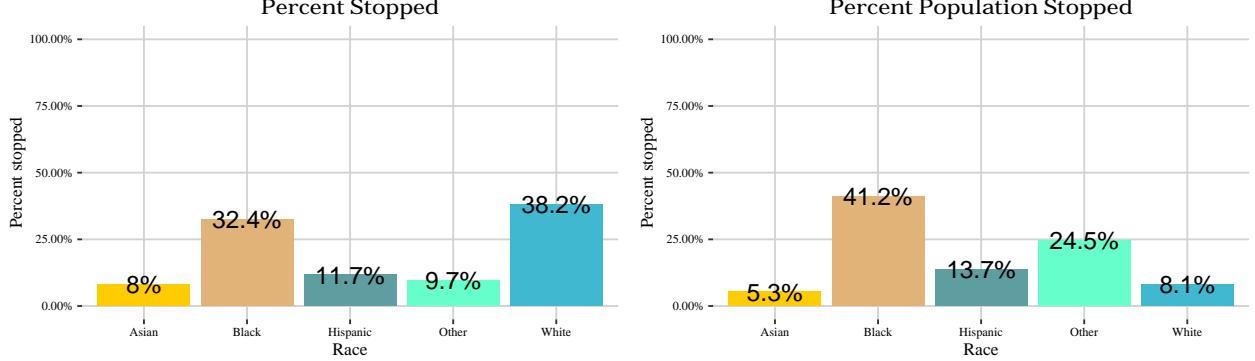


Figure 9: Police Stops in Berkeley by Race

Another disposition variable we found to be of interest was the Enforcement variable, which has the options Arrest, Citation, Warning, and Other. We were interested to see how these enforcements were applied to different groups of people. Consider Figure 10 (top), which maps the stops that led to the enforcement of arrest. The number of incidents of police stops mapped for Blacks seems higher than for the other groups. Upon examination of the chart in Figure 10 (bottom), we see that of the 441 individuals stopped and then arrested as a result of the stop during this 18-month period, Blacks were the majority. Though the Berkeley police department is prohibited from the use of biased policing practices, Figure 10 does suggest that the enforcement of arrest may be applied differently to different groups of people.

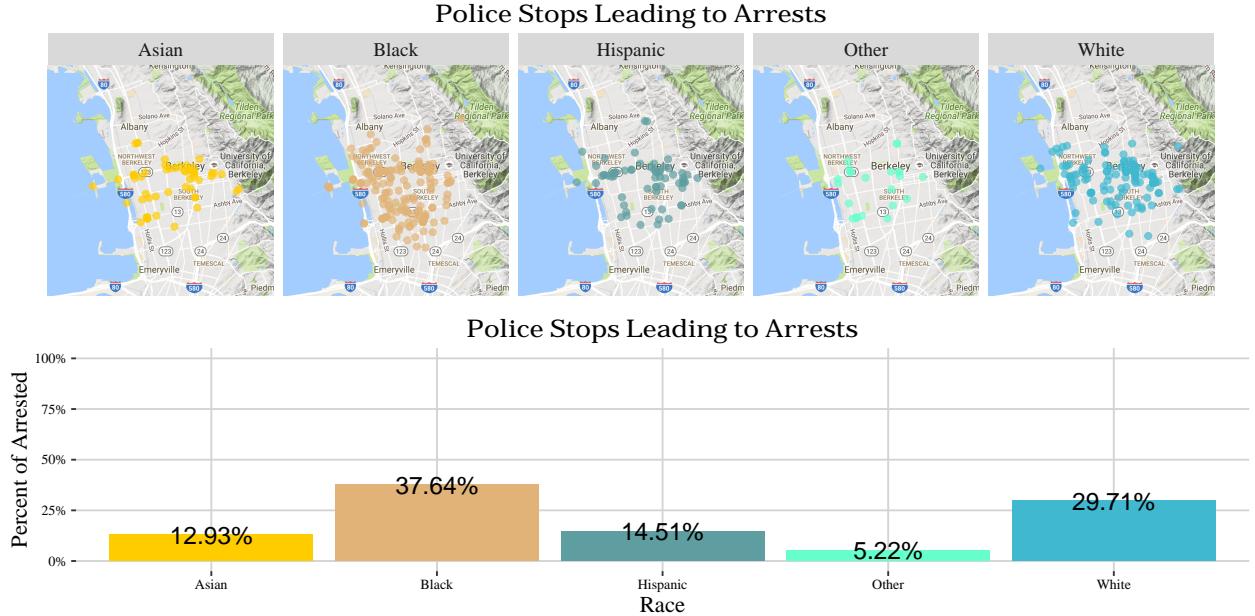


Figure 10: Police Stops Leading to Arrests

A third interesting disposition variable was the Reason disposition, which indicates the reason for the police stop in each case. The reasons available for police stops are Investigation, Other, Probation/Parole,

Reasonable Suspicion, Traffic, and Wanted. One way we explored the data initially was to see how these reasons affected different groups of people, and we would recommend exploring this idea in depth at a later time. Due to the length and scope of this paper however, we decided to end our discussion of the stop data by showing that most of the Berkeley police stops are traffic stops. Mapping the coordinates by reason for the stop, as in Figure 11, suggests that this is the most common reason for a stop in Berkeley.

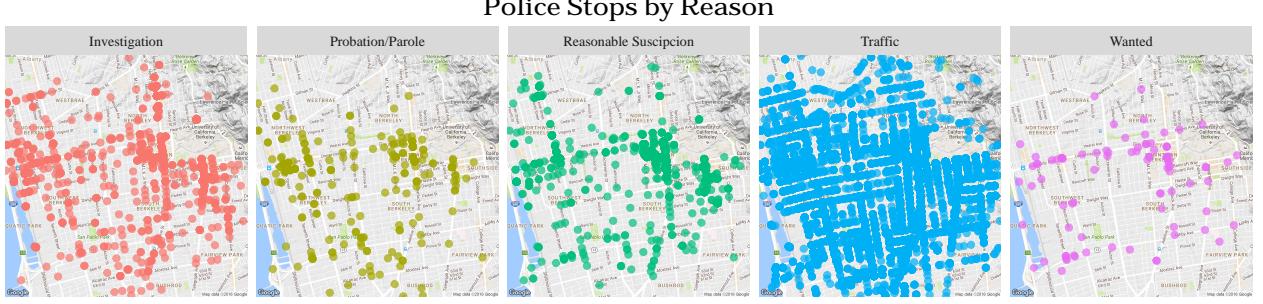


Figure 11: Police Stops by Reason

We also analyzed possible preferences of Berkeley police officers when arresting individuals. We calculated the conditional probability of an individual being arrested given a specific reason for the stop and given a specific race / age range / gender. That is to say, we calculated the probability a person will be arrested if he was stopped for a specific reason and was of a specific race / age range / gender.

For Figure 12, we can tell that the probability of a stopped individual being arrested in a stop with the reason being traffic is much lower than other stop reasons. The average conditional probability of being arrested for a traffic stop is 1.58%. Additionally, the probability of being arrested in a stop with the reason wanted is much higher than other reason. The average conditional probability of being arrested given being stopped for the wanted reason is 33.33%. Lastly, an interesting fact is that the conditional probability of being arrested in a stop for probation or parole and race Asian is 75%, which is much higher than that of any other race. Our conclusion is that Asian people are much more liable to be arrested in Berkeley during their probation or parole if stopped by police.

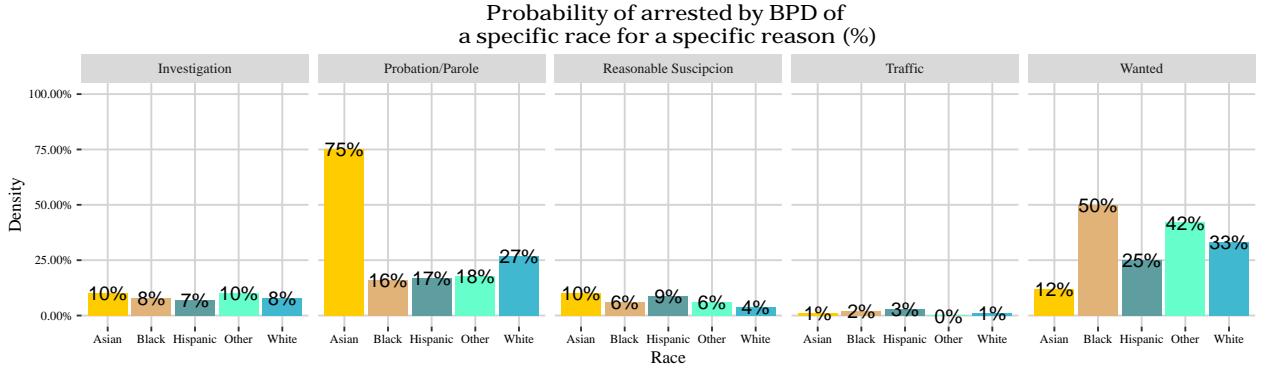


Figure 12: Conditional Probability of Being Arrested During a Police Stop: Race VS Reason

In Figure 13, the probability of being arrested as a result of a police stop in Berkeley for a traffic reason is again the lowest and wanted is the highest. Again similar to the conclusion in Figure 12, an interesting fact is that the conditional probability of being arrested in a stop with the reason probation or parole and age range 0-18 is 67.56%, which is much higher than that of any other age range for that reason. Lastly, the conditional probability of being arrested in a police stop given an individual in the age range 0-18 given any reason is more than that of any other age range. This suggests that in a police stop, teenagers are more liable to be arrested in Berkeley.

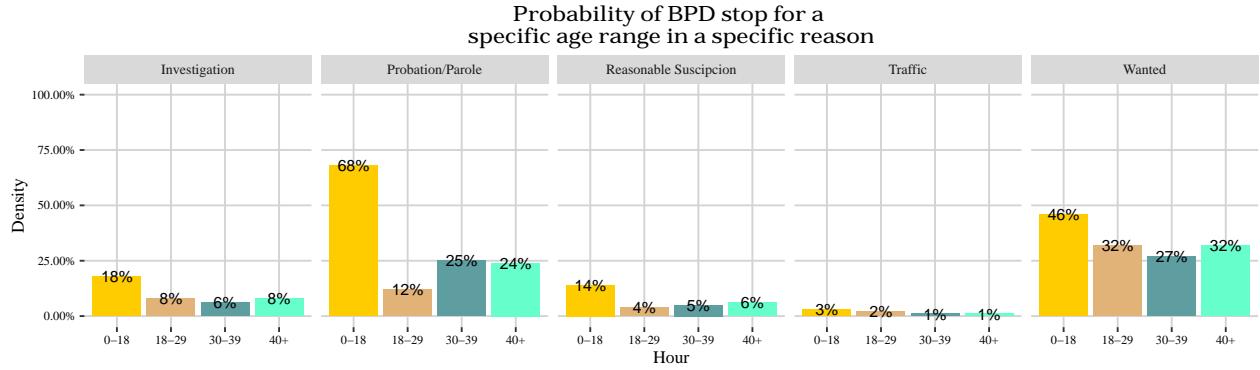


Figure 13: Conditional Probability of Being Arrested During a Police Stop: Age VS Reason

Again for Figure 14, the conditional probability of being arrested during a police stop for the stop reason of traffic is the lowest and wanted is the highest. Again similar to the conclusion above, an interesting fact is that the conditional probability of being arrested in a police stop with the stop reason probation or parole given a gender of female is 50.00%, which is much higher than that of male, which is 20.11%.

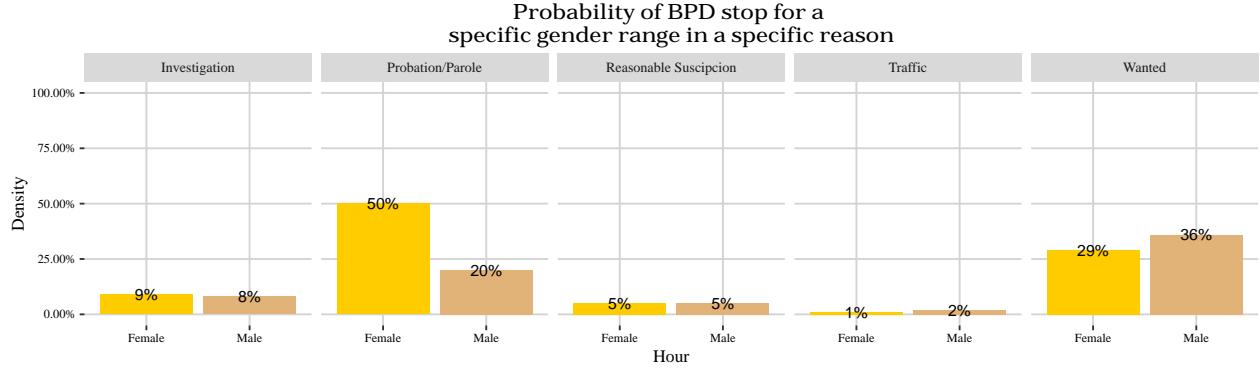


Figure 14: Conditional Probability of Being Arrested During a Police Stop: Gender VS Reason

Figure 15 is interesting because it suggests that black people are more likely to be stopped by the police in Berkeley during the hours of night. The average ratio of black people to all people stopped by the Berkeley Police Department at night is about 40%. Figure 15 also suggests that white people are more liable to be stopped by the police in Berkeley during the hours of daytime. The average ratio of white people to all people stopped by the Berkeley Police Deparment during the day is about 45%. While the ratio of Asian people, Hispanic people and other people stopped by police fluctuate during the daytime and the night, with an average ratio of 8%, 11% and 9% respectively.

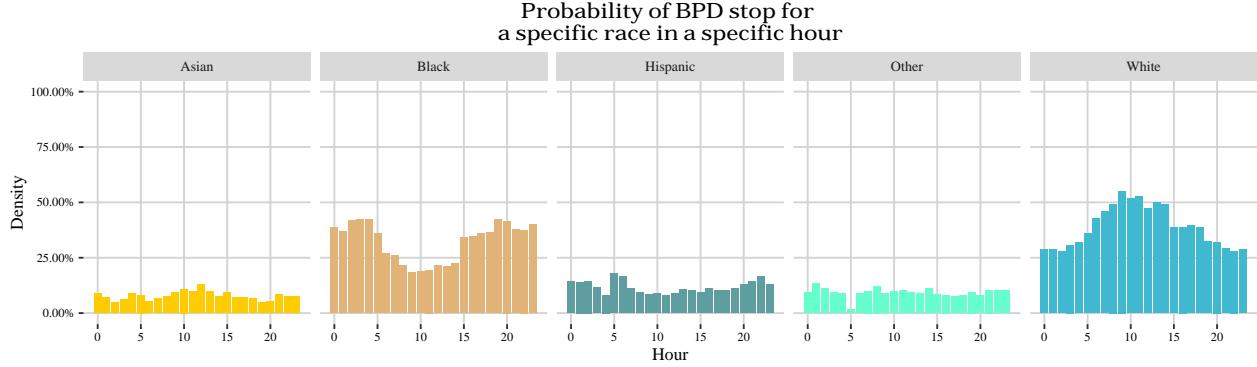


Figure 15: Conditional Probability of Being Arrested During a Police Stop: Race VS Hour

Figure 16 implies that people aged from 18 to 29 are more liable to be stopped by police in Berkeley at night. The average ratio of people aged from 18 to 29 to all people stopped by the Berkeley Police Department at night is greater than 40%. The figure also implies that people aged 40+ are more liable to be stopped by the police during the daytime. The average ratio of people aged 40+ to all people stopped by police during the day is greater than 40%. Further, the ratio of stops of people aged 0-18 and people aged 30-39 fluctuates during the daytime and the night, with an average ratio of 2.5% and 25% respectively.

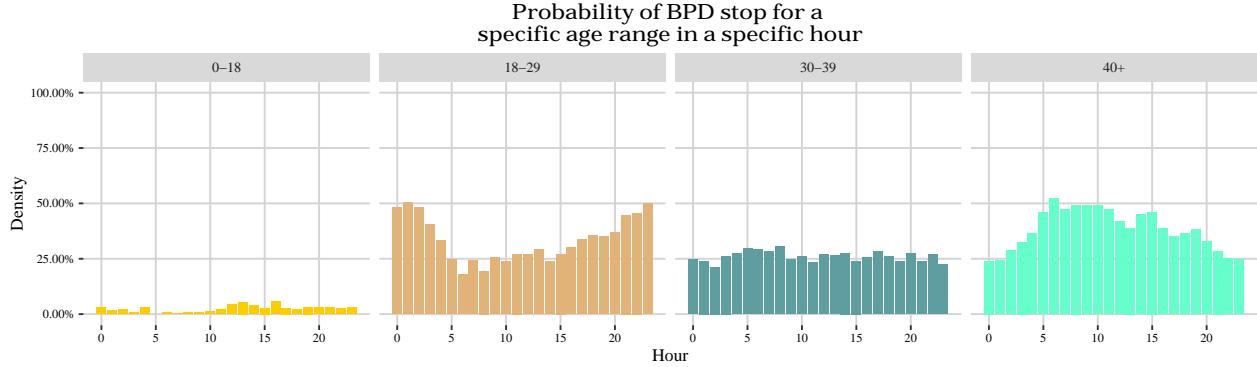


Figure 16: Conditional Probability of Arrested During a Police Stop: Age VS Hour

Berkeley Police Jail and Arrest Data Set Analysis

Moving on from the stop data, we considered the available jail data and arrest data for the city of Berkeley, which is not necessarily dependent on the stop data. Two of the most interesting variables to consider in these two data sets were also race and age. We decided to examine the percent of people arrested and the percent of people jailed in Berkeley by both age and race. Figure 17 is a visualization. Note that the following percentages are related to the number of total people jailed or the number of people arrested during a recent 30 day time frame, respectively, and we have not adjusted the data relative to the population of Berkeley. It is interesting to note that the proportions of arrests and jailings during this time are nearly identical, since these arrests and jailings occur during the same time frame. The data suggests that being jailed after being arrested in Berkeley is almost certain, and that this is true for all races. However, note that the proportion of Blacks being arrested and jailed is more than the other races, in spite of Berkeley being a city of mostly White people.

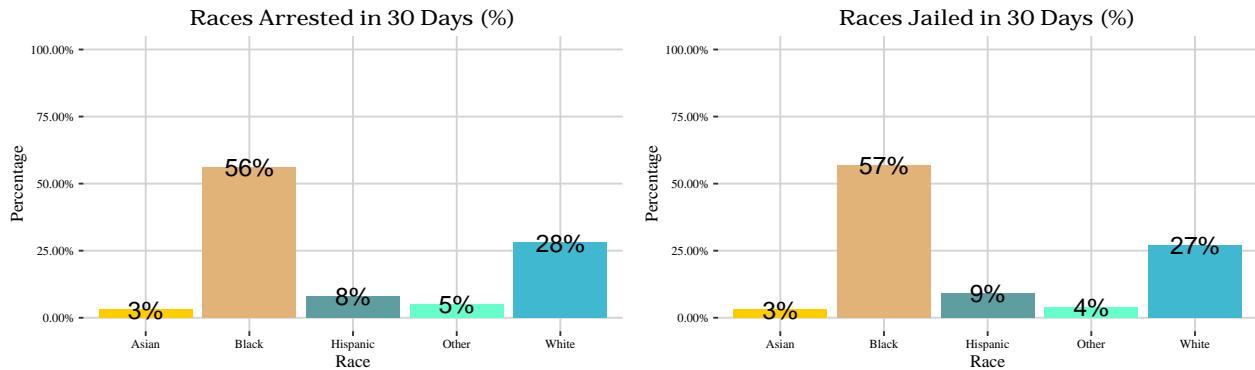


Figure 17: Arrested and Jailed in Berkeley

Problems

The data sets provided by the police (stops, arrests, and jail bookings) have several problems. First, the open data site updates the data frequently and we are concerned that the data available to be used for this paper won't be available after the data is updated, since the early data is removed. Second, concerning the methods for recording disposition data like race and gender, the Berkeley Open Data website indicates that the police officers conducting the stops are deciding these dispositions on their own, meaning that they are judging whether or not the individual apprehended is male or female, black or hispanic, etc., meaning that the data could be misrecorded. Even though there doesn't seem to be a means to get around this problem with the actual method of collecting the disposition data for stops that don't lead to citation or arrest, it is still important to address it as something that could induce error into a statistical analysis of the data.

Another problem with the stop data is the many data-entry problems we faced in the character string location column (which sometimes included only 1 street), since there is likely some stops that are not mapped correctly due to approximations made by the Google geolocation service (for example, in some cases, only having one street name instead of a corner caused Google to estimate the exact location based on few conditions). Though most of the location problems were fixed using string detection, there is likely some error in our visual representation of the stop data maps. Also, the website indicates that any specific addresses were changed to "block" address form to preserve anonymity; therefore the spatial data cannot be expected to be of a precision smaller than a city block.

In regards to the racial data from all of the police data, it is important to note that the census data we used to compare the stop data to included many more categories for race than just the five used by police (Asian, Black, Hispanic, Other, and White). To deal with this problem, other races included by the census data but not included by the police data were combined with the category of other. Additionally, the census data had a category for Two or More Races, which we also included as Other in order to compare the data sets. While the census data allowed for overlap between racial categories, there was no overlap for the racial categories in the police data. This creates a problem when comparing the two data sets, and so it should be acknowledged.

Finally, it should be noted that a more accurate analysis of Berkeley police data would also incorporate demographics for the homeless population as well as for the UC student population, both of which are not included in the census data. Including these data sets would provide a more nuanced context for the stop data.

Conclusion

The goal of this paper was to investigate and visualize Berkeley Police interactions with the community through the available police data. Although our preliminary findings unfortunately seem to suggest that certain groups of people, especially black people, are disproportionately likely to be stopped, arrested, and

jailed in the city of Berkeley, we recommend a more rigorous statistical analysis of the data be done in order to assess this claim. We also recommend an idea we had for this project to future researchers that could make the Berkeley police data more accessible. For example, one could make the data more accessible by creating a `shiny` app that allows the user to toggle between the different dispositions of the data (race, age range, enforcement, reason, etc.) and between modes of view (map or bar chart, for example) or order to play with the data in a visual way. Another `shiny` app idea we had was to filter the calls for service data by proximity to certain locations so that users could use the app to select a location and find a visual displaying the most common types of calls in that area. It would also be useful to find a way to create an app that can easily be updated when the new data is posted (every month), allowing users of the app to get the most up-to-date information.

Sources

- Stop Data (16,000)
- Arrest (200)
- Jail Bookings (250)
- Calls for Service (4,000)
- Berkeley Census Data
- Census 2010 Population and race data by county tract polygons
- Berkeley Census 2010 Tract Polygons