

INTERACTIVE DATA SCIENCE REPORT

Team: Crime Scene

Members: Aaron Ho, Natasha Ninan, Neema Nayak, Swapnil Kadakia

GitHub Link: <https://github.com/CMU-IDS-2022/final-project-crime-scene>

Video URL:

Abstract

The data science problem we aim to address is how human well-being in cities and states influences hate crime patterns over the years. Our approach to this problem is to address these issues and offer a broader exploration into the behaviors, demographics, and other defining features behind the causes of hate crimes. Our approach first provides a high-level overview of the prevalence of hate crimes in the U.S. from 1990 to 2020. Pre-processing the data, clustering, correlation, feature importance, the timeline of hate crimes, hate crime distribution across the United States, variables contributing to the hate crime rate of each state, and analyzing hate crime statistics for cities in each state are all steps in the process. This analysis aligns with our expectation of multiple factors resulting in hate crimes. Our efforts will attempt to simulate that as much as possible. Users can explore the patterns in demographics and behaviors of a particular state and the prevalence of hate crimes in that state. The expectation is that the features associated with a state (i.e., education, racial diversity, income) would help explain patterns in hate crimes in that state. In exploring this relation, we aim to clarify what influences hate crimes and inform policy actions and decisions to address these crimes in cities and states. The further study generated from the conclusions of these results would be to explore the most defining features and focus research on how these features contribute to biases and crime.

Introduction

In much of American history, there exists a pervasive and persistent tension between people of racial and ethnic groups. This tension has evolved into hate crimes that range from intimidation, vandalism, assault, and murder in recent years. Despite the best efforts of educators, policymakers, and community organizers, tensions seem to increase along with corresponding crime. The rise in recent hate crimes was quite disheartening. With many of us having loved ones who live in the U.S., we wanted to understand why these incidents happen so that we can develop areas for improvement.

Crime is generally well-studied in academic and public spaces. Crime data is readily accessible through local and federal law enforcement agencies such as the FBI. Hate-based crime is readily available for viewing through online tools created by these agencies and readily downloaded and manipulated for further analysis. However, these tools are limited in their purpose and conclusions.

For the public, current visualizations and tools for exploring hate crime incidents are often compiled into dashboards or graphs showcasing simple statistics of a particular crime. While this allows for creating a

concise and pertinent data-driven story, it is a straightforward analysis of the data that may not be very informative. It opens several problems regarding the flexibility of the data exploration with other correlated features and the robustness of any conclusions made from the data. The graphs, charts, and even tables may not reflect the particular purpose that the user has with the data.

When our team examined news articles behind recent hate crime incidents and reflected upon our own experiences, we collectively felt that it wasn't sufficient to conclude that people needed to discriminate less and have less bias. A common explanation is that education is necessary to correct bias and reduce future criminal activity. While we wholeheartedly agree that this is an essential component of addressing this issue, hate crimes are still complex occurrences of many factors. Therefore, a singular one-fits-all approach would not address hate crime patterns sufficiently.

We have seen that the causes behind crimes and specifically hate crimes, are expected to be highly multi-dimensional and, thus, impossible to assign a single reason to an incident. Additionally, many studies view the prevalence of hate crimes as attributed to trends or spikes. For example, the 2020 COVID pandemic saw a spike in Anti-Asian hate crimes due to negative stereotyping of Asian-Americans concerning COVID-19. One would expect that with the increase in high school completion or education, the hate crimes would decrease. This reality makes us realize that additional factors influence hate crimes in these cities. We are doubtful that these trends or spikes can be due to a particular reason and that perhaps, many other factors have been at work before the global pandemic.

A robust analysis of the causes behind hate crimes will need the flexibility to explore and analyze the data with additional or related datasets regarding human well-being that can help inform us about hate crime patterns. We intend to solve this problem by creating a robust application that considers various factors surrounding hate crimes to understand why these crimes happen in the first place.

Related Work

Studies on the true intention behind occurrences of hate crimes have shown that hate crime offenders are not always motivated by a single factor. They are driven by multiple factors. To fully understand the nature of these hate crimes, situational factors like location and victim-offender relationships need consideration. These factors may result in certain types of offenses and hate motivation. Research shows that hate crimes may also be the product of social environments. Hate crimes are more likely to occur where society is structured to benefit some communities like white, male, heterosexual, over others.^[1]

Studies in the psychology of offenders committing hate crimes suggest that hate crime offenders usually belong to one of four categories. Thrill-seekers are those motivated by thrill and excitement. Defensive are those motivated by a desire to protect their territory. Retaliators are those who act in retaliation for a perceived attack against their group. Missions are perpetrators who make it their mission to eradicate difference. The causes of hate crimes are not yet conclusive. However, there is evidence within social psychology that suggests that offenders might be influenced by their perception of some groups posing a threat to them. These threats are either realistic threats - competition over jobs, housing, and other

resources, and physical harm to themselves, or symbolic - threats posed to the values and social norms of people.^[1]

Research on countries with the least hate crimes enables us to identify reasons why they remained successful in minimizing such crimes. Denmark, known to be one of the safest countries, has a high level of equality and a strong sense of responsibility for social welfare. Denmark also provides its people with services and perks helping them to live comfortable lives. In addition to this, it provides its citizens with free access to healthcare, tuition-free education, and at-home care helpers for the elderly. Since everyone's basic needs are satisfied, and every citizen is respected, Denmark has maintained its low hate crime rates. Singapore has resorted to strict penalties even on minor crimes that have enabled them to keep overall hate crimes low. The government has strict controls on guns and firearms, resulting in a significant reduction in occurrences of violent crimes. The United States of America can follow suit and implement such measures within the parameters of its own socio-economic and political measures to improve its hate crime rates.^[2]

Methods

During the data preprocessing stage, we realized that to leverage some of the categorical variables of interest for our analysis, we needed to convert them to a binary representation. Features with values for a different category could be represented by one-hot-encoding (OHE). OHE allows us to convert a column with categorical variables with a binary representation. The result is an expanded version of the column with a binary encoding for whether the variable belongs to that record. We applied this principle to the OFFENDER_RACE column since each hate crime incident contained a distinct variable for the OFFENDER_RACE in the incident. Multiple races were represented with the designation OFFENDER_RACE_Multiple. Other features were trickier because they were multi-labeled features. However, we still wanted to represent them in a binary form. With OHE, we would use the OneHotEncoder class in sklearn to convert our categorical features into a binary format. To represent multi-labeled features, we would use MultiLabelBinarizer, also from sklearn. A notable example of this can be seen in our BIAS_DESC column. As we understand how hate crimes are motivated, hate crimes can be motivated by multiple biases. Therefore the BIAS_DESC can be considered a multi-label feature. Like OHE, the column would be expanded, except now multiple assignments can be made to a single incident if multiple biases are involved. The OFFENSE_NAME, VICTIM_TYPES, and LOCATION_NAME are also features that are multi-labeled, and MultiLabelBinarizer can be applied to them. Binarizing our categorical features allows us to better represent features in an interpretable and conducive way for analysis. It enables us to represent them as numerical values in our clustering algorithm and visualizations.

One set of techniques used on the FBI Hate Crimes dataset is a dimensionality reduction technique called UMAP and a clustering technique called DBSCAN. With the help of these two techniques, we created a series of charts that clustered on the data points and labeled them accordingly. The idea is to highlight similar pairwise points together so that by grouping them, we can determine what their similar attributes are. With dimensionality reduction, we can embed the data from a higher dimension into a lower-dimensional space that can still retain information or structure in the data. We knew that we

didn't want to use Principal Component Analysis (PCA) since that technique uses a linear projection of the data and tends to overemphasize the global structure of the data rather than preserve local structures. Tests using PCA yielded groupings that were rather unsatisfactory to how we expected to behave. The next step was to see if using a probabilistic interpretation of pairwise similarities would fare better for our analysis. It is with this we could use t-SNE to construct our embeddings. However, with t-SNE, our resulting 2-dimensional embeddings resulted in clusters that seemed too apparent as the pairwise similarities between points seemed to be magnified. Upon further analysis, it appeared that t-SNE was preserving much of the local structure in the data. We then also embedded our data using DensMAP, a variation of UMAP. Compared to t-SNE, DensMAP seemed to work better in terms of preserving a balance of local and global structure. Despite DensMAP performing relatively slower than t-SNE, we felt that DensMAP provided a more accurate representation of the global structure of the data. Next, we decided to use a nonparametric clustering algorithm like DBSCAN as it allows us to focus on densities instead of assuming that the data follows a Gaussian distribution. The densities computed enable us to create multi-labeled predictions for which cluster the points fall within. This prediction allowed us to extend the functionality of our visualization, as we can now iteratively add more information that can help explain the distribution of our data. We decided to work with about 1000 samples to allow our application to compile and run in a reasonable amount of time.

For pre-processing and cleaning the city_data, we loaded it as a data frame using the pandas library. We obtained our dataset from "<https://www.cityhealthdashboard.com/>". The database had more than 20 metrics about a city but we selected only those features that related to our work on hate crimes and might affect the number of crimes in the city. We have considered only those rows that contained the relevant feature in our data frame. We also dropped all the irrelevant columns and only kept those columns that informed us about the metric of a city like the metric_name and est. There were around 70 rows with null est values; those were dropped as they did not provide useful information.

The feature importance (variable importance) describes which features are relevant and helps us give a better understanding of what factors are affecting the hate crimes in a particular city. For feature importance, we started by pivoting the original 'city_data' dataset in such a way that each unique feature/metric is the column and the cities are the individual columns present in the dataset. The final dataset utilized includes the 'est' values for each city in the metric columns, as well as the state in which the city is located in the state abbr column. After this, we combined this dataset with the hate_crime dataset on the city name and state abbreviation and included the number of crimes in that city as an additional column in the final dataset. This column was used as the target column in the feature importance. The city name and the state abbreviation are further dropped for running the ML algorithm. We have chosen the Random Forest Regressor to fit our model. The X-axis contains all the features and the Y-axis contains the number of crimes. Other algorithms, such as logistic regression, delivered us identical feature importance, thus we proceeded with the random forest regressor.

APP WORKING

The App is divided into the following six sections:

1. Home - Describe the concept of our project, the reasons for choosing this issue, the question we want to address, and the goal we want to achieve through it.

The screenshot shows the app's home page. On the left, there is a sidebar with a title "Hate Crimes in US" and a list of navigation items: "Home" (which is highlighted in blue), "Exploratory Data Analysis", "Hate Crime Distribution", "Clustering", "Feature Importance", and "Exploring States & Cities". The main content area has a title "Home" and a text block about the Pittsburgh synagogue shooting. Below the text is a photograph of a memorial at the Tree of Life Synagogue, featuring several white crosses and stars with names written on them.

Figure 1: Home Page

2. Exploratory Data Analysis - Describing the Data Sources, including their provenance and why it was chosen. We also offer extra visuals for each data set used for data exploration.

For the Hate Crime data set there are two Visualizations:

2.1 Total Number of Crimes Recorded over the last 30 years in the United States represented as a sorted descending bar graph in the order of hate crime count in states.

2.2 Total Recorded Hate Crimes as per Offender's Race, which shows which races are responsible for committing the incident in a decreasing fashion. And a similar visualization for the Victim's Hate Crime Type shows the most attacked racial group in a decreasing manner.

For the Wellness Factors of Cities Dataset:

These include bar graphs for the following relevant wellness factors as per states in the United States:

1. High School Completion
2. Life Expectancy
3. Income Inequality
4. Neighborhood racial/ethnic segregation

5. Racial/Ethnic Diversity

6. Unemployment



Figure 2: Exploratory Data Analysis Page

3. Hate Crime Distribution - Here we get to explore an overall view of the Hate Crime Distribution across the country, time, and race.

3.1 Timeline of Hate Crime Across 30 Years - This line graph shows the ascent and descent of the hate crime rate over the years 1990 to 2020. Furthermore, selecting only the years you are interested in exploring on the scale provided, will display the graph only for those particular years.

3.2 U.S Map - This is an interactive choropleth map that demonstrates the distribution of hate crime statistics across each state of the United States, this map is interlinked to a bar graph that shows the top 15 states. Selecting a state you are interested in, in the US Map will highlight the respective state in the Bar Graph.

3.3 Correlation between the Offender's Race and the Victim's Hate Crime Type - Exploring which offender's race promotes hate crime the most and which victim's group is the most harmed would be an intriguing point of view. This is accomplished by using a heat map to correlate both of these values.

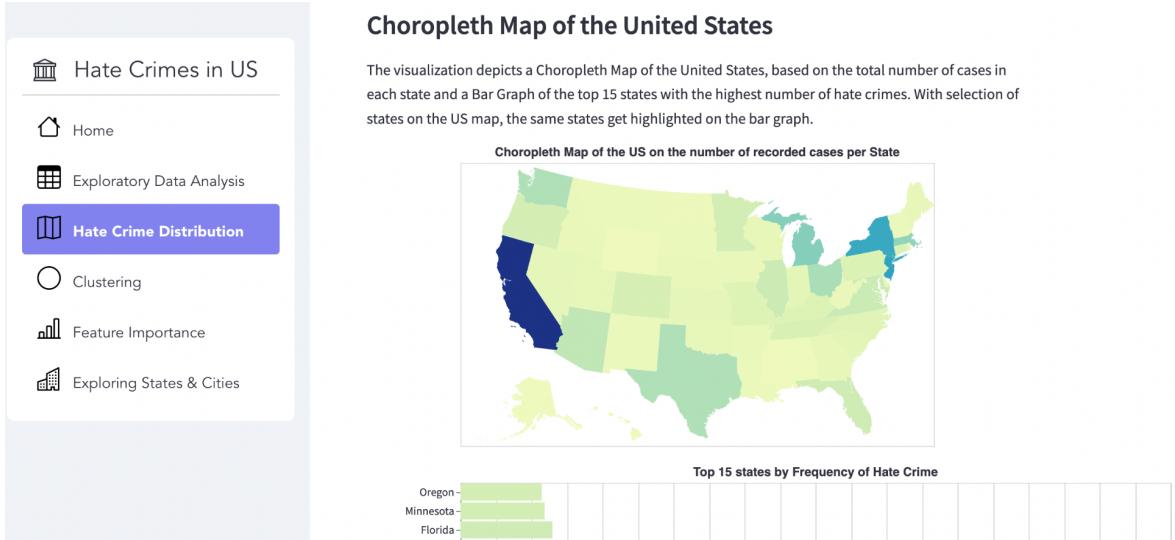


Figure 3: Hate Crime Distribution Page

4. Clustering - In this section, we will be exploring clustering techniques on features such as Bias, Location, Offender Race, Crime, and Victim Type. DBSCAN and DensMAP algorithms were used to reduce the features and cluster the data. The visualizations display clusters based on the user-selected features. The addition of more features creates sparse clusters.

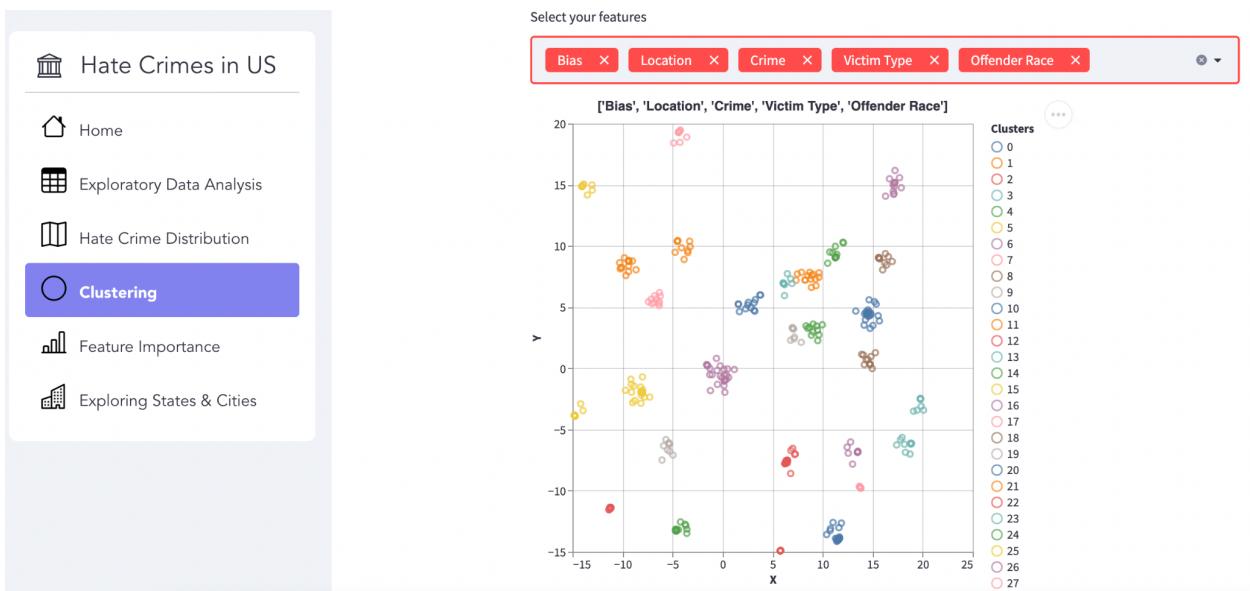


Figure 4: Clustering Page

5. Feature Importance - In this section, we will be exploring the important features that impact hate crimes in US states. The bar graph shows the top features that influence hate crimes in the US. This graph can be customized to be displayed for the states and the feature the user is interested in.

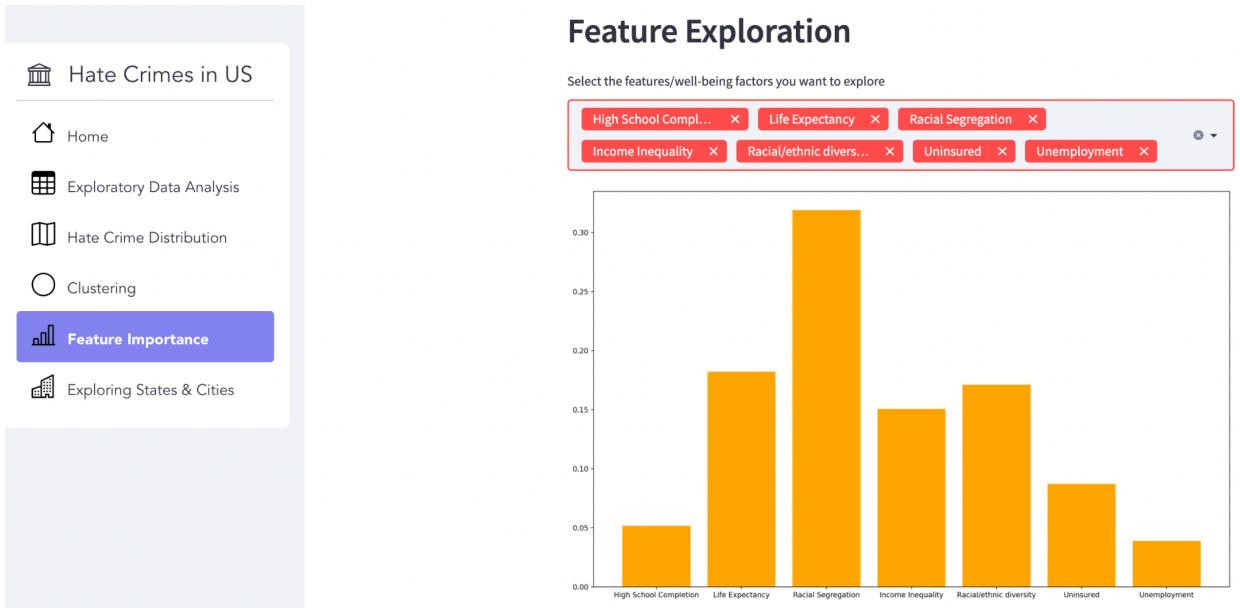


Figure 5: Feature Importance Page

- Exploring States and Cities - In this section, we will be exploring the states of the US. We will first explore the distribution of well-being factors in the state using a pie chart. Next, we will represent the overall hate crime cases in the state over the past 3 decades in the form of a line graph. This is followed by a bar graph that compares the hate crime rates across the various cities in the state. The line graph and bar graph are interlinked. This selection of a particular time period in the line graph displays the crime rates per city for those years. These visualizations can be further explored for the cities as well.

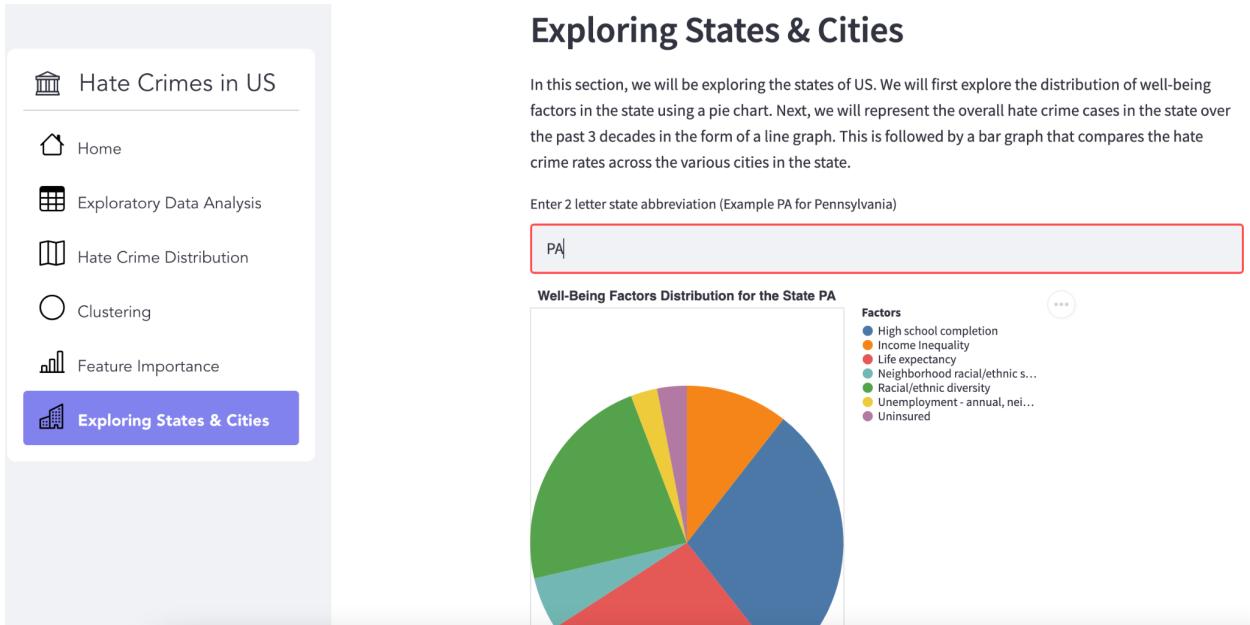


Figure 6: Exploring States and Cities

RESULTS

Through this project, we want to know how socioeconomic factors, particularly human well-being, have influenced hate crime patterns in the United States through the years. To begin, we'll look at the clustering of hate crimes, the relationship between the race of the offender and the race of the victim, and the feature importance of the well-being factors. Then there's the hate crime distribution and timeline across the United States, which leads to a closer look at each state, including aspects like well-being, hate crimes reported over time, and hate crime distribution across cities. Each state's depiction can also be studied further by its cities. These interactive visualizations and the insights we gained from them to answer our question are discussed in the following sections.

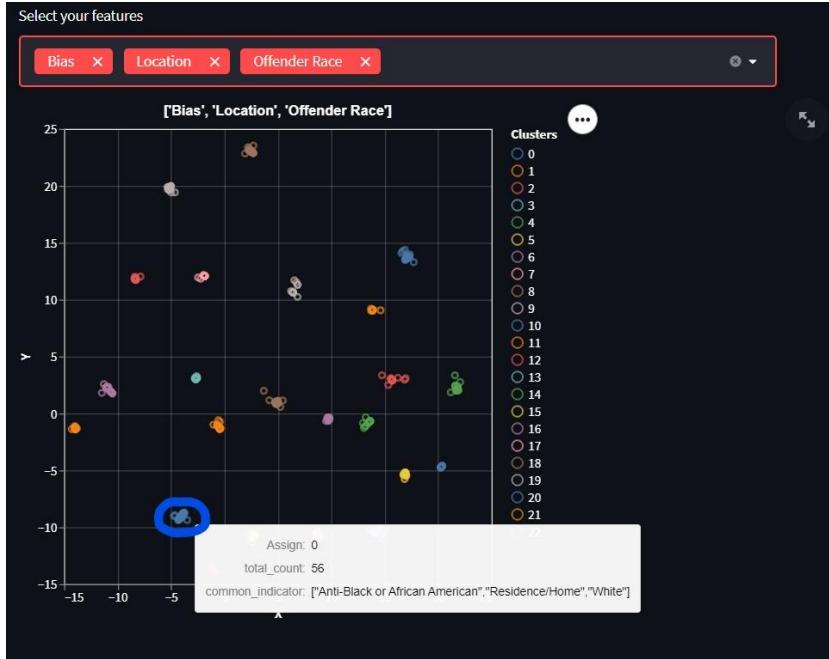


Figure 7: Clustering on Hate Crimes

With the help of our clustering algorithm, we discerned the structure based on the groupings in the hate crimes dataset. Despite using a subsample of 1000 random points, our clustering algorithm was able to produce interpretable clusters that made sense. While each subsequent addition of features creates sparser clusters, we felt it was more informative to look at the clustering with two to three features because it cultivated a better understanding of the global and local structure of the data. Our analysis into the groupings as demonstrated by the 'Figure: Clustering on Hate Crimes' showed that Anti-Black or African American was the major bias expressed in most of these hate crime incidents, with these incidents occurring at Residences/Homes in the form of verbal assault or intimidation. Given that many of these incidents occur at Residential/Homes as intimidation or verbal assault, we ascertain that these may be relatively isolated incidents that do not exhibit overly aggressive violent behavior, nor are they publicly displayed. As a result, many of these victims experience this type of aggression in their own homes and communities. We may even find that these incidents may usually go unnoticed and unreported. From our analysis, we may be underrepresenting the number of crimes like this simply through the limitation of our dataset.

Is there a correlation between the offender's race and the hate crime type of the victim? Answering this question would provide insight into the motivations and psychology of attacks on specific ethnic groups. The Victim's Ethnicity is represented on the x-axis, while the Offender's Race is shown on the y-axis in this heat map. The graph shows that most hate crimes are committed against African Americans, followed by attacks on the LGBTQ+ community, as corroborated by clustering. We also see that the most common offender is white, which may be since white people make up the bulk of the population in the

United States. Crimes against African Americans committed by white offenders are the most serious type of hate crime, according to the below visualization.

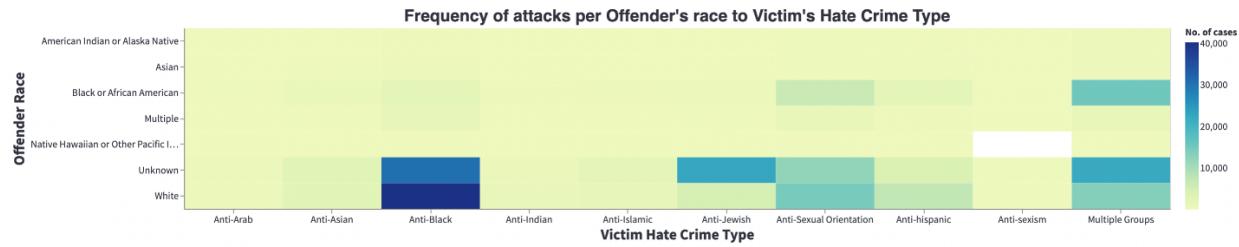


Figure 8: Correlation between Offender's Race and Hate Crime Type

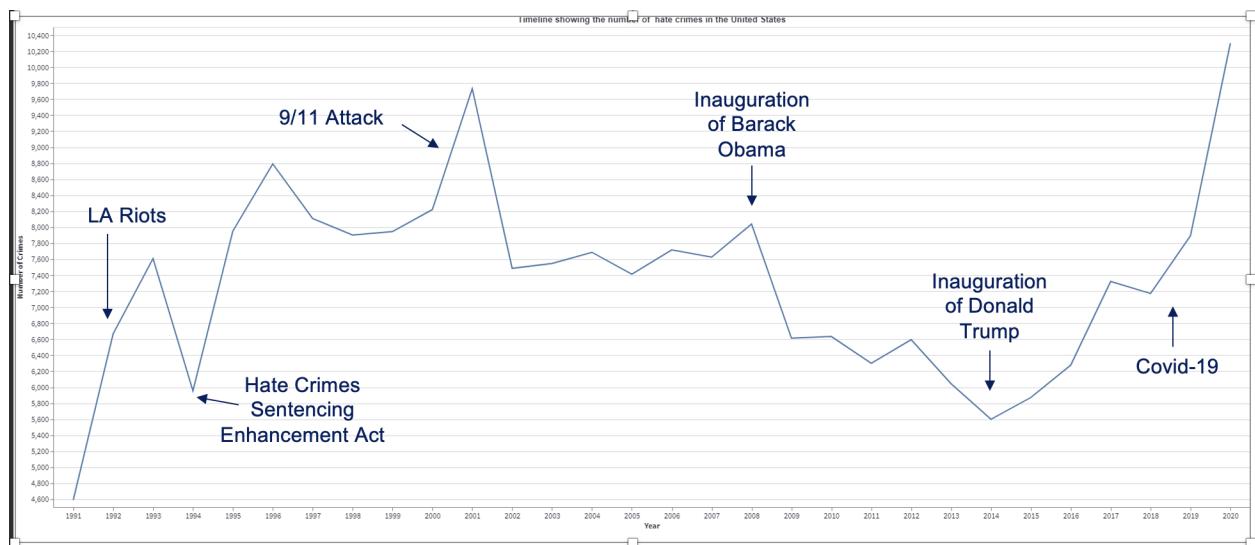


Figure 9: Hate Crime Statistics Over The Past 30 Years

The above visualization depicts the evolution of hate crime statistics in the United States over the last 20 years. The x-axis represents every year from 1991 to 2020, and the y-axis is the total number of cases. The sharp spikes and drops in hate crimes have been linked to major events that occurred in the United States. Some of these have had a good impact on reducing hate crime, while others have dramatically increased hate crime. Riots, presidential elections, policy decisions, and global pandemics are some of these events. The most concerning aspect of this graph is the never-before-seen surge in the number of incidents after the Covid-19 Pandemic began in 2019. This surge opens up the possibility of future research into the extent of the relationship between the occurrence of these significant events and the number of recorded hate crimes. This is discussed in further detail in the upcoming sections.

Choropleth Map of the US on the number of recorded cases per State

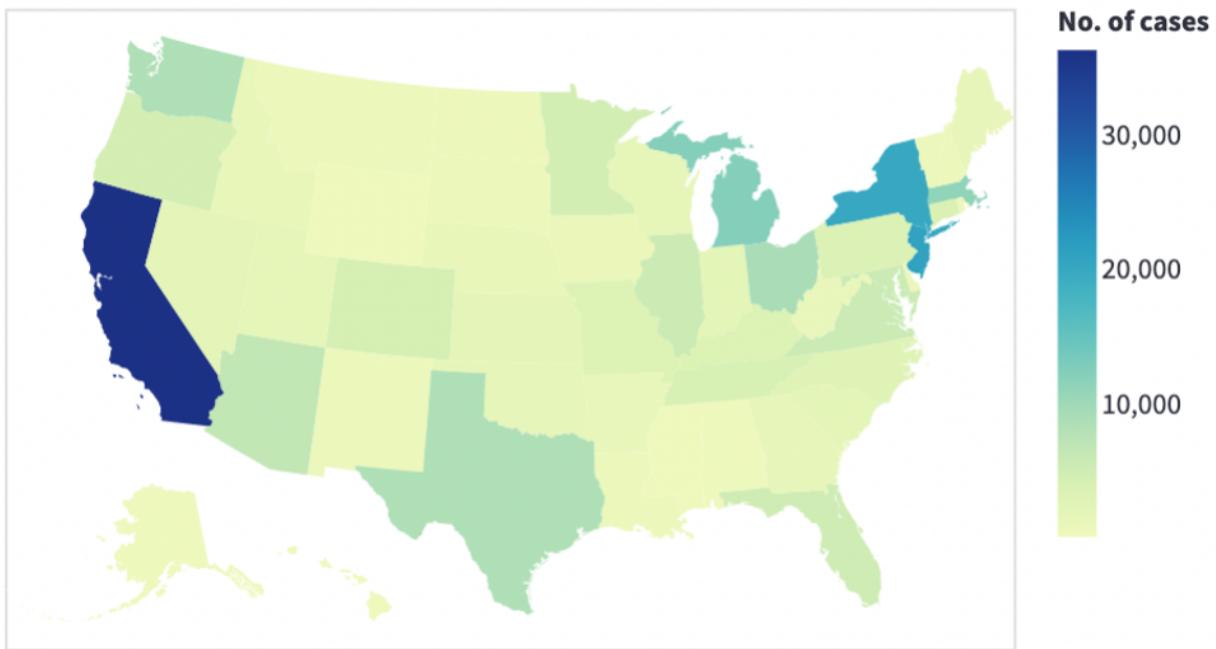


Figure 10: Hate Crime Statistics Over The Past 30 Years

The color variation in the above Choropleth Map of the United States is based on the total number of cases in each state. California has the highest number of hate crime instances. We will look into the numerous socio-economic reasons that have led to California being the most targeted state for hate crimes.

CASE STUDY IN CALIFORNIA

The state of California's well-being factors is depicted in the Pie-Chart below. The following are the factors that were analyzed:

- High School Completion: The completion rate refers to the percentage of pupils who enroll in a high school program and complete it.
- Income Inequality: The degree to which income is distributed unequally throughout a population is income inequality. The more unequal the distribution, the greater the income disparity.
- Life Expectancy: Life expectancy is a statistical measure of how long a person is likely to live based on their birth year, present age, and other demographic parameters such as gender.
- Neighborhood racial/ethnic segregation: The spatial separation of two or more social groups within a given geographic area, such as a municipality, a county, or a metropolitan area, is referred to as residential segregation.
- Racial/ethnic diversity: The recognition and celebration of racial diversity are known as racial diversity. Diversity promotes differences within and between racial identities, recognizing the

intersectionality of various groups such as "ethnicity, gender...age, national origin, religion, handicap, sexual orientation, financial level, education, marital status, language, and physical appearance."

- Unemployment: Unemployment is a circumstance in which a person who is actively looking for a job is unable to find work. Unemployment is a significant indicator of the economy's health.
- Uninsured: The population that does not have insurance.

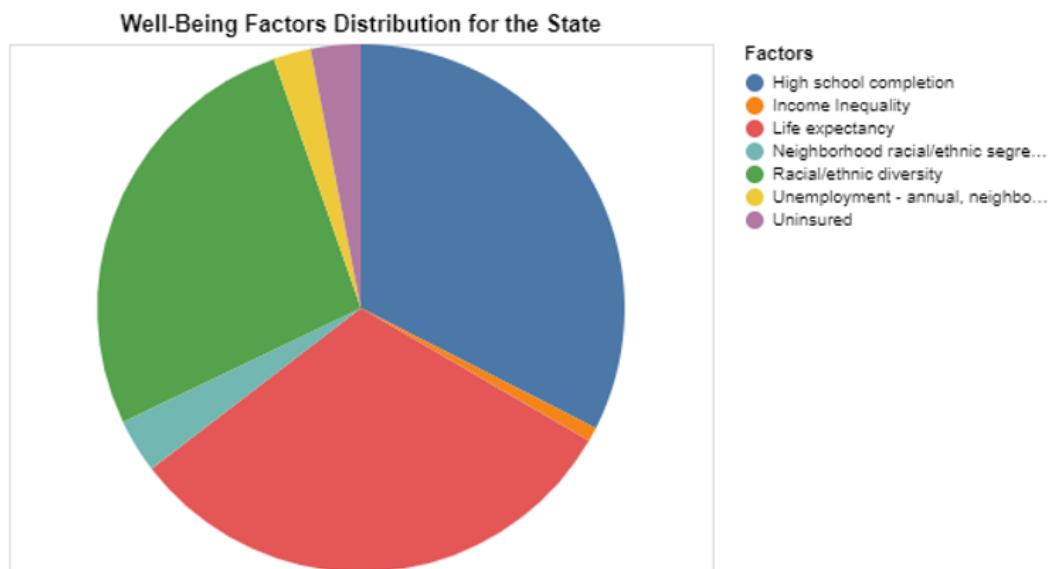


Figure 11: Well Being Factors Distribution in California

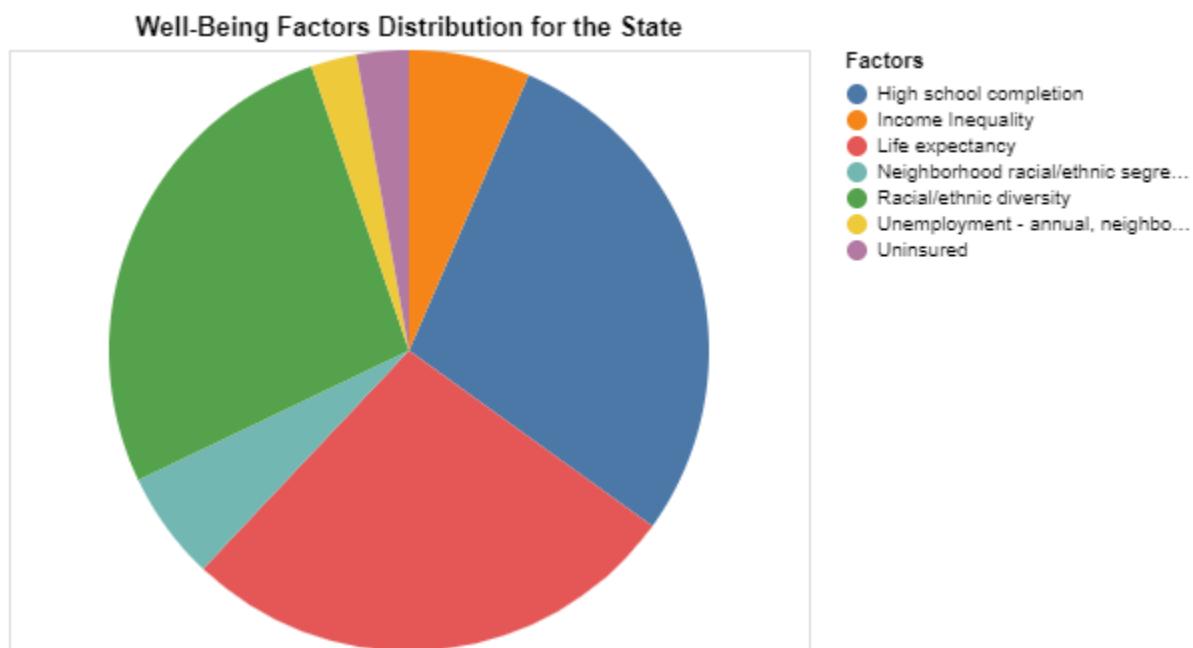


Figure 12: Well Being Factors Distribution in New York

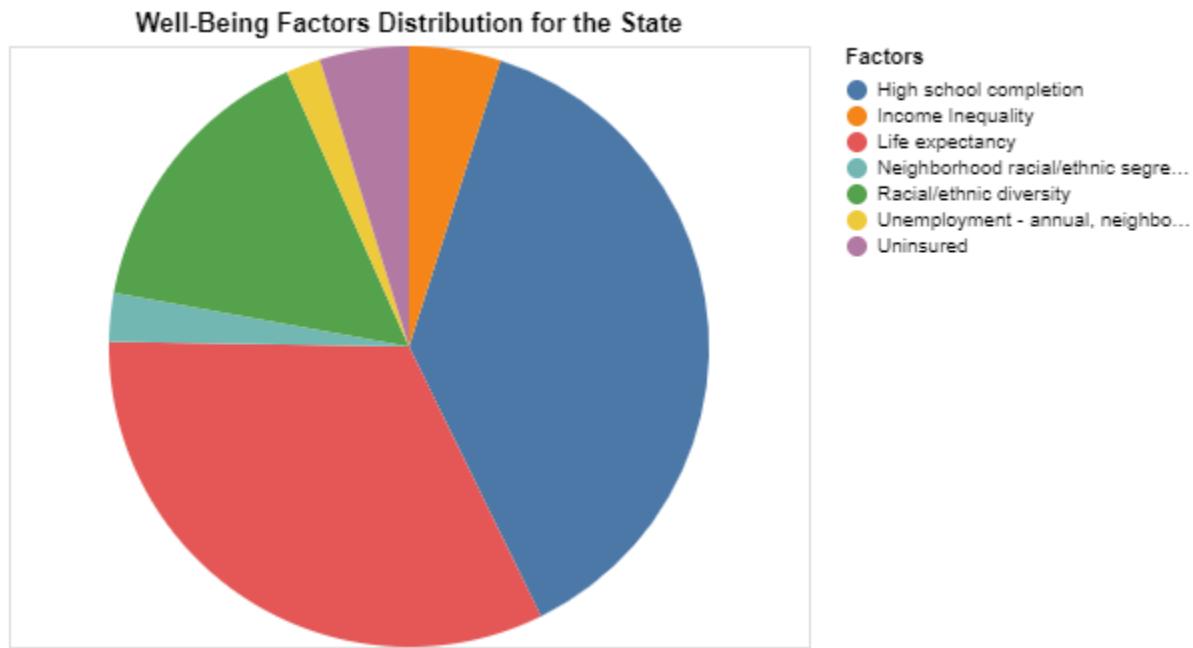


Figure 13: Well Being Factors Distribution in Wyoming

While comparing crime across California, New York, and Wyoming, it is clear that racial/ ethnic diversity and neighborhood racial/ethnic segregation is the differentiating factor. So, with an increase in overall diversity within a state, the rate of hate crimes increases. Hence, states like California and New York with higher racial diversity have more crimes than Wyoming with a lesser racial diversity.

Second, we look for patterns in the variations in the number of cases in California from 1991 to 2020. As a result, we may deduce that after 1997, the years 2001 and 2002 had the highest number of instances, while 2014 had the lowest. This rise after 2014 can be related to the outbreak of the covid-19 epidemic, which resulted in racial discrimination against Asians.

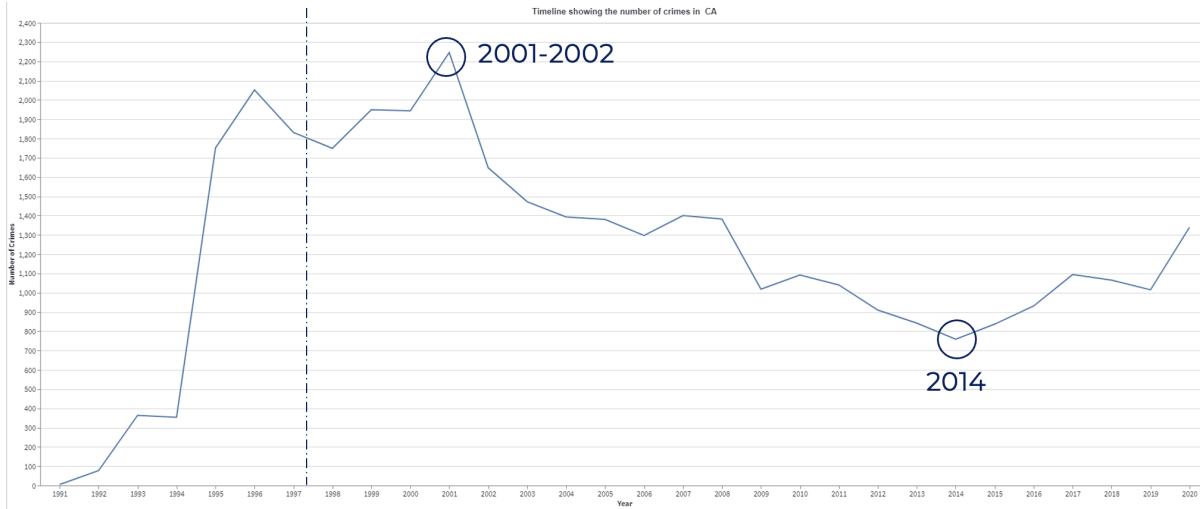


Figure 14: Hate Crime Statistics Over Time in California

Lastly, we created a sorted bar graph of the distribution of the total number of cases across the top 10 cities in California to see how the hate crime rate varies across the state. In contrast to other cities, more urban locations such as Los Angeles, San Francisco, and San Diego are predicted to have a significantly high number of incidents as a result of this research. This is related to the feature importance, which stated that the more ethnic diversity, segregation, and life expectancy there are, the greater the hate crime rate.

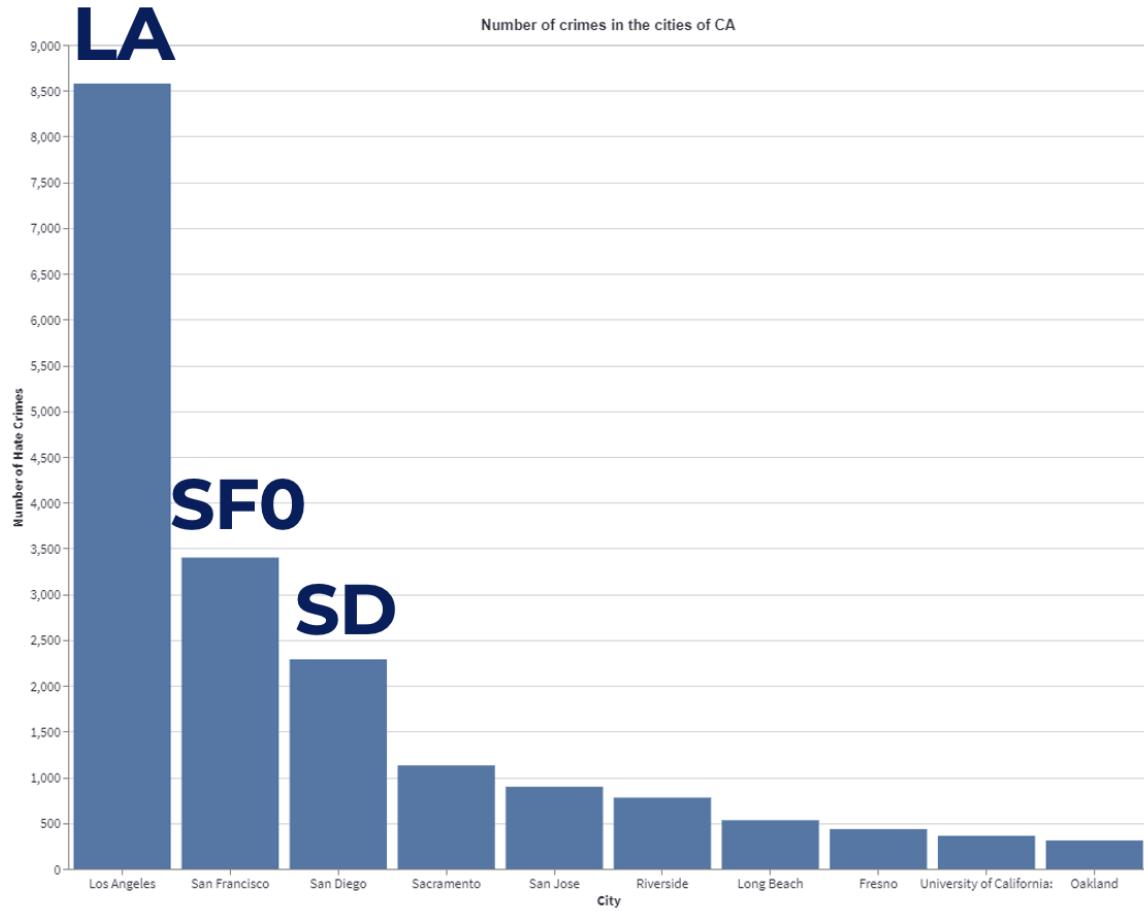


Figure 15: Exploring Cities in California

Hate crimes against African Americans are remained frequent across the United States, closely followed by anti-LGBTQ+ hate crimes, according to the clustering and correlation. In addition, while there has been a steady decline in hate crime in recent years, the pandemic has once again pushed the surge in incidents to a dangerously steep level. California has the highest number of instances in the United States, owing to its diverse racial diversity and segregation throughout its neighborhoods. Furthermore, the majority of cases in each state occur in the most populous metropolitan areas. As a result, we might deduce that hate crimes are more likely in urban regions with a more diversified population.

Discussion

By analyzing hate crimes over the past three decades, it is evident that socio-economic factors play a significant role in the number of crimes committed per region. High racial/ ethnic diversities often lead to a very diverse crowd resulting in social tensions. This increase leads to increased cases of threats and intimidation against minority communities. Within the United States of America, the white community forms a majority. As a result, the prominent race of the offender is white.

Through clustering techniques, it was clear that intimidation and verbal assault are more prominent than physical assault or harm of any form. The cause of this is verbal disputes majorly happen at home and within neighborhoods where physical harm is not a preferred mode of solving disputes.

Another insight gained from our project was that important incidents were correlated to crime rates. The inauguration of President Obama significantly reduced crime rates in the US. On the other hand, the COVID-19 pandemic resulted in a surge in hate crimes. Anti-sexual orientation crimes have increased over the past few years. This increase could be a result of the LGBTQ Acts that came into place, causing a lot of social disturbance based on the values of society.

One of the most interesting insights gained from this project was that despite movements such as black lives matter that tried to create awareness of the African American community, the crimes against this community have been a majority. This majority shows that social awareness is not enough in reducing crimes against this community. Rather, stricter enforcements need to be placed. Another insight gained from the project is that the race of the majority of offenders is unknown. This could be because the offenders have not been caught, the victim is unable to identify the race of the offender, or the offender's race was not reported with the crime.

One of our observations while implementing the project was that in the clustering, adding more features resulted in more dispersion. Clusters were not as dense as when the number of considered features was less. Another observation during the feature importance analysis was that not all features were significant to hate crime statistics contribution. Unemployment and high school graduation have little impact on hate crime, which is surprising given that education and employment are two of society's most crucial elements.

Future Work

Our application has much potential for extensive research and study into the topic of hate crimes. The application we developed gives a relatively high-level overview of the trends in hate crime patterns and the factors most correlated with those patterns. As a result, we hope that our application can be used to direct research to understand why hate crimes occur and what are the plausible solutions for it.

An area of research that can be further enhanced is to explore the psychology behind hate crimes and how different factors and reasonings form together based on aggression and bias. With this, first-hand accounts and personal testimony can help us better understand why and how people form their biases. By better understanding offenders and victims, we can develop solutions that challenge ingrained thought processes. As discussed before, hate crimes are multidimensional issues that are affected by several socioeconomic factors brought on by the increasing density of diverse people groups. Therefore, we can expect that answers may be very diverse and that hate crime may be driven by a myriad of issues.

Another extension for this application is to look at effective solutions that go beyond what factors are most correlated with hate crime patterns. A simple solution would be to look at education and conclude that education is necessary to combat hate crime. However, from our analysis, hate crimes typically happen in areas where the average level of education is very high. Therefore, our solutions must confront notions in a way that creates real experiences that lie contrary to their own biases. An idea that came up in our discussions was to promote and encourage cultural exchanges so that communities can learn from one another. It forces a level of uneasiness that comes with operating beyond a comfort zone. Yet, we believe that this is an effective strategy to help confront hate-based crimes and aggression. We hope that our application can help explore additional strategies and develop effective programs and policies that can help address tensions in communities.

Overall, we expect that our application gives a good starting point in further exploration of the reasonings and motivations behind hate crimes. We hope that further research can be applied in this study and that the tools and algorithms developed here can be used to further understanding. Furthermore, our team hopes to continue to extend our work by examining and analyzing hate crimes in other countries, applying the same methodology we have here within another country's context. The hope is to compile a larger dashboard of how hate crimes behave in other countries and compare across the world based on similar features and characteristics.

Additionally, our application can be refined by capturing the effect of other demographic metrics or simply with a larger dataset of hate crime incidents. A methodology that we wanted to apply, but didn't find appropriate was to run a classification-based predictive model to determine if someone was more likely to commit a hate crime or be a victim of a hate crime. While we do acknowledge the ethical concerns, we explored the possibility that this task can be applied holistically to cities or regions to determine their susceptibility to a high rate of hate crimes. We aim to do so by looking at public policies, laws, acts, and other federal and state actions that can contribute to a high or lower rate of hate crimes. To do so, we will need additional data on these policies and we need to research additional algorithms and methodologies to analyze these types of data. We believe that this can enhance the insights and learnings of our current application and further our understanding of the characteristics of hate crimes. In the future, we hope that we can effectively combat hate crimes and biases and foster loving and supportive communities everywhere.

References

- [1] EHRC. (2016). Research report 102: Causes and motivations of hate crime. <https://www.equalityhumanrights.com/sites/default/files/research-report-102-causes-and-motivations-of-hate-crime.pdf>
- [2] World Population Review. (2021). Safest Countries in the World 2020. Worldpopulationreview.com. <https://worldpopulationreview.com/country-rankings/safest-countries-in-the-world>