# Socioeconomic Determinants of Hate Crimes in the United States*

## A Bayesian Hierarchical Analysis

Yuxuan Wei

November 29, 2024

This study provides the relationship between socioeconomic factors and hate crimes across U.S. states using a Bayesian hierarchical framework. We analyze a dataset capturing hate crime incidence alongside state-level socioeconomic indicators such as median income, unemployment rates, and income inequality. Results reveal that higher income inequality and unemployment rates are associated with elevated hate crime rates, while higher median income is inversely related. These findings emphasize the need for targeted socioeconomic interventions to mitigate hate crime prevalence.

## 1 Introduction

Hate crimes remain a significant social issue in the United States, threatening community safety and equality. These crimes disproportionately affect marginalized groups and undermine social trust. Existing research identifies socioeconomic factors—such as income inequality, unemployment, and education—as potential contributors to hate crime prevalence. However, most studies overlook state-specific variations in these relationships, limiting the understanding of how regional disparities influence hate crime patterns. This paper uses data sourced from FiveThirtyEight (2017), which provides detailed state-level hate crime statistics alongside socioeconomic indicators, to explore these relationships rigorously.

This paper investigates the effect of state-level socioeconomic factors on hate crime rates. The estimand is the average marginal effect of variables, including income inequality, unemployment, and education levels, on hate crimes per 100,000 population. A Bayesian hierarchical model is employed to account for regional heterogeneity, allowing for more precise estimation of these effects.

---

*Code and data are available at: [https://github.com/wyx827/Hate_Crimes_Socioeconomic_Analysis].

The results demonstrate that states with higher income inequality and unemployment rates report more hate crimes, while higher median income and education levels are associated with reduced hate crime prevalence. These findings emphasize the socioeconomic disparities that contribute to hate crime patterns, suggesting the importance of addressing inequality and unemployment at the state level.

Understanding the socioeconomic determinants of hate crimes is essential for informing policies aimed at reducing their prevalence. This analysis provides evidence for the role of structural inequality and economic conditions in shaping hate crime patterns, offering a foundation for targeted policy interventions to foster safer and more equitable communities.

The remainder of this paper is structured as follows. Section 2 describes the dataset, variables, and preprocessing steps. Section 3 outlines the methodological framework, including the Bayesian hierarchical model. Section 4 presents the results, accompanied by tables and visualizations. Section 5 discusses the implications, limitations, and potential future directions of the analysis. Section A concludes with a summary of the key findings and recommendations.

## 2 Data

### 2.1 Overview

This study employs R (R Core Team 2023) for cleaning and analyzing the dataset, leveraging libraries such as `dplyr` (Wickham et al. 2019), `ggplot2` (Wickham 2016), `corrplot` (Wei and Simko 2024), and `arrow` (Foundation 2023). The dataset, sourced from FiveThirtyEight (FiveThirtyEight 2017), investigates the relationships between socioeconomic variables and hate crime prevalence across the United States. It aggregates data for all 50 states and the District of Columbia.

After preprocessing, which involved removing missing values, standardizing variables, and ensuring data consistency, the final dataset includes **51** observations for **7** predictors: `median_income`, `unemployment_rate`, `metro_area_population_share`, `high_school_education_share`, `gini_index`, and `non_white_population_share`. The outcome variable, `hate_crimes_per_100k_splc`, represents hate crime prevalence reported by the SPLC. Together, these variables form the foundation for examining socioeconomic influences on hate crime rates.

---

## 2.2 Measurement and Considerations

Understanding how socioeconomic conditions translate into measurable data is vital for interpreting the results. Below, we discuss the key measurement aspects and their implications.

- **Phenomena Representation**: Hate crime incidents, inherently complex societal phenomena, are represented in the dataset as `hate_crimes_per_100k_splc`, which quantifies the number of hate crimes reported per 100,000 residents in each state. This measure was derived from state-level aggregates reported by the SPLC, ensuring consistency across regions.

- **Data Collection Context**: Socioeconomic variables such as income, unemployment, and education are proxies for broader economic and social conditions. These data were collected from credible government sources and harmonized by FiveThirtyEight. The inclusion of standardized metrics, such as `gini_index` (inequality measure) and `metro_area_population_share`, ensures comparability across states despite regional disparities.

- **Challenges in Measurement**: Hate crime data, while critical, is subject to biases in underreporting and differences in enforcement. States with higher transparency or stricter enforcement may report higher rates, not necessarily indicating a greater prevalence. Similarly, socioeconomic indicators such as `non_white_population_share` could interact with other unmeasured factors, complicating causal interpretations.

- **Validation and Limitations**: Standardization ensures data comparability, but reliance on a single year of data limits trend analysis. This temporal snapshot captures static relationships, not dynamics over time. Additionally, state-level aggregation masks within-state heterogeneity, such as urban vs. rural differences in socioeconomic conditions and hate crime reporting.

---

## 2.3 Variables of Interest

This section introduces the dataset's key variables, including detailed descriptions, visualizations, and insights into their relevance. Visualizations illustrate distributions, relationships, and variability to contextualize their potential influence on hate crime rates.

### 2.3.1 Outcome Variable: Hate Crimes Per 100k SPLC

The outcome variable, representing the number of hate crimes per 100,000 residents as reported by the SPLC, its distribution, shown in Figure 1, exhibits a skewed distribution. Most states report hate crime rates between 0.1 and 0.5 per 100,000 residents. However, the presence of an outlier with a much higher rate highlights significant regional disparities. This variation underscores the importance of investigating contributing factors, such as socioeconomic, demographic, and policy-related influences, to understand the drivers of hate crimes.

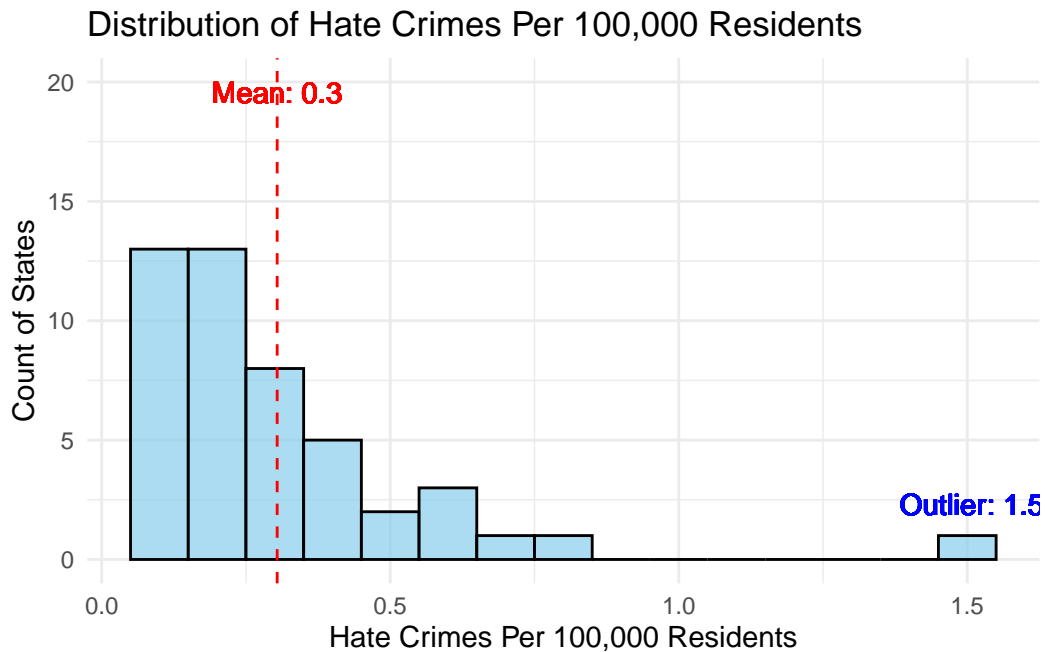## Distribution of Hate Crimes Per 100,000 Residents



Figure 1: Distribution of hate crimes per 100,000 residents, as reported by the SPLC, illustrating state-level variation. Most states report between 0.1 and 0.5 hate crimes per 100,000 residents, with a notable outlier indicating a disproportionately higher hate crime rate.

## 2.4 Predictor variables

### 2.4.1 Median Income

Median income represents the annual income of the median household in each state, serving as a proxy for economic prosperity. As shown in Figure 2, shows the distribution of median household income across states. There is a relatively normal distribution with a mean of $54,802 and a median of $54,310. The narrow spread suggests a clustering of household incomes around this value, with limited variation across states. This clustering might indicate

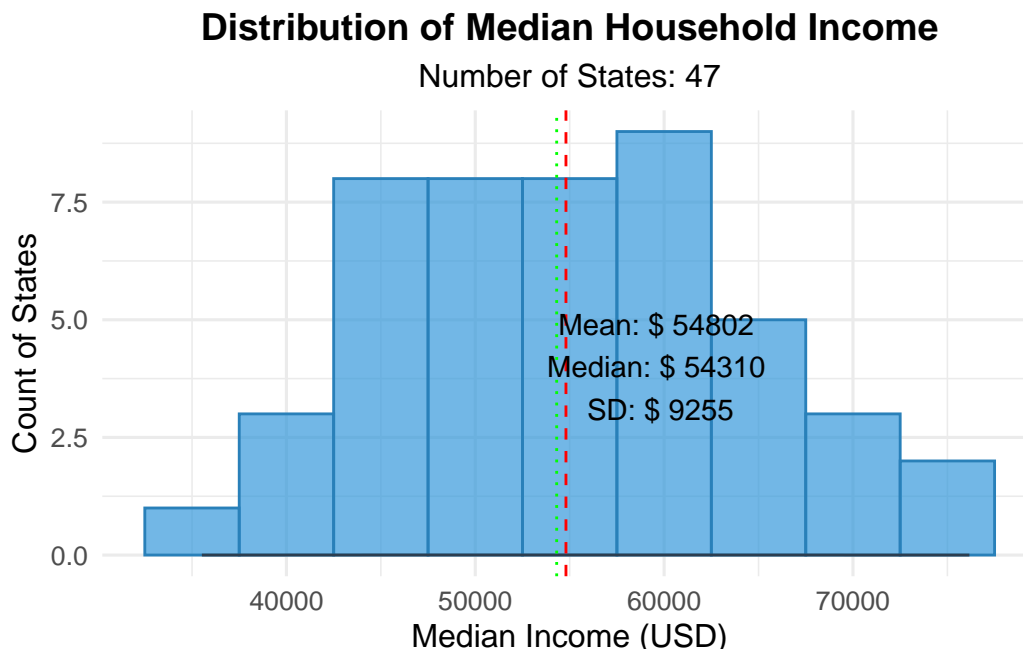consistency in income policies or similar economic opportunities in various regions of the United States.

## Distribution of Median Household Income

### Number of States: 47



Figure 2: Distribution of median household income across states, illustrating the variation in income levels. The mean household income is $54,802, while the median is slightly lower at $54,310, indicating a slightly right-skewed distribution. The standard deviation of $9,255 highlights variability across states.

### 2.4.2 Unemployment Rate

The unemployment rate measures the percentage of individuals unemployed on a seasonal basis. As shown in Figure 3, it shows an extremely narrow spread, with the mean and median unemployment rates both around 0.05%. This low variability suggests that unemployment rates are consistently low across states. The majority of states fall within a very small range of unemployment, which implies effective employment policies or consistent labor market dynamics nationwide.

### 2.4.3 Metro Area Population Share

Metro area population share reflects the proportion of a state's population residing in metropolitan areas, which means states with higher urbanization levels may experience unique social dynamics influencing hate crime prevalence. As shown in Figure 4, it shares reveals that the majority of states have a population share in metro areas concentrated
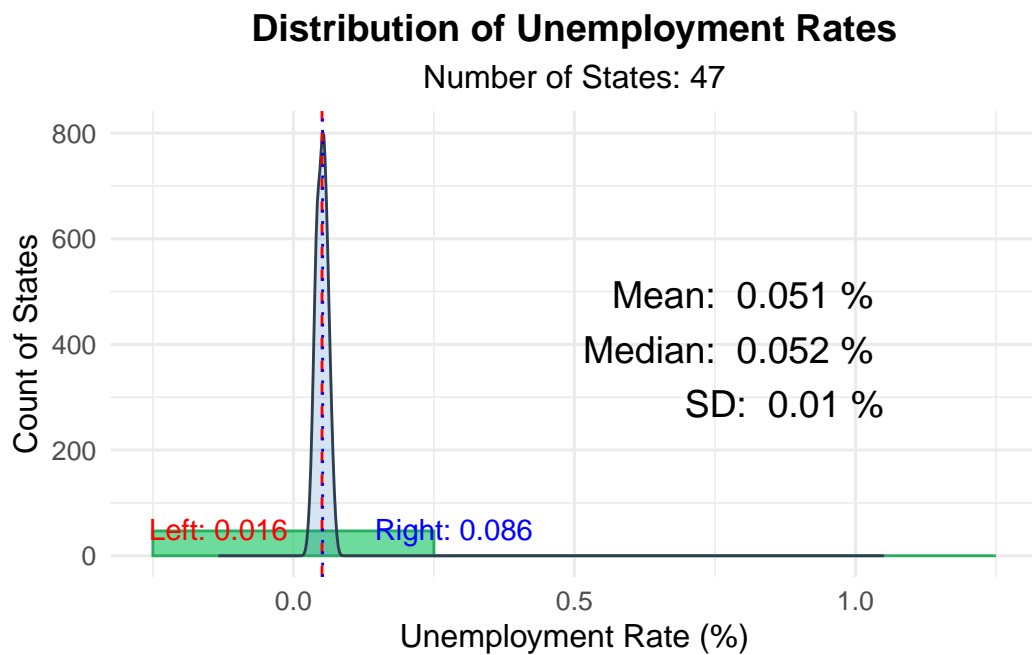
Figure 3: Distribution of unemployment rates among states. The unemployment rate is centered around 0.05%, showing very low variability among states, suggesting uniform employment opportunities across regions.

around 0.77, with a standard deviation of 0.2. This suggests that many states have a high proportion of their populations in metropolitan areas, reflecting the urban-centric population distribution and possibly highlighting urbanization trends in the U.S.



**Metro Area Population Share Distribution**
Number of States: 47

Mean: 0.77 %
Median: 0.8 %
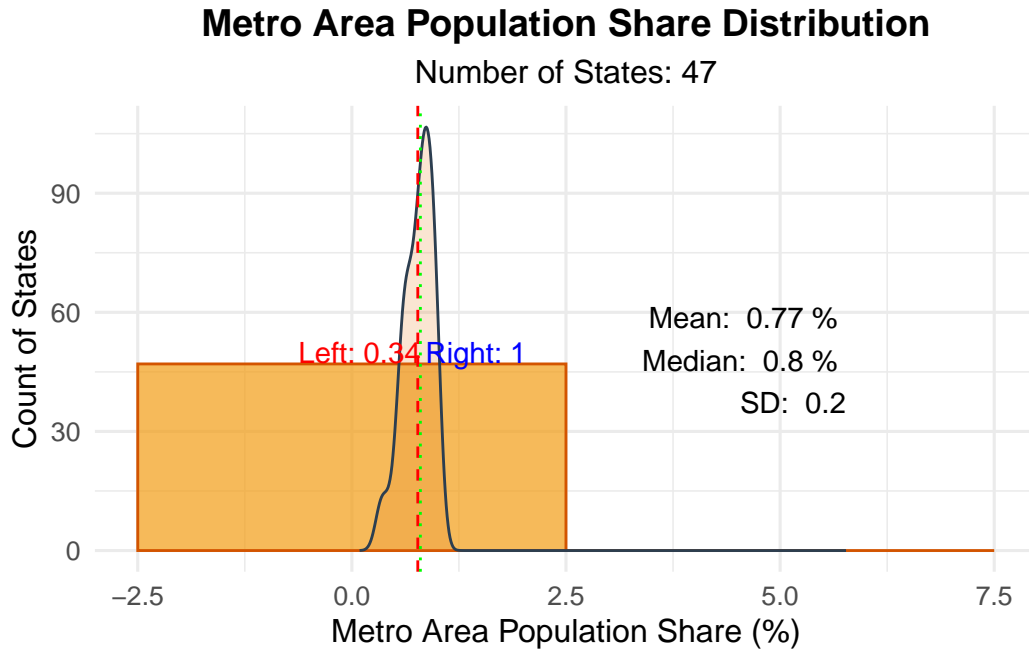SD: 0.2

Left: 0.34 Right: 1

Figure 4: Distribution of metro area population share for states, showing the extent to which populations are concentrated in metropolitan areas. Most states have metro area shares clustered around 0.77, indicating a high proportion of urbanized population.

### 2.4.4 Gini Index

The Gini Index measures income inequality, with higher values indicating greater disparity. As shown in Figure 5, it is centered around a mean of 0.456 and a median of 0.455. The relatively narrow spread in Gini values across states suggests that income inequality is consistent across the nation, with most states experiencing moderate levels of income disparity. This consistency may reflect similar economic structures or social policies.

### 2.4.5 Non-White Population Share

This variable represents the proportion of non-white residents in each state. Higher non-white population shares may correlate with higher hate crime rates due to increased exposure to discriminatory behavior. As shown in Figure 6, this distribution of non-white population share indicates that most states have around 30% non-white residents. This suggests a considerable level of ethnic diversity across the nation. The mean of 0.32 and median of 0.3 show a moderate
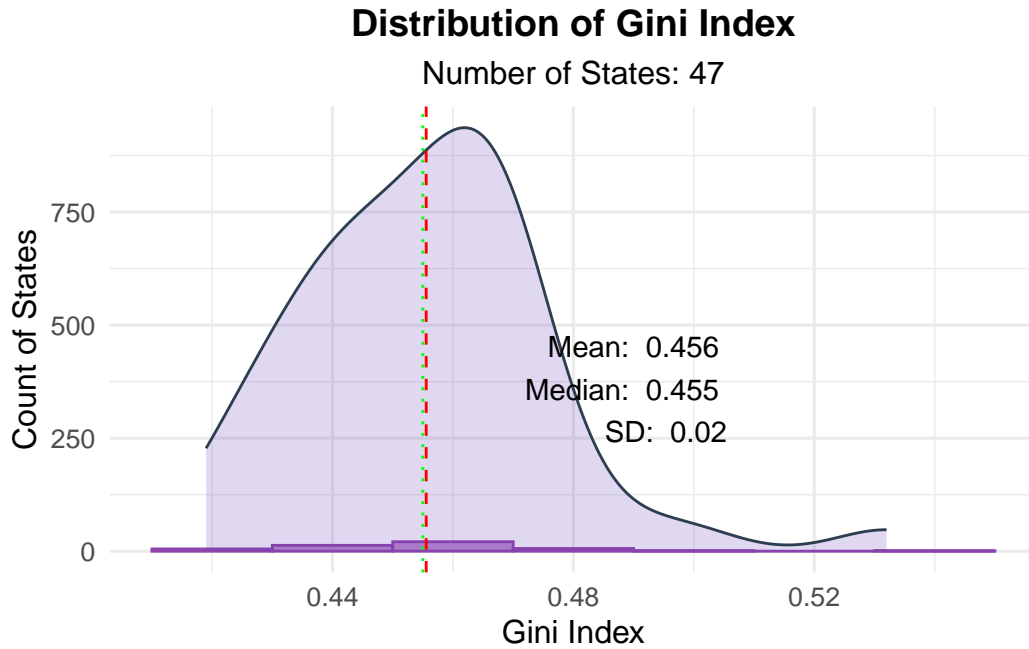
Figure 5: Distribution of Gini Index values across states, indicating income inequality. Most states exhibit a Gini Index value around 0.455, suggesting moderate levels of income inequality.

concentration of non-white population, reflecting varying levels of demographic diversity in different states.

## 2.5 Summary Statistics and Variable Relationships

The relationships between predictors and hate crimes were explored using correlation matrices, summarized in Figure 7. The correlation heatmap summarizes the relationships between various variables, providing insight into potential predictors of hate crimes per 100k SPLC. Notable correlations include the positive relationship between high school education share and metro area population share (r = 0.57), and between Gini Index and unemployment rate (r = 0.49). The relationship between hate crimes per 100k SPLC and non-white population share is weak, suggesting limited direct influence of ethnic diversity on hate crime prevalence at the state level.

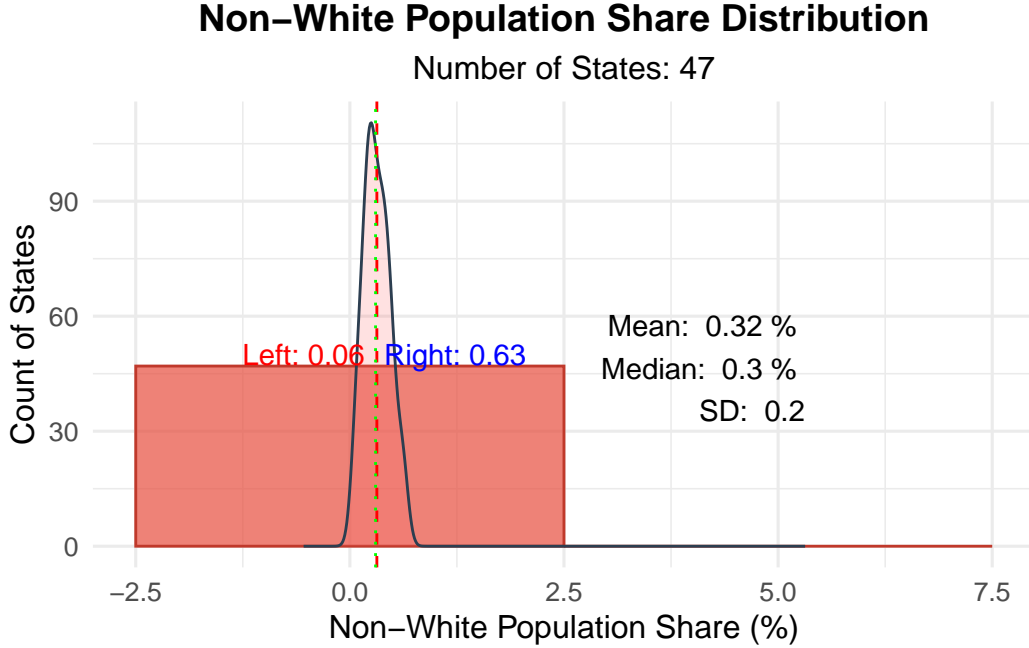## Non–White Population Share Distribution

Number of States: 47



Figure 6: Distribution of non-white population shares across states. The majority of states have a non-white population share around 0.3, reflecting regional diversity levels.

## 3 Model

### 3.1 Model Selection

To analyze the relationship between socioeconomic factors and hate crimes, a **Bayesian hierarchical model** was chosen. This model accounts for state-level heterogeneity by including random effects for states, allowing for variability in hate crimes that might be influenced by unobserved state-specific factors. The Bayesian framework also allows incorporating prior knowledge and managing uncertainty effectively.

Alternative models were considered, including:

1. **Logistic Regression**: This model was used for scaling hate crimes within the range (0, 1). However, logistic regression is primarily suited for binary or bounded outcomes and may not capture the full variability of continuous predictors in this dataset.

2. **Simple Linear Regression**: While linear regression can estimate the relationship between hate crimes and individual predictors, it does not accommodate hierarchical structures or address multilevel variability, making it less robust for this dataset.

9

Figure 7: Correlation heatmap of predictor variables, showing relationships among median household income, unemployment rate, metro area population share, high school education share, Gini index, and hate crimes per 100k SPLC. (income = median_income, unemployment = unemployment_rate, population = metro_area_population_share, education = high_school_education_share, Gini = gini_index, non_white_population = non_white_population_share, hate_crimes = hate_crimes_per_100k_splc, avg_hatecrimes = avg_hatecrimes_per_100k_fbi)

The **Bayesian hierarchical model** was selected as it provides the flexibility to model the complex relationships between hate crimes and predictors while considering the state-level effects.

Model's diagnostics, alternative models considered, and full regression output, ensuring transparency and replicability of the analysis. are included in Appendix C.

## 3.2 Bayesian Hierarchical Model Overview

The model predicts hate crimes per 100,000 residents (`hate_crimes_per_100k_splc`) using the following predictors:

- **Median Household Income (`median_income`)**: Captures the economic status of residents within each state.
- **Unemployment Rate (`unemployment_rate`)**: Measures the proportion of the population without jobs.
- **Metro Area Population Share (`metro_area_population_share`)**: Reflects the proportion of residents living in metropolitan areas.
- **High School Education Share (`high_school_education_share`)**: Indicates the proportion of residents with at least a high school diploma.
- **Gini Index (`gini_index`)**: Measures income inequality within each state.
- **Non-White Population Share (`non_white_population_share`)**: Reflects the demographic diversity within each state.
- **State (`state`)**: Included as a random effect to capture state-specific variability.

The model takes the form:

$$[ \text{Hate Crimes}\_i = \_0 + \_1 \ \text{Income}\_i + \_2 \ \text{Unemployment}\_i + \_3 \ \text{Population}\_i + \_4 \ \text{Education}\_i + \_5 \ \text{Gini}\_i + \_6 \ \text{NonWhite}\_i + u\_i + \_i ]$$

Where:

- ( \_0) is the intercept term.
- ( \_1, \_2, …, \_6) are coefficients for the fixed effects (predictors).
- (u\_i) represents the random effects for state-level variations.
- ( \_i Normal(0, ^2)) is the residual error term.

### 3.2.1 Assumptions and Priors

Key assumptions for the Bayesian hierarchical model include:

1. **Linearity**: The relationship between predictors and the outcome is linear.
2. **Homoscedasticity**: The residuals have constant variance.
3. **Random Effects**: State-level random effects capture unobserved variability.

4. **Normality**: Errors are normally distributed.
5. **Priors**: Normal priors with mean 0 and standard deviation 2.5 are used for coefficients and intercepts, allowing for flexibility while avoiding extreme values.

The prior distributions are as follows:

$$\beta_j \sim \text{Normal}(0, 2.5), \quad \forall j \tag{1}$$

## 3.3 Interpretation of Coefficients

The Bayesian hierarchical model estimates the relationship between socioeconomic predictors and hate crimes per 100,000 residents. Below is the interpretation of each coefficient:

### 3.3.1 Fixed Effects

- **Intercept (($\beta_0$))**: The intercept represents the baseline predicted hate crimes per 100,000 residents when all predictors (median income, unemployment rate, metro area population share, high school education share, Gini index, and non-white population share) are at their mean or baseline values. This serves as the reference point for interpreting the effects of the predictors.

- **Median Household Income (($\beta_1$))**: This coefficient indicates how a \$1,000 increase in median household income affects the predicted hate crimes per 100,000 residents. A negative coefficient suggests that higher income is associated with fewer hate crimes, potentially reflecting improved economic stability and reduced socioeconomic tension.

- **Unemployment Rate (($\beta_2$))**: This coefficient captures the effect of a 1% increase in the unemployment rate on the predicted hate crimes per 100,000 residents. A positive coefficient implies that higher unemployment is associated with more hate crimes, possibly due to heightened social stress and economic insecurity.

- **Metro Area Population Share (($\beta_3$))**: This coefficient represents the effect of a 1% increase in the proportion of residents living in metropolitan areas on the predicted hate crimes. A positive coefficient could indicate that urban areas, with their higher population density, may experience more hate crimes due to greater diversity and interpersonal interactions.

- **High School Education Share (($\beta_4$))**: This coefficient reflects how a 1% increase in the proportion of residents with at least a high school diploma impacts the predicted hate crimes. A negative coefficient suggests that higher education levels are associated with fewer hate crimes, highlighting the role of education in promoting tolerance and reducing prejudice.

- **Gini Index (( $\beta_5$ ))**: The Gini index measures income inequality. This coefficient shows how a 0.01 increase in the Gini index affects hate crimes. A positive coefficient indicates that higher income inequality correlates with more hate crimes, suggesting that economic disparities may contribute to social tensions.

- **Non-White Population Share (( $\beta_6$ ))**: This coefficient captures the effect of a 1% increase in the non-white population share on hate crimes. A positive coefficient may reflect increased targeting of minority groups in states with greater diversity.

### 3.3.2 Random Effects

- **State-Level Effects (($u_i$))**: The random effects account for state-specific variability in hate crime rates that cannot be explained by the predictors. These effects capture unobserved factors unique to each state, such as cultural, political, or enforcement differences.

### 3.3.3 Residual Variance (( $\sigma^2$ ))

The residual variance represents the variability in hate crimes not accounted for by the model. A smaller residual variance indicates that the model explains more of the variability in hate crimes.

---

## 3.4 Example Interpretation

If the coefficient for **unemployment rate** (( $\beta_2$ )) is (0.25), then a 1% increase in unemployment rate would lead to an increase of (0.25) hate crimes per 100,000 residents, holding all other predictors constant. Similarly, if the coefficient for **median household income** (( $\beta_1$ )) is (-0.05), then a $1,000 increase in income would reduce hate crimes by (0.05) per 100,000 residents.

## 3.5 Implications

The coefficients provide insights into how socioeconomic factors impact hate crimes. For example:

- Policies targeting economic inequality (reducing the Gini index) and increasing educational attainment may help reduce hate crimes.

- Urbanization and population diversity may require targeted interventions to address potential tensions.

These interpretations guide the understanding of how predictors influence hate crimes and inform potential policy recommendations.

## 3.6 Model Justification

The following are interpretations of the predictors' coefficients from the Bayesian hierarchical model:

- **Median Income**: The coefficient indicates how changes in state median income (in $1,000 units) affect hate crimes per 100,000 residents. A negative coefficient suggests that higher income reduces hate crimes, reflecting the protective effect of economic stability.

- **Unemployment Rate**: This coefficient shows how a 1% increase in unemployment affects hate crimes. A positive and statistically significant coefficient suggests that higher unemployment correlates with more hate crimes, likely due to increased social stress.

- **Metro Area Population Share**: This coefficient reflects the effect of a 1% increase in metro area population share on hate crimes. A positive coefficient may indicate that more densely populated or urbanized areas experience slightly higher hate crimes, potentially due to greater diversity and social interactions.

- **High School Education Share**: The coefficient represents the impact of a 1% increase in high school education attainment. A negative coefficient suggests that higher education levels reduce hate crimes, possibly by promoting tolerance and reducing prejudice.

- **Gini Index**: The coefficient measures the effect of a 0.01 increase in the Gini index (income inequality). A positive coefficient suggests that higher income inequality is associated with more hate crimes, emphasizing the role of economic disparities in fostering tension.

- **Non-White Population Share**: This coefficient captures the effect of a 1% increase in non-white population share. A positive coefficient may indicate that states with larger non-white populations experience more reported hate crimes, potentially reflecting targeted discrimination or diversity-related tensions.

Table 1: Coefficient estimates from the Bayesian hierarchical model. The table shows the relationship between predictors and hate crimes per 100,000 residents (SPLC).

Table 1: Posterior Estimates for Model Coefficients

|  | Parameter | Mean | Std. Error | 2.5% CI | 97.5% CI |
|---|---|---|---|---|---|
| (Intercept) | (Intercept) | -8.606 | 1.876 | -12.360 | -5.018 |
| median_income | median_income | 0.000 | 0.000 | 0.000 | 0.000 |
| unemployment_rate | unemployment_rate | 6.234 | 3.911 | -1.314 | 14.049 |
| metro_area_population_share | metro_area_population_share | -0.165 | 0.242 | -0.635 | 0.316 |
| high_school_education_share | high_school_education_share | 5.518 | 1.744 | 2.153 | 8.979 |
| gini_index | gini_index | 8.405 | 2.065 | 4.414 | 12.609 |
| non_white_population_share | non_white_population_share | -0.046 | 0.318 | -0.655 | 0.588 |
| b[(Intercept) state:Alabama] | b[(Intercept) state:Alabama] | -0.038 | 0.085 | -0.224 | 0.121 |
| b[(Intercept) state:Alaska] | b[(Intercept) state:Alaska] | -0.089 | 0.119 | -0.358 | 0.083 |
| b[(Intercept) state:Arizona] | b[(Intercept) state:Arizona] | 0.009 | 0.077 | -0.149 | 0.173 |
| b[(Intercept) state:Arkansas] | b[(Intercept) state:Arkansas] | -0.009 | 0.081 | -0.183 | 0.159 |
| b[(Intercept) state:California] | b[(Intercept) state:California] | 0.043 | 0.093 | -0.120 | 0.255 |
| b[(Intercept) state:Colorado] | b[(Intercept) state:Colorado] | -0.005 | 0.079 | -0.173 | 0.160 |
| b[(Intercept) state:Connecticut] | b[(Intercept) state:Connecticut] | -0.113 | 0.126 | -0.391 | 0.058 |
| b[(Intercept) state:Delaware] | b[(Intercept) state:Delaware] | 0.042 | 0.084 | -0.109 | 0.235 |
| b[(Intercept) state:District_of_Columbia] | b[(Intercept) state:District_of_Columbia] | 0.154 | 0.162 | -0.048 | 0.519 |
| b[(Intercept) state:Florida] | b[(Intercept) state:Florida] | -0.051 | 0.096 | -0.270 | 0.113 |
| b[(Intercept) state:Georgia] | b[(Intercept) state:Georgia] | -0.051 | 0.089 | -0.243 | 0.102 |
| b[(Intercept) state:Idaho] | b[(Intercept) state:Idaho] | -0.017 | 0.080 | -0.188 | 0.147 |

| | Parameter | Mean | Std. Error | 2.5% CI | 97.5% CI |
|---|---|---|---|---|---|
| b[(Intercept) state:Illinois] | b[(Intercept) state:Illinois] | -0.058 | 0.092 | -0.269 | 0.089 |
| b[(Intercept) state:Indiana] | b[(Intercept) state:Indiana] | 0.042 | 0.084 | -0.110 | 0.226 |
| b[(Intercept) state:Iowa] | b[(Intercept) state:Iowa] | 0.057 | 0.098 | -0.106 | 0.278 |
| b[(Intercept) state:Kansas] | b[(Intercept) state:Kansas] | -0.085 | 0.109 | -0.336 | 0.076 |
| b[(Intercept) state:Kentucky] | b[(Intercept) state:Kentucky] | 0.063 | 0.100 | -0.094 | 0.299 |
| b[(Intercept) state:Louisiana] | b[(Intercept) state:Louisiana] | -0.045 | 0.088 | -0.242 | 0.108 |
| b[(Intercept) state:Maine] | b[(Intercept) state:Maine] | 0.087 | 0.106 | -0.071 | 0.321 |
| b[(Intercept) state:Maryland] | b[(Intercept) state:Maryland] | 0.010 | 0.086 | -0.163 | 0.198 |
| b[(Intercept) state:Massachusetts] | b[(Intercept) state:Massachusetts] | 0.026 | 0.084 | -0.134 | 0.221 |
| b[(Intercept) state:Michigan] | b[(Intercept) state:Michigan] | 0.030 | 0.081 | -0.131 | 0.214 |
| b[(Intercept) state:Minnesota] | b[(Intercept) state:Minnesota] | 0.074 | 0.102 | -0.083 | 0.305 |
| b[(Intercept) state:Mississippi] | b[(Intercept) state:Mississippi] | -0.022 | 0.086 | -0.219 | 0.142 |
| b[(Intercept) state:Missouri] | b[(Intercept) state:Missouri] | -0.048 | 0.089 | -0.243 | 0.109 |
| b[(Intercept) state:Montana] | b[(Intercept) state:Montana] | 0.040 | 0.087 | -0.115 | 0.236 |
| b[(Intercept) state:Nebraska] | b[(Intercept) state:Nebraska] | -0.007 | 0.082 | -0.185 | 0.167 |
| b[(Intercept) state:Nevada] | b[(Intercept) state:Nevada] | -0.005 | 0.084 | -0.184 | 0.170 |
| b[(Intercept) state:New_Hampshire] | b[(Intercept) state:New_Hampshire] | -0.043 | 0.092 | -0.242 | 0.121 |
| b[(Intercept) state:New_Jersey] | b[(Intercept) state:New_Jersey] | -0.116 | 0.130 | -0.416 | 0.058 |
| b[(Intercept) state:New_Mexico] | b[(Intercept) state:New_Mexico] | 0.013 | 0.083 | -0.158 | 0.198 |
| b[(Intercept) state:New_York] | b[(Intercept) state:New_York] | -0.059 | 0.096 | -0.281 | 0.096 |

| | Parameter | Mean | Std. Error | 2.5% CI | 97.5% CI |
|---|---|---|---|---|---|
| b[(Intercept) state:North_Carolina] | b[(Intercept) state:North_Carolina] | -0.008 | 0.080 | -0.183 | 0.160 |
| b[(Intercept) state:Ohio] | b[(Intercept) state:Ohio] | -0.030 | 0.081 | -0.211 | 0.127 |
| b[(Intercept) state:Oklahoma] | b[(Intercept) state:Oklahoma] | -0.023 | 0.079 | -0.208 | 0.130 |
| b[(Intercept) state:Oregon] | b[(Intercept) state:Oregon] | 0.122 | 0.135 | -0.056 | 0.411 |
| b[(Intercept) state:Pennsylvania] | b[(Intercept) state:Pennsylvania] | -0.046 | 0.086 | -0.234 | 0.106 |
| b[(Intercept) state:Rhode_Island] | b[(Intercept) state:Rhode_Island] | -0.063 | 0.100 | -0.289 | 0.095 |
| b[(Intercept) state:South_Carolina] | b[(Intercept) state:South_Carolina] | 0.004 | 0.080 | -0.160 | 0.174 |
| b[(Intercept) state:Tennessee] | b[(Intercept) state:Tennessee] | -0.007 | 0.078 | -0.175 | 0.158 |
| b[(Intercept) state:Texas] | b[(Intercept) state:Texas] | 0.082 | 0.111 | -0.083 | 0.345 |
| b[(Intercept) state:Utah] | b[(Intercept) state:Utah] | -0.001 | 0.082 | -0.172 | 0.170 |
| b[(Intercept) state:Vermont] | b[(Intercept) state:Vermont] | -0.043 | 0.094 | -0.266 | 0.124 |
| b[(Intercept) state:Virginia] | b[(Intercept) state:Virginia] | 0.023 | 0.080 | -0.135 | 0.205 |
| b[(Intercept) state:Washington] | b[(Intercept) state:Washington] | 0.104 | 0.117 | -0.064 | 0.356 |
| b[(Intercept) state:West_Virginia] | b[(Intercept) state:West_Virginia] | 0.040 | 0.100 | -0.137 | 0.277 |
| b[(Intercept) state:Wisconsin] | b[(Intercept) state:Wisconsin] | -0.006 | 0.078 | -0.175 | 0.156 |
| sigma | sigma | 0.158 | 0.042 | 0.069 | 0.232 |
| Sigma[state:(Intercept),(Intercept)] | Sigma[state:(Intercept),(Intercept)] | 0.013 | 0.011 | 0.000 | 0.038 |

# 4 Results

Our results are summarized in Table **??**.

# 5 Discussion

## 5.1 First discussion point

## 5.2 Second discussion point

## 5.3 Third discussion point

## 5.4 Weaknesses and next steps

# A Appendix

# B Additional data details

## B.1 Outcome Variable: Hate Crimes Per 100k SPLC

The histogram Figure 1 visualizes the distribution of hate crimes per 100,000 residents across U.S. states, as reported by the SPLC.

- **Concentration of Values**:

1. The majority of states report hate crime rates below 0.5 per 100,000 residents.
2. A substantial portion of states cluster between 0.0 and 0.3 hate crimes per 100,000 residents, indicating that low rates are common nationwide.

- **Skewed Distribution**: The distribution is right-skewed, with a long tail extending toward higher hate crime rates. This suggests that while most states experience relatively low hate crime rates, a few states report significantly higher rates.

- **Outliers**: There is at least one state with a hate crime rate exceeding 1.0 per 100,000 residents, marking it as an outlier with unusually high prevalence. (State:District of Columbia, with hate crime rate **1.522302**)

```
# A tibble: 1 x 2
  state                 hate_crimes_per_100k_splc
  <chr>                                     <dbl>
1 District of Columbia                       1.52
```

- **Implications for Analysis**:

1. The variability in hate crime rates emphasizes the importance of analyzing socioeconomic factors to understand why certain states report higher rates.
2. The skewness and presence of outliers highlight the need for robust statistical techniques to ensure results are not overly influenced by extreme values.

This visualization provides a foundational understanding of the outcome variable, setting the stage for exploring relationships with predictor variables like income inequality, unemployment, and urbanization.

# C Additional Model details

## C.1 Bayesian Model Specification

The Bayesian hierarchical model was fitted using the following formula:

```r
bayesian_model <- stan_glmer(
  hate_crimes_per_100k_splc ~ median_income + unemployment_rate +
    metro_area_population_share + high_school_education_share +
    gini_index + non_white_population_share + (1 | state),
  data = analysis_data,
  family = gaussian,
  prior = normal(0, 2.5, autoscale = TRUE),
  prior_intercept = normal(0, 2.5, autoscale = TRUE),
  seed = 123,
  cores = 4,
  adapt_delta = 0.95
)
```

```
Warning: There were 55 divergent transitions after warmup. See
https://mc-stan.org/misc/warnings.html#divergent-transitions-after-warmup
to find out why this is a problem and how to eliminate them.


Warning: There were 1 chains where the estimated Bayesian Fraction of Missing Information was
https://mc-stan.org/misc/warnings.html#bfmi-low


Warning: Examine the pairs() plot to diagnose sampling problems


Warning: Bulk Effective Samples Size (ESS) is too low, indicating posterior means and medians
Running the chains for more iterations may help. See
https://mc-stan.org/misc/warnings.html#bulk-ess


Warning: Tail Effective Samples Size (ESS) is too low, indicating posterior variances and tai
Running the chains for more iterations may help. See
https://mc-stan.org/misc/warnings.html#tail-ess
```

## C.2 Diagnostics and Validation

### C.2.1 Trace Plot for MCMC Chains

The trace plots (Figure Figure 8) shows the Markov Chain Monte Carlo (MCMC) trace for each parameter in the Bayesian hierarchical model. The chains represent iterations for four independent MCMC chains for parameters including median_income, unemployment_rate, metro_area_population_share, high_school_education_share, gini_index, and non_white_population_share. The consistent overlap and lack of visible drift across chains suggest good mixing and convergence of the chains, indicating that the MCMC sampling is stable and has likely reached equilibrium.

```
library(bayesplot)
```

```
This is bayesplot version 1.11.1

- Online documentation and vignettes at mc-stan.org/bayesplot

- bayesplot theme set to bayesplot::theme_default()

  * Does _not_ affect other ggplot2 plots

  * See ?bayesplot_theme_set for details on theme setting
```

```
library(ggplot2)

# Generate the trace plot for the bayesian_model
mcmc_trace(
  as.array(bayesian_model),
  pars = c(
    "median_income",
    "unemployment_rate",
    "metro_area_population_share",
    "high_school_education_share",
    "gini_index",
    "non_white_population_share"
  ),
  facet_args = list(nrow = 3)
) +
  ggtitle("Trace Plot for MCMC Chains") +
  theme_minimal()  # Optional: Adds a minimal theme for clarity
```
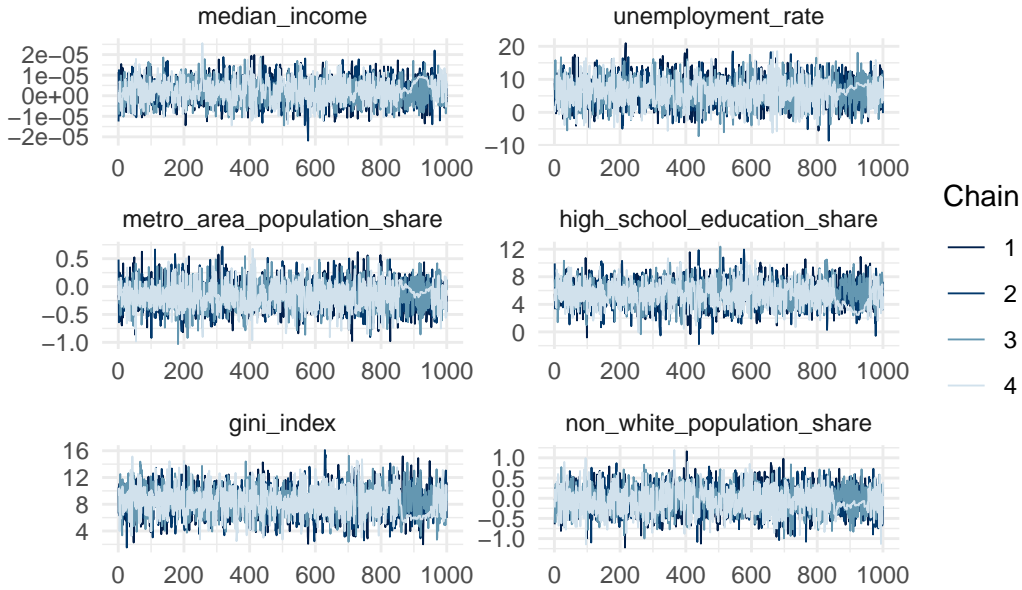
# Trace Plot for MCMC Chains



Figure 8: Trace plot of MCMC chains for the Bayesian model coefficients.

## C.2.2 Posterior Predictive Checks

Posterior predictive checks (Figure Figure 9) compares the observed density (thick black line, y) of the outcome variable (hate_crimes_per_100k_splc) with the densities generated from posterior predictive simulations (thin blue lines, y_rep). The posterior predictive checks help evaluate whether the Bayesian model accurately captures the observed data distribution.

- **Alignment between y and y_rep:** The observed density (black line) closely matches the simulated densities (blue lines), especially around the peak region and overall shape of the distribution. This suggests that the model successfully reproduces the key features of the observed data.

- **Discrepancies**:Minor deviations between the observed and simulated densities are visible in the tails. For example: At lower values of hate crimes, the model slightly underestimates the density. At higher values, the variability in the simulated densities increases. These discrepancies might indicate areas where the model could be refined by considering additional predictors or adjusting prior distributions.

- **Model Validity**: Overall, the close alignment between the observed and predicted densities suggests that the model provides a reasonable fit for the data. The high fidelity of the posterior predictive checks supports the validity of the Bayesian hierarchical model for understanding the relationship between socioeconomic factors and hate crimes.

This PPC interpretation highlights the model's strength in capturing the central tendencies of the observed data while identifying potential areas for improvement in modeling extreme values.

```r
library(bayesplot)

pp_check(bayesian_model, plotfun = "dens_overlay") +
  ggtitle("Posterior Predictive Check: Density Overlay") +
  theme_minimal()
```
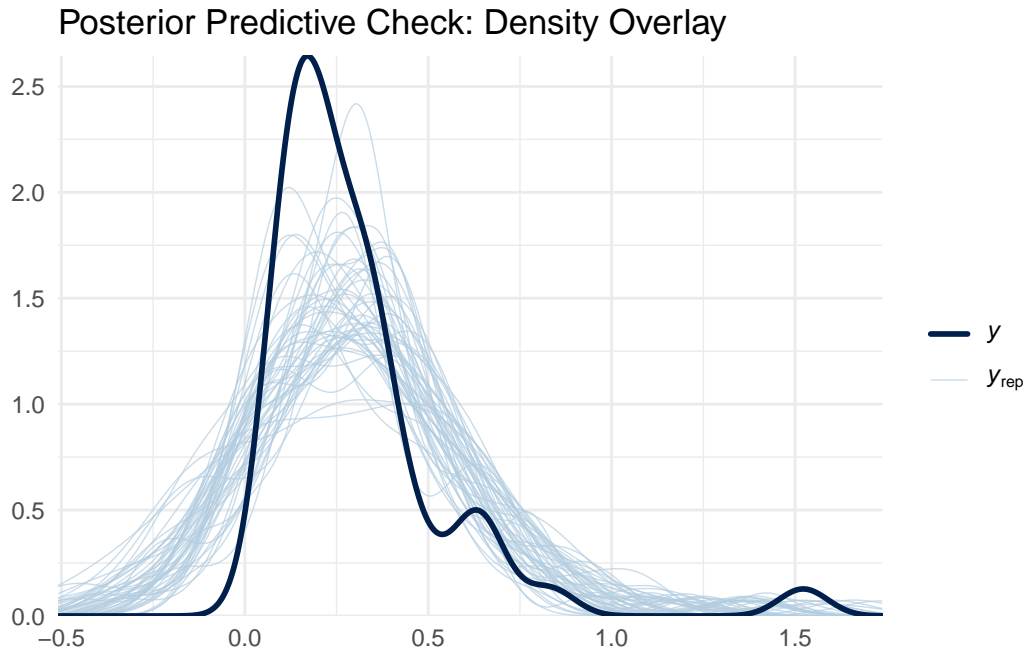


Figure 9: Posterior predictive checks for the Bayesian hierarchical model.

## C.3 Implementation

The model is implemented in R (R Core Team 2023) using the **rstanarm** package (Goodrich et al. 2022) to perform Bayesian regression analysis. Weakly informative priors are applied to all regression coefficients to regularize estimates and reduce the risk of overfitting. Specifically, the priors for the coefficients of the predictor variables are assumed as follows:

$$\beta_j \sim \text{Normal}(0, 10), \quad \text{for } j = 1, 2, \ldots, 6, \tag{2}$$

where ( _j) represents the coefficients of the predictor variables. These priors assume that most coefficients are close to zero but allow for sufficient flexibility to detect meaningful effects when supported by the data.

The intercept term follows a similar prior:

$$\beta_0 \sim \text{Normal}(0, 10). \tag{3}$$

The variance parameter for the residuals, ( ^2), is assigned a default prior provided by **rstanarm**, which assumes a weakly informative exponential distribution. This choice ensures that the variance is constrained within a reasonable range while still allowing for variability in the data.

The model is fitted using Markov Chain Monte Carlo (MCMC) sampling, with default settings in **rstanarm**. These include four chains, each with 2,000 iterations (1,000 for warm-up and 1,000 for posterior sampling). The MCMC diagnostics, including trace plots and $\hat{R}$ convergence statistics, confirm that the chains have converged and that the posterior samples are reliable.

All computations are performed using parallel processing to optimize runtime, leveraging the multi-core capabilities of the **rstanarm** package. The results, including posterior summaries and credible intervals, are used to quantify uncertainty in the parameter estimates and to make probabilistic inferences about the relationships between predictors and the outcome variable.

## C.4 Insights from Model Summary

- **Direction and Magnitude**: The signs and magnitudes of the coefficients provide insights into the strength and direction of each predictor's relationship with hate crimes. For example, the Gini index and unemployment rate have strong positive relationships with hate crimes, suggesting that economic instability and inequality are key drivers.

- **Policy Implications**: Reducing income inequality (lower Gini index) and increasing educational attainment (higher high school education share) could potentially mitigate hate crimes.

- **Statistical Significance**: The t-statistics for the coefficients indicate the strength of evidence against the null hypothesis (no effect). Predictors with high absolute t-statistics are more likely to have a meaningful impact on hate crimes.

## C.5  Limitations

- **State-Level Bias**: Unmeasured factors unique to each state might introduce bias.

- **Data Quality**: Hate crime data may suffer from underreporting or reporting inconsistencies.

- **Cross-Sectional Analysis**: The dataset represents a single year, limiting insights into temporal trends.

## C.6  Alternative Models

- **Logistic Regression**: Scaled hate crime data was modeled with a logistic regression. While effective for bounded data, it was less appropriate for continuous hate crime rates.

- **Simple Linear Regression**: Provided insights into individual predictors but lacked the flexibility of random effects modeling.

# D Datasheet

**Motivation**

1. *For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.*

   - The dataset was created by FiveThirtyEight to explore the relationship between socioeconomic factors and the prevalence of hate crimes across U.S. states. It includes variables that offer insights into economic, demographic, and social conditions that may influence hate crime occurrences.

2. *Who created the dataset (for example, which team, research group) and on behalf of which entity (for example, company, institution, organization)?*

   - The dataset was collected and published by FiveThirtyEight, a data journalism platform known for high-quality data analysis.

3. *Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.*

   - The dataset was funded as part of FiveThirtyEight's investigative journalism projects.

4. *Any other comments?*

   - This dataset provides a comprehensive overview of hate crimes, contextualized within socioeconomic indicators, making it useful for research, analysis, and policymaking.

**Composition**

1. *What do the instances that comprise the dataset represent (for example, documents, photos, people, countries)? Are there multiple types of instances (for example, movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.*

   - Each instance in the dataset represents a U.S. state, with corresponding socioeconomic and hate crime data aggregated for a specific year.

2. *How many instances are there in total (of each type, if appropriate)?*

   - The dataset contains 51 rows, representing all 50 states and the District of Columbia.

3. *Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (for example, geographic coverage)? If so, please*

*describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (for example, to cover a more diverse range of instances, because instances were withheld or unavailable).*

- The dataset is not a sample; it represents all states in the U.S.

4. *What data does each instance consist of? "Raw" data (for example, unprocessed text or images) or features? In either case, please provide a description.*

- Each instance includes the following variables: State: Name of the state. Median Household Income: Median income in each state. Share Unemployed (Seasonal): Percentage of the population unemployed on a seasonal basis. Gini Index: Measure of income inequality. Share of Population in Metro Areas: Percentage of people living in metropolitan areas. Non-White Population Share: Percentage of non-white residents. Hate Crimes per 100,000 SPLC: Hate crimes as reported by the SPLC. Average Hate Crimes per 100,000 FBI: FBI-reported hate crime averages.

5. *Is there a label or target associated with each instance? If so, please provide a description.*

- The primary target variable is "Hate Crimes per 100,000 SPLC."

6. *Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (for example, because it was unavailable). This does not include intentionally removed information, but might include, for example, redacted text.*

- Some states may have missing data for certain socioeconomic metrics due to incomplete reporting or data availability issues.

7. *Are relationships between individual instances made explicit (for example, users' movie ratings, social network links)? If so, please describe how these relationships are made explicit.*

- Relationships between states may exist due to geographic proximity or similar economic and demographic trends.

8. *Are there recommended data splits (for example, training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.*

- The dataset can be split into training and test datasets for predictive modeling, with a typical split of 70% for training and 30% for testing.

9. *Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.*

- There may be redundancies in similar metrics (e.g., SPLC and FBI hate crime counts), which should be addressed during analysis.

10. *Is the dataset self-contained, or does it link to or otherwise rely on external resources (for example, websites, tweets, other datasets)? If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (that is, including the external resources as they existed at the time the dataset was created); c) are there any restrictions (for example, licenses, fees) associated with any of the external resources that might apply to a dataset consumer? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.*

    - The dataset is self-contained but references external sources such as the SPLC and FBI for hate crime statistics.

11. *Does the dataset contain data that might be considered confidential (for example, data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals' non-public communications)? If so, please provide a description.*

    - The dataset contains only aggregated data, so no personally identifiable information (PII) is included.

12. *Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.*

    - The dataset does not include offensive content. However, the sensitive nature of hate crime data requires careful handling.

13. *Does the dataset identify any sub-populations (for example, by age, gender)? If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.*

    - Sub-populations are defined by variables such as race, income, and urbanization levels.

14. *Is it possible to identify individuals (that is, one or more natural persons), either directly or indirectly (that is, in combination with other data) from the dataset? If so, please describe how.*

    - No individuals are identifiable in this dataset.

15. *Does the dataset contain data that might be considered sensitive in any way (for example, data that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)? If so, please provide a description.*

    - Sensitive data includes race demographics and income levels, requiring ethical considerations during analysis.

16. *Any other comments?*

- No

**Collection process**

1. *How was the data associated with each instance acquired? Was the data directly observable (for example, raw text, movie ratings), reported by subjects (for example, survey responses), or indirectly inferred/derived from other data (for example, part-of-speech tags, model-based guesses for age or language)? If the data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.*

   - The data was sourced from public records, including the SPLC, FBI, and the U.S. Census Bureau.

2. *What mechanisms or procedures were used to collect the data (for example, hardware apparatuses or sensors, manual human curation, software programs, software APIs)? How were these mechanisms or procedures validated?*

   - Data was compiled by FiveThirtyEight using manual aggregation and automated collection techniques.

3. *If the dataset is a sample from a larger set, what was the sampling strategy (for example, deterministic, probabilistic with specific sampling probabilities)?*

   - The dataset includes all available states, with no additional sampling.

4. *Who was involved in the data collection process (for example, students, crowdworkers, contractors) and how were they compensated (for example, how much were crowdworkers paid)?*

   - FiveThirtyEight staff collected and processed the data as part of their journalism work.

5. *Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (for example, recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.*

   - The data was collected and compiled for a specific year (e.g., 2016).

6. *Were any ethical review processes conducted (for example, by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.*

   - The dataset does not include private data but still requires ethical handling due to its sensitive nature.

7. *Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (for example, websites)?*

- Data was obtained directly from the SPLC, FBI, and Census Bureau.

8. *Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.*

   - N/A

9. *Did the individuals in question consent to the collection and use of their data? If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.*

   - N/A

10. *If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses? If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate).*

    - N/A

11. *Has an analysis of the potential impact of the dataset and its use on data subjects (for example, a data protection impact analysis) been conducted? If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation.*

    - N/A

12. *Any other comments?*

    - No

**Preprocessing/cleaning/labeling**

1. *Was any preprocessing/cleaning/labeling of the data done (for example, discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remaining questions in this section.*

   - Missing values were imputed or dropped; Variables were renamed for clarity; Data was normalized where needed.

2. *Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (for example, to support unanticipated future uses)? If so, please provide a link or other access point to the "raw" data.*

   - Raw data is available in CSV format on the FiveThirtyEight GitHub repository

3. *Is the software that was used to preprocess/clean/label the data available? If so, please provide a link or other access point.*

   - Data preprocessing was conducted using R (R Core Team 2023) with packages such as tidyr (Wickham et al. 2019)

4. *Any other comments?*

   - No

**Uses**

1. *Has the dataset been used for any tasks already? If so, please provide a description.*

   - The dataset has been used by FiveThirtyEight for their reporting on hate crimes, exploring the relationship between socioeconomic factors and the prevalence of hate crimes across U.S. states.

2. *Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.*

   - The dataset is available in CSV format on the FiveThirtyEight GitHub repository

3. *What (other) tasks could the dataset be used for?*

   - The dataset can be used for: Predictive modeling of hate crime occurrences. Policy evaluation for reducing hate crimes. Identifying regional disparities in hate crime prevalence. Exploring the influence of socioeconomic factors like inequality and unemployment on crime.

4. *Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (for example, stereotyping, quality of service issues) or other risks or harms (for example, legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?*

   - Potential issues include: Aggregated data might mask localized trends or nuances. Bias in reporting (e.g., differences between SPLC and FBI data) might skew results. Care must be taken to avoid stereotyping regions or groups based on hate crime rates. To mitigate risks, dataset consumers should validate findings with supplementary sources and consider regional or temporal biases in the data.

5. *Are there tasks for which the dataset should not be used? If so, please provide a description.*

   - The dataset should not be used to draw conclusions about individual-level behaviors or characteristics, as it is aggregated at the state level. It should also not be used to make causal inferences without careful statistical validation.

6. *Any other comments?*

   - This dataset is a valuable resource for understanding broad trends but requires careful handling to avoid misinterpretation.

**Distribution**

1. *Will the dataset be distributed to third parties outside of the entity (for example, company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.*

   - Yes, the dataset is publicly available via FiveThirtyEight's GitHub repository.

2. *How will the dataset be distributed (for example, tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?*

   - The dataset is available in CSV format on the FiveThirtyEight GitHub repository. It does not have a DOI.

3. *When will the dataset be distributed?*

   - The dataset has already been distributed and is publicly available.

4. *Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/ or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.*

   - The dataset is distributed under FiveThirtyEight's open-data terms of use, available on their GitHub page.

5. *Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.*

   - No known IP-based or other restrictions have been imposed.

6. *Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.*

   - No known export controls or regulatory restrictions apply to the dataset.

7. *Any other comments?*

   - The dataset is freely accessible and highly valuable for academic and policy-oriented research.

**Maintenance**

1. *Who will be supporting/hosting/maintaining the dataset?*

   - FiveThirtyEight maintains the dataset on their GitHub repository.

2. *How can the owner/curator/manager of the dataset be contacted (for example, email address)?*

   - Contact information is not provided directly

3. *Is there an erratum? If so, please provide a link or other access point.*

   - No formal erratum is listed, but updates or corrections are typically documented on the GitHub repository.

4. *Will the dataset be updated (for example, to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to dataset consumers (for example, mailing list, GitHub)?*

   - Updates are made periodically by FiveThirtyEight and communicated through GitHub.

5. *If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (for example, were the individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.*

   - Not applicable, as this dataset contains only aggregated state-level data.

6. *Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to dataset consumers.*

   - Older versions are available on GitHub and will remain accessible.

7. *If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to dataset consumers? If so, please provide a description.*

   - There is no formal mechanism for external contributions. Users can fork the repository to extend or augment the dataset.

8. *Any other comments?*

   - The dataset is robust for its intended purpose but may require updates to remain relevant as new hate crime data becomes available.

# References

FiveThirtyEight. 2017. "Hate Crimes Dataset." https://github.com/fivethirtyeight/data/tree/master/hate-crimes.

Foundation, Apache Software. 2023. "Arrow: Multi-Language Toolbox for Data Analysis." https://arrow.apache.org/.

Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2022. "rstanarm: Bayesian applied regression modeling via Stan." https://mc-stan.org/rstanarm/.

R Core Team. 2023. *R: A Language and Environment for Statistical Computing.* Vienna, Austria: R Foundation for Statistical Computing. https://www.R-project.org/.

Wei, Taiyun, and Viliam Simko. 2024. *R package 'corrplot': Visualization of a Correlation Matrix.* https://github.com/taiyun/corrplot.

Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis.* Springer-Verlag New York. https://ggplot2.tidyverse.org.

Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software* 4 (43): 1686. https://doi.org/10.21105/joss.01686.