

Socioeconomic Determinants of Hate Crimes in the United States*

A Bayesian Hierarchical Analysis

Yuxuan Wei

November 28, 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 6 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, using libraries such as `dplyr` (Wickham et al. 2019), `ggplot2` (Wickham 2016), `arrow` (Foundation 2023), and `corrplot` (`corrplot?`). The dataset, sourced from (FiveThirtyEight 2017), explores the relationships between socioeconomic variables and hate crime prevalence in the United States. It includes aggregated data for all 50 states and the District of Columbia.

After preprocessing, which involved filtering missing values and standardizing key variables, the final analysis dataset consists of **51** observations across **7** key variables: `state`, `median_income`, `unemployment_rate`, `metro_area_population_share`, `high_school_education_share`, `gini_index`, and `non_white_population_share`. These variables capture essential socioeconomic conditions and their potential influence on hate crime rates, measured as `hate_crimes_per_100k_splc`.

2.2 Measurement and Considerations

The dataset focuses on measuring the influence of socioeconomic factors on hate crimes. Hate crime data is aggregated from credible sources, including the SPLC and FBI, and represents the prevalence of reported hate crimes per 100,000 residents.

2.2.1 Data Source and Quality Assurance

The dataset was obtained from FiveThirtyEight(FiveThirtyEight 2017), a trusted source recognized for its rigorous data quality standards. Each state-level observation documents essential socioeconomic and hate crime metrics, including income, unemployment, and population composition.

To ensure comparability across states, socioeconomic indicators were standardized and checked for consistency. Hate crime data, while robust, is subject to underreporting and reporting bias, requiring careful interpretation.

2.2.2 Temporal and Geographic Context

The dataset provides a snapshot of state-level hate crime and socioeconomic data for a single year. While this ensures temporal consistency, it limits the ability to analyze trends over time. Additionally, regional disparities, such as differences in reporting standards or enforcement, may influence observed hate crime rates.

2.3 Variables of Interest

The analysis examines the relationships between `hate_crimes_per_100k_splc` (the outcome variable) and the following socioeconomic predictors. Below, we describe each variable and present visualizations to illustrate its distribution and relevance.

2.3.1 Outcome Variable: Hate Crimes Per 100k SPLC

The outcome variable is the number of hate crimes per 100,000 residents, as reported by the SPLC. Its distribution, shown in Figure 1, highlights variation across states, with most states reporting between 0.1 and 0.5 hate crimes per 100,000 residents. The variability suggests that some states experience disproportionately higher rates of hate crimes.

2.4 Predictor variables

3 Model

The goal of our modelling strategy is twofold. Firstly,...

Here we briefly describe the Bayesian analysis model used to investigate... Background details and diagnostics are included in Appendix C.

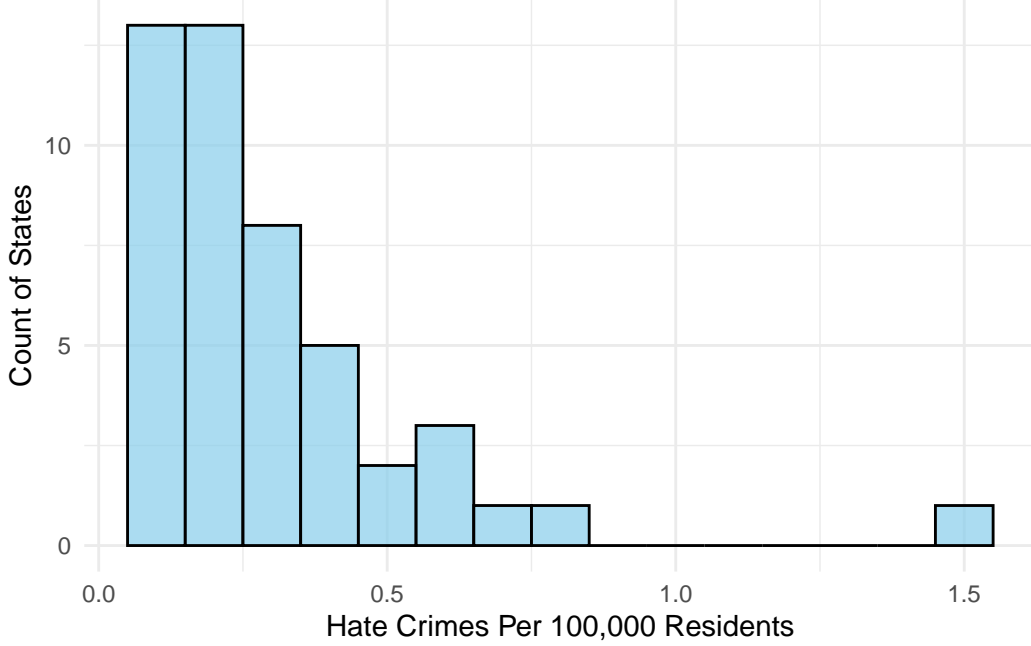


Figure 1: Distribution of Hate Crimes Per 100,000 Residents Across States

3.1 Model set-up

Define y_i as the number of seconds that the plane remained aloft. Then β_i is the wing width and γ_i is the wing length, both measured in millimeters.

$$y_i | \mu_i, \sigma \sim \text{Normal}(\mu_i, \sigma) \quad (1)$$

$$\mu_i = \alpha + \beta_i + \gamma_i \quad (2)$$

$$\alpha \sim \text{Normal}(0, 2.5) \quad (3)$$

$$\beta \sim \text{Normal}(0, 2.5) \quad (4)$$

$$\gamma \sim \text{Normal}(0, 2.5) \quad (5)$$

$$\sigma \sim \text{Exponential}(1) \quad (6)$$

We run the model in R (R Core Team 2023) using the `rstanarm` package of Goodrich et al. (2022). We use the default priors from `rstanarm`.

3.1.1 Model justification

We expect a positive relationship between the size of the wings and time spent aloft. In particular...

We can use maths by including latex between dollar signs, for instance θ .

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

C Additional Model details

C.1 Posterior predictive check

In `?@fig-ppcheckandposteriorvsprior-1` we implement a posterior predictive check. This shows...

In `?@fig-ppcheckandposteriorvsprior-2` we compare the posterior with the prior. This shows...

Examining how the model fits, and is affected
by, the data

C.2 Diagnostics

`?@fig-stanareyouokay-1` is a trace plot. It shows... This suggests...

`?@fig-stanareyouokay-2` is a Rhat plot. It shows... This suggests...

Checking the convergence of the MCMC algo-
rithm

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/>.
- 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 Golemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.