

Developing and validating an anti-Asian area racism index at the county level in the contiguous United States 2020 - 2021

Alexander Hohl^{1*}, Ming Wen^{2,3,4}, Guangzhen Wu², Yue Zhang⁵, Zhenlong Li⁶, Dejun Su⁷

¹School of Environment, Society & Sustainability, The University of Utah, Salt Lake City, UT, USA

²Department of Sociology, The University of Utah, Salt Lake City, UT, USA

³Department of Sociology, The University of Hong Kong, Hong Kong SAR, China

⁴Research Hub of Population Studies, The University of Hong Kong, Hong Kong SAR, China

⁵School of Medicine, The University of Utah, Salt Lake City, UT, USA

⁶Department of Geography, The Pennsylvania State University, University Park, PA, USA

⁷Department of Health Promotion, University of Nebraska Medical Center, Omaha, NE, USA

*Corresponding author: alexander.hohl@ess.utah.edu

Abstract

Historical narratives and the "model minority" myth have obscured the realities of anti-Asian racism in the United States. The escalation of prejudice and hate crimes against Asian Americans during the COVID-19 pandemic further underscored the need for robust measures to quantify this phenomenon. This study proposes a novel county-level index specifically designed to capture the multifaceted nature of anti-Asian racism. The index integrates a diverse data set including geotagged Twitter data assessing public attitudes and potential hate speech directed toward Asian Americans, anti-Asian hate crime records from the Federal Bureau of Investigation's Uniform Crime Reporting System, Google Search Trends data about anti-Asian stereotypes and narratives, and alien land bills denoting context for discriminatory policies against Asian immigrants at the state level. We employed Principal Component Analysis to combine these data sources into a single, composite index. A validation of the index using nationally representative survey data indicates that two of the three identified principal components significantly predict area racism against Asian respondents. This study offers a nuanced understanding of anti-Asian racism and has the potential to inform targeted interventions, the allocation of resources for community support and educational initiatives, and can be instrumental for policymakers in identifying areas with heightened anti-Asian bias. Additionally, the index serves as a foundation for future research, facilitating the exploration of correlations between anti-Asian racism and various health and social outcomes. While limitations exist regarding data subjectivity and availability, this index represents a significant advancement in measuring anti-Asian racism at the county level. It paves the way for a more comprehensive understanding of this critical issue and the development of effective strategies to combat racial injustice and address related geographic disparities.

1. Introduction

Racial bias manifests as a system that disadvantages individuals based solely on their race (DuBois, 2003). It is a well-documented determinant of health disparities, impacting entire racial and ethnic groups beyond just isolated incidents (Gee & Ford, 2011). Studies consistently demonstrate poorer health outcomes among those experiencing racial discrimination (Alhusen et al., 2016; Berger & Sarnyai, 2015). While research has explored individual expressions of prejudice and bigotry, a critical gap lies in addressing the broader, systemic nature of racism (Bailey et al., 2017, 2021). Societal norms and attitudes reflecting this bias, often referred to as "aggregate racism," can be even more detrimental to public health (Cobbinah & Lewis, 2018). This pervasive bias creates chronic stress and fear within targeted communities, leading to a population-level increase in negative health outcomes (Iwamoto & Liu, 2010).

Asians and Asian Americans are often overlooked in studies of racism due to their status as a "model minority" (Horse, 2021; Saito, 1997). Despite the vast contributions Asians make to our society, their presence at the highest levels of government, business, academia, and popular media remains disproportionately low (Huynh et al., 2011). This lack of representation fuels the persistent stereotype of Asian Americans as perpetual outsiders, inherently different from the mainstream (Wu, 2023). This perception creates a pervasive sense of "otherness," where their belonging in American society feels conditional, despite their deep roots and achievements (Sabharwal et al., 2022).

The emergence of COVID-19 in the United States of America (USA) triggered a marked increase in anti-Asian discrimination (Nguyen et al., 2020; Tessler et al., 2020). Data from social media platforms reveal a significant uptick in anti-Asian sentiment, including hate speech (Hohl et al., 2022), while organizations like Stop AAPI Hate have documented over 11,000 incidents targeting Asian Americans and Pacific Islanders that occurred from January 2020 to December 2022 (Stop AAPI Hate, 2023). This escalation reflects a long-standing pattern of racial bias against Asians in the USA, evidenced by a 77% rise in reported anti-Asian hate crimes in 2020 compared to the previous year (The United States Department of Justice, 2020). Research has documented a correlation between the use of pejorative terms like "Chinese virus" by prominent individuals and subsequent spikes in anti-Asian hashtags on social media platforms (Hswen et al., 2021). Even more, researchers found that hate speech on social media fuels hate crimes targeting minorities (Müller & Schwarz, 2023).

Traditionally, research on racism and health relied heavily on self-reported experiences of discrimination (Groos et al., 2018). However, these surveys have limitations, such as recall bias and difficulty capturing subtler forms of racism (Gelman et al., 2007; Krieger et al., 2011). In addition, questionnaire-based data collection may be difficult to implement in a timely manner, especially during times of global crisis (Hohl et al., 2024). The approach to data collection has dramatically transformed in recent times, driven primarily by groundbreaking developments like Web 2.0 and Big Data (Huang et al., 2024). Abundant new human and earth observation data provides a feasible way to gain such information in near real-time (Jordan et al., 2018; Nguyen et al., 2021). Such data may be collected passively (without active participant involvement) as social media posts (e.g. Twitter/X), online search queries (Google), or actively (participants actively provide data) by government authorities (e.g. Federal Bureau of Investigation Uniform Crime Reporting Program, FBI UCR) or civil rights organizations (e.g. APA Justice) (Singleton et al., 2018). Although these novel data sources each have biases, the goal is to collectively provide an alternative that offers a presumably more objective view of the phenomenon of interest.

Therefore, this research aims to create a county-level composite index of anti-Asian area racism for the contiguous USA. By leveraging multiple data sources, we create a data product that can inform interventions that promote resilience and empower Asian Americans to combat racial injustice. This paper is structured as follows: We describe data collection and aggregation in Section 2.1 (Data), along with

Principal Components Analysis (PCA) and the anti-Asian area racism index design in Section 2.2 (Methods). Section 3 (Results) present the resulting index and its spatial distribution. Section 4 (Discussion) provides a synopsis of the study, discusses its limitations, and highlights its relevance. Finally, Section 5 (Conclusions) offers remarks on the usage and future prospects of the index. To ensure the reproducibility of our study, we have shared all data and code in an open-source repository.¹

2. Data and Methods

2.1 Data

We collected human and earth observation data from four different sources (Twitter, FBI UCR, Google Search Trends, Alien Land Bills), which cover a wide range of anti-Asian hate. We aggregated all data to the county level within the contiguous USA ($n = 3108$) of 2020.

2.1.1 Twitter

We collected 17,385,878 geotagged tweets using the public Twitter (now X) Streaming API. We filtered out tweets from non-human sources, such as automated weather updates, employment advertisements, and promotional content. For instance, tweets from sources like TweetMyJOBS and CareerArc were excluded (Martín et al., 2020). We used a catalog of Twitter sources indicative of human-authored tweets, identified through manual inspection (Li et al., 2021; Martín et al., 2020). Our analysis only included tweets from these verified sources, sent from the contiguous USA, written in English, sent between December 12th, 2019, and August 1st, 2022, and including the topic of COVID-19 (Hohl et al., 2022). This period starts shortly before the arrival of the virus in the USA (January 2020) and ends approximately after fourth wave where Omicron subvariants BA.4 and BA.5 drove infections (Centers for Disease Control and Prevention (CDC), 2024).

We categorized the tweets as either hateful or non-hateful based on the inclusion of specific keywords associated with anti-Asian hate within the tweet body (see Appendix Table A1). These keywords are indicative of anti-Asian hate both generally and in the context of COVID-19 and were compiled from three different sources: 1) A study on the spatial distribution of anti-Asian hate on social media during COVID-19 in the USA (Hohl et al., 2022) 2) hatebase.org, a now-retired service designed to assist organizations and online communities in detecting, monitoring, and analyzing hate speech, and 3) a glossary of anti-Asian terms by the Committee of 100, a non-profit organization of Chinese Americans (Committee of 100, 2022). We chose this method for its simplicity and absence of training data requirement, rather than employing more complex machine learning classifiers.

In addition, we performed sentiment analysis on the tweet bodies using the Valence Aware Dictionary and sEntiment Reasoner (VADER), a sentiment analysis tool specifically attuned to social media (Hutto & Gilbert, 2014). For each tweet, VADER produces sentiment scores in four categories: Negative, neutral, positive, and compound. The negative, neutral, and positive scores represent the ratios of text that fall into each category, meaning they sum up to 1. The compound score is calculated by adding the valence scores of each word in the lexicon, adjusting them according to the rules, and then normalizing the result to a range between -1 (most extremely negative) and +1 (most extremely positive).

¹ <https://github.com/alexandster/AAPI-hate-index>

Each tweet has location information as a latitude/longitude coordinate pair, which allowed us to compute the counts of hateful and non-hateful tweets, as well as county-level sums and averages of sentiment scores. We achieved this through spatial join of the tweet point locations and the 2020 U.S. Census county geometries (United States Census Bureau, 2020), utilizing PostgreSQL/PostGIS spatial database software (Strobl, 2008).

2.1.2 Federal Bureau of Investigation Uniform Crime Reporting Program (FBI UCR) Hate Crimes

Hate crimes are defined as “crimes in which the perpetrators acted based on a bias against the victim’s race, color, religion, or national origin,” as well as “actual or perceived sexual orientation, gender identity, disability, or gender” (Federal Bureau of Investigation, 2021). We obtained hate crime data from the FBI UCR for the years 2020 and 2021 (United States Department of Justice. Federal Bureau of Investigation, 2023). These data include the county identifier (5-digit FIPS code), bias motivation, as well as the number of victims and offenders for each hate crime. We selected all records with anti-Asian bias motivation, allowing us to compute the counts of hate crime victims and offenders for each county in the contiguous USA.

2.1.3 Google Search Trends

Provided by Google, Google Search Trends is a tool that allows users to see how frequently specific search terms are entered into Google's search engine over a certain period. It provides insights into the popularity of search queries, showing patterns and trends in search behavior globally or within specific regions. It also allows for comparing search interests across geographic regions, whereas users can choose between state-level, metro-level, and city-level. Search trend values are calculated on a scale from 0 to 100, with 100 representing the location with the highest popularity as a fraction of total searches in that location. This prevents locations with the highest search volume from always being ranked the highest.

We chose to obtain scores for the top 5 most frequent anti-Asian hate terms in our Twitter data (chinavirus, wuhanvirus, chinesevirus, ccpvirus, chinaliedpeoplepledied), and averaged the values to quantify overall interest in anti-Asian hate terms. The reason for restricting the number of hate terms is the current impossibility of automated retrieval of Google Search Trends data, necessitating manual retrieval, which limits the scalability of data collection. We chose to obtain Google Search Trends data at the metro area-level for the years 2020 and 2021. The metro area-level strikes a balance between spatial granularity and data sparseness. Metros are geographical areas that generally correspond to metropolitan areas. In our study, each county within the same metro area was assigned the metro area-level search trend score.

2.1.4 Alien Land Bills

We collected the status of alien land bills for all 49 states in the contiguous USA and the District of Columbia from Asian Pacific American (APA) Justice, an organization dedicated to advocating for civil rights, social justice, and equitable treatment of APAs. Alien land bills are legislation or regulations concerning land ownership or use by non-citizens or extraterrestrials, often referred to as "alien" in legal terms. We used APA Justice data that classifies each state as either “Passed alien land bills into state law”, “Died in current/recent legislative session”, “Introduced alien land bills”, or “No known alien land bills introduced” as of 2023 (APA Justice, 2023a, 2023b). All counties of a state are assigned the corresponding state-level classification.

2.2 Methods

2.2.1 Variable Normalization

We derived 21 variables (Table 1) from the data outlined in Section 2.1. We performed normalization by either the resident population (*tweets_rate*, *hateful_tweets_rate*, *user_rate*, *vics_rate*, *off_rate*), the total number of tweets (*hateful_tweets_prop*), the number of users (*tweets_per_user*, *hateful_tweets_per_user*), or the number of hateful users (*hateful_tweets_per_hateful_users*; hateful users - those who sent at least one hateful tweet). Since the tweets included timestamps, we quantified the temporal variation in tweet volume as the difference in tweet volume between the highest and lowest temporal units (months for *monthly_variation*, weeks for *weekly_variation*) divided by the highest volume. We derived two sets of sentiment analysis-based variables: the first set averages the sentiment scores within a county (*sent_neg*, *sent_neu*, *sent_pos*, *sent_compound*), whereas the second set takes their sum (*sent_neg_sum*, *sent_neu_sum*, *sent_pos_sum*) to quantify the volume of the corresponding sentiment. While the *GST.interest* variable is already normalized by Google, *ALB* is a state-level ordinal variable that takes on four different values: “Passed alien land bills into state law”, “Died in current/recent legislative session”, “Introduced alien land bills”, or “No known alien land bills introduced”.

Maps of the spatial distribution of each variable are found in the Appendix (Appendix Figures A1 - A21).

Data	Variable	Notes	Calculation
Twitter	<i>tweets_rate</i>	Number of tweets per population	tweets / population
	<i>hateful_tweets_prop</i>	Hateful tweet proportion	hateful tweets / total tweets
	<i>hateful_tweets_rate</i>	Number of hateful tweets per population	hateful tweets / population
	<i>tweets_per_user</i>	Number of tweets per user	tweets / users
	<i>user_rate</i>	Number of users per population	users / population
	<i>hateful_tweets_per_user</i>	Number of hateful tweets per user	hateful tweets / users
	<i>hateful_tweets_per_hateful_users</i>	Number of hateful tweets per hateful users	hateful tweets / hateful users
	<i>monthly_variation</i>	Monthly variation in the number of tweets	(max(n) - min(n))/max(n)
	<i>weekly_variation</i>	Weekly variation in the number of tweets	(max(n) - min(n))/max(n)
	<i>sent_neg</i>	Average negative sentiment [0 - 1]	
	<i>sent_neu</i>	Average neutral sentiment [0 - 1]	
	<i>sent_pos</i>	Average positive sentiment [0 - 1]	
	<i>sent_compound</i>	Average compound sentiment [-1 - 1]	
FBI UCR	<i>sent_neg_sum</i>	Sum of negative sentiment	
	<i>sent_neu_sum</i>	Sum of neutral sentiment	
	<i>sent_pos_sum</i>	Sum of positive sentiment	
FBI UCR	<i>vics_rate</i>	Number of hate crime victims per population	victims / population

	<i>off_rate</i>	Number of hate crime offenders per population	offenders / population
	<i>vics_off</i>	Number of victims per offender	victims / offenders
Google Search Trends	<i>GST.interest</i>	Average Google Search Trends interest	
Alien Lands Bills	<i>ALB</i>	Alien Land Bills	

Table 1: List of variables and their descriptions.

2.2.2 Principal Components Analysis (PCA)

We performed PCA to create the anti-Asian area racism index. PCA is a statistical procedure used in exploratory data analysis and machine learning to simplify the complexity of high-dimensional data while retaining trends and patterns (Pearson, 1901). It transforms the original variables into a new set of uncorrelated variables known as principal components. These components are ordered so that the first few retain most of the variation in the original variables. The steps to compute PCA include 1) standardization to eliminate the effect of the scale of the data, 2) covariance matrix computation to understand how the input variables relate to each other, 3) eigenvalue (variance explained by each component) and eigenvector calculation to provide magnitude and direction of the principal components, 4) principal components selection of the top k eigenvectors, and 5) data projection, transforming the data into the new feature space (Jolliffe, 2014). PCA has the benefits of 1) dimensionality reduction, which helps reduce computational cost and complexity of downstream analysis, as well as data visualization and interpretation, and 2) noise reduction by identifying the components that capture the most variance and ignoring the less significant ones (Izenman, 2008). PCA has been applied in many domains, including the social sciences (Anselin et al., 2007; Bitter et al., 2007), ecology (Bastianoni et al., 2008), environmental sciences (Parinet et al., 2004; Tran et al., 2002), and geography (Demšar et al., 2013).

In this study, we performed PCA on the 21 variables described in Section 2.1. We selected the number of principal components by drawing the scree plot (Lewith et al., 2010), which plots the eigenvalues from largest to smallest. Together, the selected components constitute our anti-Asian racism index. The number of components is identified by looking for an “elbow”, a point where the curve flattens out. This point suggests the number of components capturing most of the variance while minimizing data redundancy. Then, we plotted the cosine score (Cos2) that quantifies variable representation for each selected principal component (Kassambara, 2017). This allows for characterizing each principle component and helps with interpretation. Lastly, we mapped principal component values back to their respective counties and drew choropleth maps to illustrate their spatial distribution.

2.2.3 Index Validation

We validated our anti-Asian racism index by comparing it to the Health, Ethnicity and Pandemic (HEAP) Survey (Su et al., 2022), conducted by the National Opinion Research Center (NORC) at the University of Chicago ($N = 2,506$). Racism experience was measured using a survey question: “Have you personally experienced any discrimination or unfair treatment because of your racial or ethnic background during the COVID-19 pandemic?” HEAP data includes demographic variables such as race/ethnicity and geographic identifiers such as county FIPS codes. NORC employed stratified sampling and post-stratification

weighting methodologies to construct this nationally representative survey of U.S. adults aged 18 years and older.

For each county with respondents in the survey, we calculated the weighted sum of experienced racism ('Yes' responses) from Asian respondents to the survey question above. We converted the sums to rates per 100,000 population, using the 2020 census data. Then, we performed linear regression, using the three dimensions of our anti-Asian racism index to predict the rate of experienced racism. We created four models: three univariate regressions (one for each dimension), to evaluate the effect of each component on experienced racism, and one multiple regression that included all three dimensions, to evaluate their relative effect. A good fit indicates that the index is capable of representing area-level anti-Asian racism, as measured by independent survey data.

3. Results

The spatial distribution of the 21 variables is illustrated in Appendix Figures A1 – A21. It is noteworthy that the variables have varying degrees of sparseness, from high (hate crime-related variables, Appendix Figures A17 -A19) to moderate (Twitter-derived variables, Appendix Figures A1 – A16), to complete spatial coverage (Google Search Trends, Appendix Figure A20, and Alien Land Bills, Appendix Figure A21).

3.1 Area-level Anti-Asian Racism Index

We identified the first three principal components, which jointly explain 81.94% of the variation of the 21 input variables (PC 1: 42.35%, PC 2: 33.89%, PC 3: 5.7%), and therefore, were selected as our anti-Asian area racism index at the county level in the contiguous USA. We decided against a strict interpretation of the elbow rule and to include the third component (PC 3), as it gives us a wider coverage of anti-Asian racism topics (see next section).

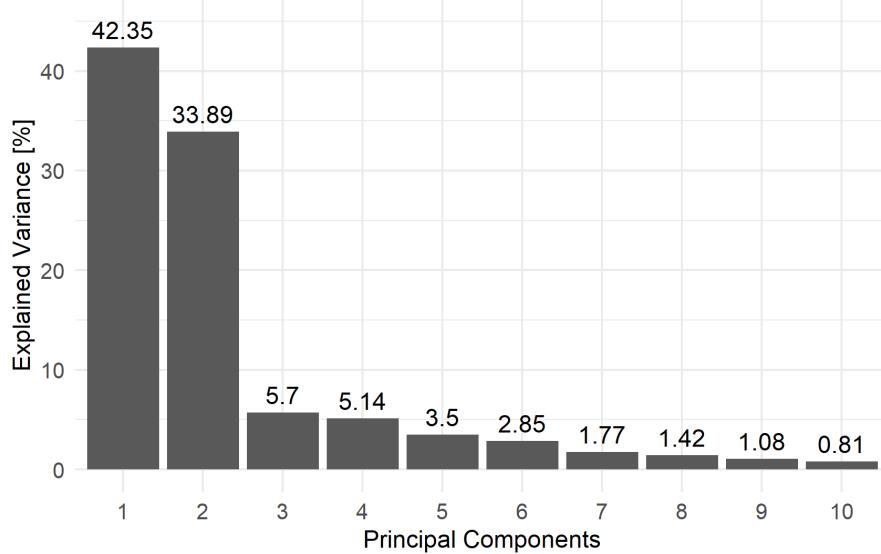


Figure 1: Scree plot of principal components.

The variable representation plot of the first component shows that the variables *sent_neu* and *GST.interest* are best represented (Figure 2), with a cosine score of 0.0892 and 0.0604, respectively. With the exception of the *ALB* variable, all top 5 represented variables of this component stem from the Twitter

data. Therefore, we call the first dimension of our anti-Asian area racism index the “Social media temporal variation (SMTV)” component. The spatial distribution of the first component shows low values in the Midwest, medium values in the West, the Rocky Mountains, and the Southeast, and high values in the Northeast, and the states of Oregon, New Mexico, Nevada, Minnesota, and Wisconsin (Figure 3).

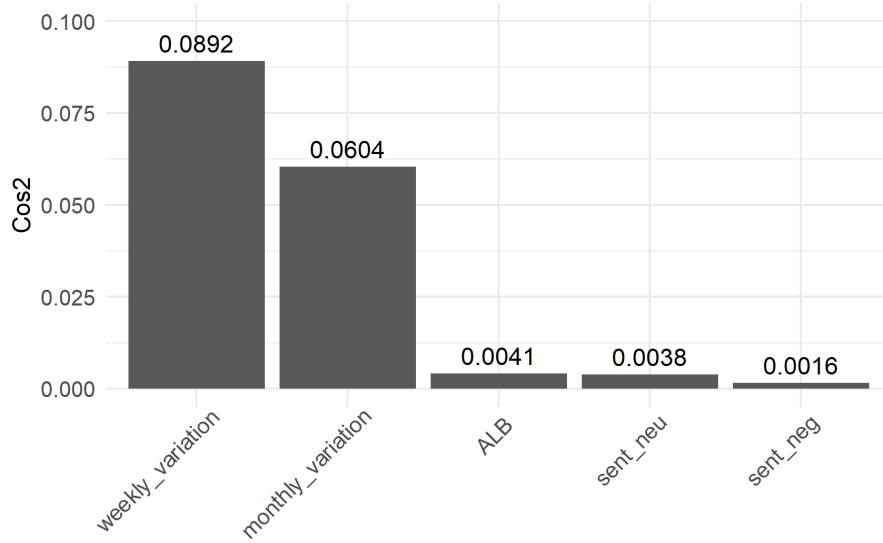


Figure 2: Variable representation plot of the first principal component.

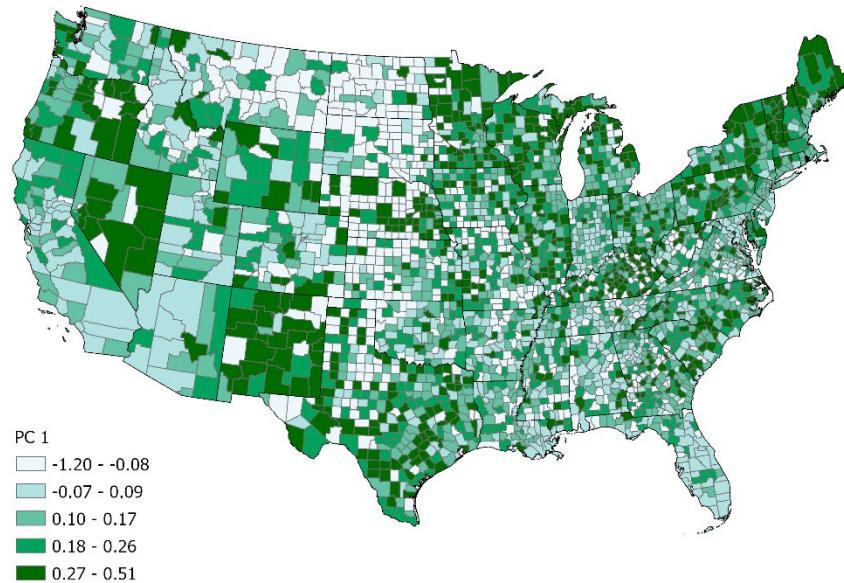


Figure 3: The spatial distribution of the first principal component.

The second principal component is dominated by the *ALB* variable, which has a cosine score of 0.1247 (Figure 4). We also find Twitter-derived variables and Google Search Interest among the top 5, but their cosine scores are substantially lower (*monthly_variation*: 0.0016; *weekly_variation*: 0.0016; *GST.interest*: 0.0004; *sent_neg*: 0.0001). Therefore, we call the second component of our anti-Asian area racism index the “Alien land bills (ALB)” component. The spatial distribution of the second component shows

variation predominantly at the state level (Figure 5). This is not surprising, given that *ALB*, the dominant variable in this component, is at the state level. The minor variation among counties of the same state exemplifies the other variables in this component, which are not irrelevant, but of lesser importance. We observe higher values in the states of Idaho, Montana, Utah, Oklahoma, Arkansas, Louisiana, Indiana, Tennessee, Alabama, Virginia, and Florida, while the states of Oregon, Nevada, Nebraska, Minnesota, Wisconsin, Kentucky, Pennsylvania, Connecticut, Massachusetts, Vermont, and Maine exhibit lower values.

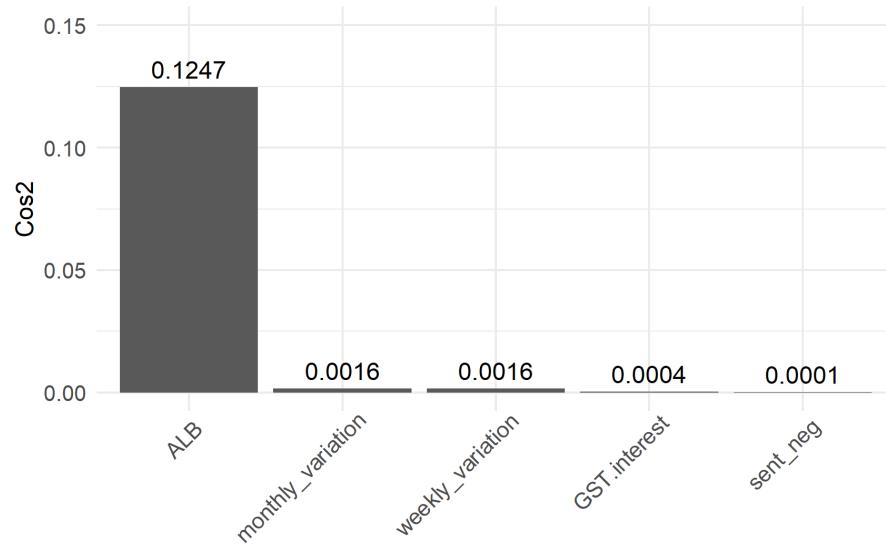


Figure 4: Variable representation plot of the second principal component.

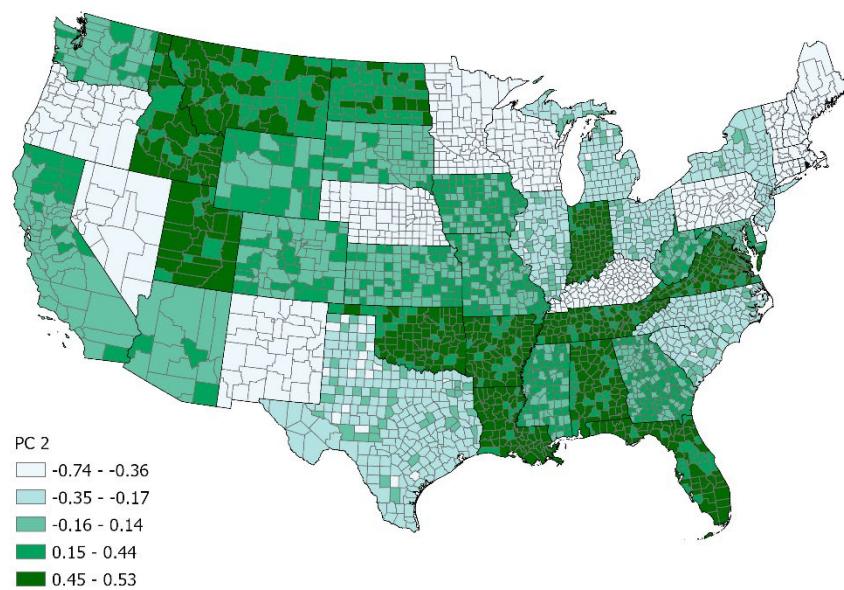


Figure 5: The spatial distribution of the second principal component.

The third principal component is dominated by *sent_neu* and *GST.interest* (Figure 6), with cosine scores of 0.0071 and 0.0054, respectively. The remaining variables among the top 5 are *weekly_variation*, *sent_compound*, and *monthly_variation*, but have substantially lower scores (0.0021, 0.0019, 0.0018, respectively). Therefore, we call this component the “Sentiment-Google Search Interest” (SGSI) component. The spatial distribution of this component is characterized by lower values in the west coast, northeast, and Florida, and higher values in parts of Texas, Mississippi, and the boundary region of Missouri, Illinois, Kentucky, and Tennessee (Figure 7).

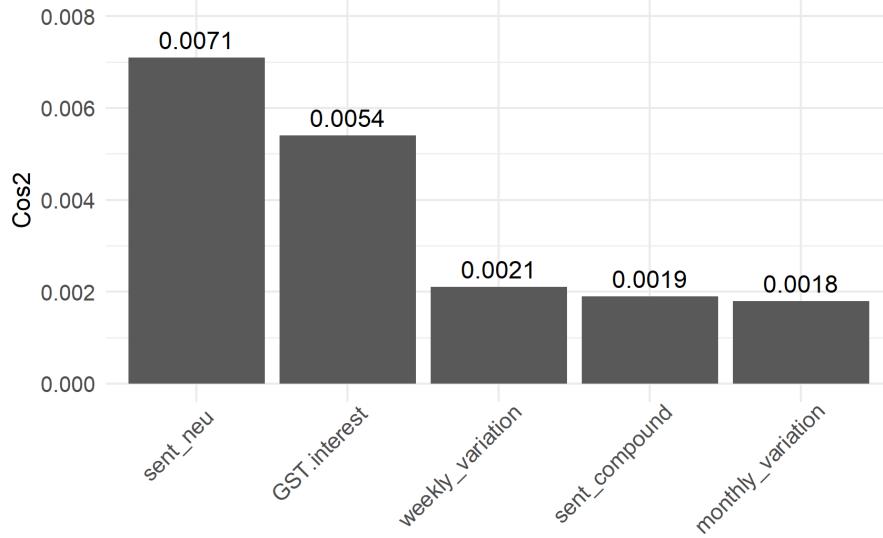


Figure 6: Variable representation plot of the third principal component.

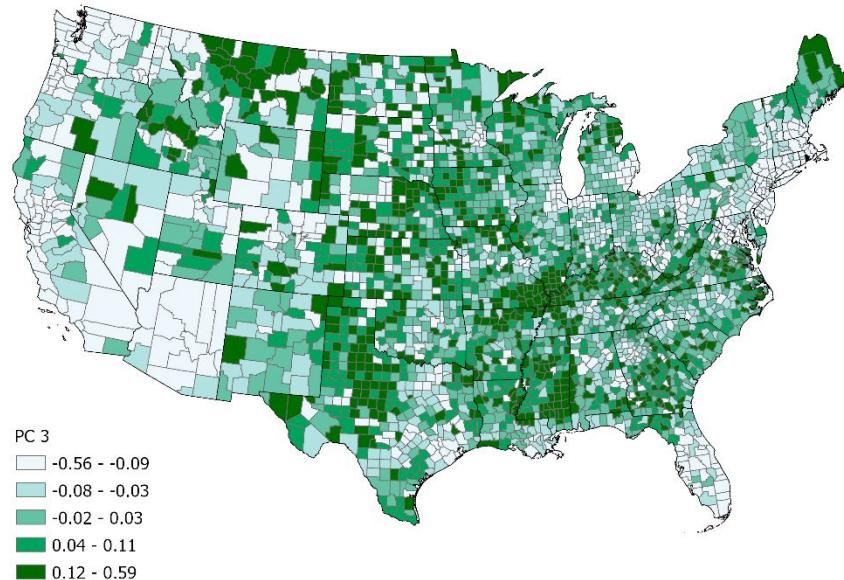


Figure 7: The spatial distribution of the third principal component.

3.2 Validation

Due to the clustered spatial distribution of HEAP survey responses and their zero-inflated pattern, we used 96 counties for validation. The weighted sum of experienced racism averaged over all counties, was 0.324. The regression modeling showed that the first (SMTV) and third (SGSI) dimensions of our anti-Asian racism index are significant predictors of experienced racism (Models 1 & 3, Table 2). This pattern was mirrored in the multiple regression model (Model 4), where the same dimensions were significant again. Model fit was generally low among univariate models ($0.0071 \leq \text{adjusted } R^2 \leq 0.0853$), but Model 3 performed best. Model 4 performed better (adjusted $R^2 = 0.1756$), and therefore explains a greater proportion of the variance in experienced racism.

Model	Variable	Estimate	Std. Error	Pr(> t)	Adjusted R ²
1	Intercept	0.0333	0.0181	0.0701	0.0741
	PC 1	0.4472	0.1524	0.0042 **	
2	Intercept	0.0697	0.0131	7.55e-07 ***	0.0071
	PC 2	-0.0520	0.0400	0.197	
3	Intercept	0.1301	0.0224	8.89e-08 ***	0.0853
	PC 3	0.3641	0.1159	0.0022 **	
4	Intercept	0.0839	0.0256	0.0014 **	0.1756
	PC 1	0.6160	0.1950	0.0021 **	
	PC 2	0.0494	0.0494	0.32	
	PC 3	0.3942	0.1104	0.0005 ***	

Table 2. Validation results. Significance codes: ‘***’: $p \leq 0.001$, ‘**’: $p \leq 0.01$

4. Discussion

In this study, we created a county-level anti-Asian area racism index for the contiguous USA. To that end, we collected data from four different sources (Twitter/X, Google Search Trends, FBI UCR, and APA Justice) and used PCA to reduce the initial 21 variables to 3 components, which capture the breadth of anti-Asian racism and hence, constitute our index. We provide important details on the amount of variation captured by each component, the representation of the initial variables, and the spatial distribution of all variables, as well as the resulting index. We validated the index against independent EAP survey data of experienced racism and found that they align to some extent. Lastly, we ensure the reproducibility of our work by providing all data and codes involved in the design of the index.

As seen in Section 3 (Results), our anti-Asian racism area index is driven by 1) social media temporal variation, 2) alien land bills, and 3) sentiment & Google search interest. These choices make intuitive sense and shed light on a variety of anti-Asian racism topics. On the other hand, anti-Asian hate crimes from the FBI UCR data did not receive substantial representation in our index. Hate crimes are extreme manifestations of hate and therefore, occur very sparsely. This is evident from their spatial distribution (Appendix Figures A17 – A19), which suspectedly caused the lack of representation in the anti-Asian area racism index. Perhaps, using a different data source, such as the National Crime Victimization Survey (Kena & Thompson, 2021) could address the issue of data sparseness. In addition, it is worth taking a closer look at the Twitter-derived variables. While variables related to the temporal variation of tweets and sentiment are represented in the resulting three dimensions, the Twitter-derived variables that quantify hate speech are not. This seems counterintuitive, but we suspect this result is due to the same reason: data sparseness. For instance, the spatial distributions of the *hateful_tweets_rate* and *hateful_tweets_prop* variables (Appendix Figure A2, A3) show a substantial number of counties that have a value of zero.

Our validation efforts show that the area-level anti-Asian racism index aligns with experienced racism captured by the HEAP survey to some degree. The associations were strongest for social media temporal variation and sentiment and Google search interest. Regional social media and Google search interest constitute unique indicators of an area's racial animus, and they are becoming increasingly relevant as internet penetrates every aspect of life (Stephens-Davidowitz, 2014). Analysis of geotagged data from social media or Google search is thus important for capturing the complex concept of anti-Asian racism in the USA and its geographic distributions.

This work has limitations: First, the data sources and derived variables that constitute the anti-Asian area racism index are subjectively chosen by the authors. We believe our choices are well justified, but we have not performed any form of validation so far. Such validation could mean comparing the index to other independent datasets of anti-Asian area racism, which we plan to do in future work. Second, Elon Musk's acquisition of Twitter and its renaming to X in October 2022 resulted in changes to content moderation policies, which were followed by an increase in hate speech on the platform (Benton et al., 2022). This exemplifies the downsides of some crowd-sourced data collection: Very few people, sometimes a single person, can decide on the direction of the platform, which can significantly impact the resulting data. This is true for the Google Search Trends data collection (Section 2.1.3) as well. Third, the keyword-based approach to identify hate speech ignores linguistic context and therefore, can misclassify tweets. For instance, the sentence "Don't call it Chinese virus, call it SARS-CoV-2" would be falsely classified as hateful due to the presence of the term "Chinese virus", which is in our list of hate terms. We plan to address this issue using more sophisticated methods to detect hate speech, such as deep learning-based large-language models (Alkomah & Ma, 2022).

Understanding how local contexts influence health outcomes is crucial to addressing health inequalities among marginalized groups in the USA. Researchers play a vital role in tackling barriers to racial equity, such as structural racism, by studying and intervening in its impact on minority health and well-being. This project contributes to the identification of high-risk areas for anti-Asian racism and provides evidence for place-based interventions to support vulnerable communities in combating racial injustice. Findings may guide the development of place-based interventions, like awareness campaigns, adjustments of communication strategies, hate crime prevention, and public safety resource allocation. Ultimately, structural interventions, such as increasing the political power of marginalized groups and changing public perceptions, are essential to addressing the root causes of structural racism. These interventions should be evidence-based and involve collaboration across various organizations and sectors.

5. Conclusions

We developed a composite anti-Asian racism index at the county level for the contiguous USA, valid for the years 2020 and 2021, based on the period of data collection. As this index is place-based, it can be easily linked to other place-based measures and observations, such as census data, hospital records, or crime data. We hope this data product will help identify associations between anti-Asian racism and various societal phenomena, such as health, policy, and the economy, ultimately aiming to improve the livelihoods of those affected by racism. We are sharing the codes and data used in the design of the index and encourage other researchers to utilize and refine our methods. Our index can be recalculated in future years if the necessary data continues to be available.

References

- Alhusen, J. L., Bower, K. M., Epstein, E., & Sharps, P. (2016). Racial discrimination and adverse birth outcomes: an integrative review. *Journal of Midwifery & Women's Health*, 61(6), 707–720.
- Alkomah, F., & Ma, X. (2022). A literature review of textual hate speech detection methods and datasets. *Information*, 13(6), 273.
- Anselin, L., Sridharan, S., & Gholston, S. (2007). Using exploratory spatial data analysis to leverage social indicator databases: the discovery of interesting patterns. *Social Indicators Research*, 82, 287–309.
- APA Justice. (2023a, October 10). *Alien Land Bills Map*.
https://www.apajustice.org/uploads/1/1/5/7/115708039/20231010_statebystatemapx1.jpg
- APA Justice. (2023b, October 20). *Alien Land Bills Table*.
https://www.apajustice.org/uploads/1/1/5/7/115708039/20231020_alienlandbillscan.pdf
- Bailey, Z. D., Feldman, J. M., & Bassett, M. T. (2021). How structural racism works—racist policies as a root cause of US racial health inequities. In *New England Journal of Medicine* (Vol. 384, Issue 8, pp. 768–773). Mass Medical Soc.
- Bailey, Z. D., Krieger, N., Agénor, M., Graves, J., Linos, N., & Bassett, M. T. (2017). Structural racism and health inequities in the USA: evidence and interventions. *The Lancet*, 389(10077), 1453–1463.
- Bastianoni, S., Pulselli, F. M., Focardi, S., Tiezzi, E. B. P., & Gramatica, P. (2008). Correlations and complementarities in data and methods through Principal Components Analysis (PCA) applied to the results of the SPIIn-Eco Project. *Journal of Environmental Management*, 86(2), 419–426.
- Benton, B., Choi, J.-A., Luo, Y., & Green, K. (2022). Hate speech spikes on Twitter after Elon Musk acquires the platform. *School of Communication and Media, Montclair State University*.
- Berger, M., & Sarnyai, Z. (2015). “More than skin deep”: stress neurobiology and mental health consequences of racial discrimination. *Stress*, 18(1), 1–10.
- Bitter, C., Mulligan, G. F., & Dall'erba, S. (2007). Incorporating spatial variation in housing attribute prices: a comparison of geographically weighted regression and the spatial expansion method. *Journal of Geographical Systems*, 9, 7–27.
- Centers for Disease Control and Prevention (CDC). (2024). *CDC Museum COVID-19 Timeline*.
<https://www.cdc.gov/museum/timeline/covid19.html>
- Cobbinah, S. S., & Lewis, J. (2018). Racism & Health: A public health perspective on racial discrimination. *Journal of Evaluation in Clinical Practice*, 24(5), 995–998.
- Committee of 100. (2022). *We Belong: A Glossary of anti-Asian Terms and Tropes*.
https://www.committee100.org/wp-content/uploads/2022/02/Anti-Asian-Glossary-v2_Updated-Late-2022.pdf
- Demšar, U., Harris, P., Brunsdon, C., Fotheringham, A. S., & McLoone, S. (2013). Principal component analysis on spatial data: an overview. *Annals of the Association of American Geographers*, 103(1), 106–128.
- DuBois, W. E. B. (2003). The health and physique of the Negro American. *American Journal of Public Health*, 93(2), 272–276.
- Federal Bureau of Investigation. (2021). *Hate Crimes*. <https://www.fbi.gov/investigate/civil-rights/hate-crimes>
- Gee, G. C., & Ford, C. L. (2011). Structural racism and health inequities: Old issues, New Directions 1. *Du Bois Review: Social Science Research on Race*, 8(1), 115–132.
- Gelman, A., Fagan, J., & Kiss, A. (2007). An analysis of the New York City police department’s “stop-and-frisk” policy in the context of claims of racial bias. *Journal of the American Statistical Association*, 102(479), 813–823.
- Groos, M., Wallace, M., Hardeman, R., & Theall, K. P. (2018). Measuring inequity: a systematic review of methods used to quantify structural racism. *Journal of Health Disparities Research and Practice*, 11(2), 13.
- Hohl, A., Choi, M., Medina, R., Wan, N., & Wen, M. (2024). COVID-19: adverse population sentiment and place-based associations with socioeconomic and demographic factors. *Spatial Information Research*, 32(1), 73–84.

- Hohl, A., Choi, M., Yellow Horse, A. J., Medina, R. M., Wan, N., & Wen, M. (2022). Spatial distribution of hateful tweets against Asians and Asian Americans during the COVID-19 pandemic, November 2019 to May 2020. *American Journal of Public Health*, 112(4), 646–649.
- Horse, A. J. Y. (2021). Anti-Asian racism, xenophobia and Asian American health during COVID-19. In *The COVID-19 Crisis* (pp. 195–206). Routledge.
- Hswen, Y., Xu, X., Hing, A., Hawkins, J. B., Brownstein, J. S., & Gee, G. C. (2021). Association of “# covid19” versus “# chinesevirus” with anti-Asian sentiments on Twitter: March 9–23, 2020. *American Journal of Public Health*, 111(5), 956–964.
- Huang, X., Wang, S., Yang, D., Hu, T., Chen, M., Zhang, M., Zhang, G., Biljecki, F., Lu, T., & Zou, L. (2024). Crowdsourcing Geospatial Data for Earth and Human Observations: A Review. *Journal of Remote Sensing*, 4, 0105.
- Hutto, C., & Gilbert, E. (2014). Vader: A parsimonious rule-based model for sentiment analysis of social media text. *Proceedings of the International AAAI Conference on Web and Social Media*, 8(1), 216–225.
- Huynh, Q.-L., Devos, T., & Smalarz, L. (2011). Perpetual foreigner in one’s own land: Potential implications for identity and psychological adjustment. *Journal of Social and Clinical Psychology*, 30(2), 133–162.
- Iwamoto, D. K., & Liu, W. M. (2010). The impact of racial identity, ethnic identity, Asian values, and race-related stress on Asian Americans and Asian international college students’ psychological well-being. *Journal of Counseling Psychology*, 57(1), 79.
- Izenman, A. J. (2008). *Modern multivariate statistical techniques* (Vol. 1). Springer.
- Jolliffe, I. (2014). Principal Component Analysis. In *Wiley StatsRef: Statistics Reference Online*. <https://doi.org/https://doi.org/10.1002/9781118445112.stat06472>
- Jordan, S. E., Hovet, S. E., Fung, I. C.-H., Liang, H., Fu, K.-W., & Tse, Z. T. H. (2018). Using Twitter for public health surveillance from monitoring and prediction to public response. *Data*, 4(1), 6.
- Kassambara, A. (2017). *Practical guide to principal component methods in R: PCA, M (CA), FAMD, MFA, HCPC, factoextra* (Vol. 2). Sthda.
- Kena, G., & Thompson, A. (2021). *Hate Crime Victimization, 2005–2019*. https://bjs.ojp.gov/sites/g/files/xyckuh236/files/media/document/hcv0519_1.pdf
- Krieger, N., Waterman, P. D., Kosheleva, A., Chen, J. T., Carney, D. R., Smith, K. W., Bennett, G. G., Williams, D. R., Freeman, E., & Russell, B. (2011). Exposing racial discrimination: implicit & explicit measures—the my body, my story study of 1005 US-born black & white community health center members. *PloS One*, 6(11), e27636.
- Lewith, G. T., Jonas, W. B., & Walach, H. (2010). *Clinical research in complementary therapies: Principles, problems and solutions*. Elsevier Health Sciences.
- Li, Z., Huang, X., Ye, X., Jiang, Y., Martin, Y., Ning, H., Hodgson, M. E., & Li, X. (2021). Measuring global multi-scale place connectivity using geotagged social media data. *Scientific Reports*, 11(1), 14694.
- Martín, Y., Cutter, S. L., Li, Z., Emrich, C. T., & Mitchell, J. T. (2020). Using geotagged tweets to track population movements to and from Puerto Rico after Hurricane Maria. *Population and Environment*, 42(1), 4–27.
- Müller, K., & Schwarz, C. (2023). From hashtag to hate crime: Twitter and antiminority sentiment. *American Economic Journal: Applied Economics*, 15(3), 270–312.
- Nguyen, T. T., Criss, S., Dwivedi, P., Huang, D., Keralis, J., Hsu, E., Phan, L., Nguyen, L. H., Yardi, I., & Glymour, M. M. (2020). Exploring US shifts in anti-Asian sentiment with the emergence of COVID-19. *International Journal of Environmental Research and Public Health*, 17(19), 7032.
- Nguyen, T. T., Huang, D., Michaels, E. K., Glymour, M. M., Allen, A. M., & Nguyen, Q. C. (2021). Evaluating associations between area-level Twitter-expressed negative racial sentiment, hate crimes, and residents’ racial prejudice in the United States. *SSM-Population Health*, 13, 100750.

- Parinet, B., Lhote, A., & Legube, B. (2004). Principal component analysis: an appropriate tool for water quality evaluation and management—application to a tropical lake system. *Ecological Modelling*, 178(3–4), 295–311.
- Pearson, K. (1901). LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11), 559–572.
- Sabharwal, M., Becerra, A., & Oh, S. (2022). From the Chinese Exclusion Act to the COVID-19 pandemic: A historical analysis of “otherness” experienced by Asian Americans in the United States. *Public Integrity*, 24(6), 535–549.
- Saito, N. T. (1997). Model minority, yellow peril: Functions of foreignness in the construction of Asian American legal identity. *Asian LJ*, 4, 71.
- Singleton, A., Spielman, S., & Folch, D. (2018). *Urban Analytics* (1st ed.). Sage.
- Stephens-Davidowitz, S. (2014). The cost of racial animus on a black candidate: Evidence using Google search data. *Journal of Public Economics*, 118, 26–40.
- Stop AAPI Hate. (2023). *Righting wrongs: how civil rights can protect asian Americans & pacific islanders against racism*. <https://stopaapihate.org/wp-content/uploads/2023/09/23-SAH-CivilRightsReport-F.pdf>
- Strobl, C. (2008). PostGIS. In S. Shekhar & H. Xiong (Eds.), *Encyclopedia of GIS* (pp. 891–898). Springer US. https://doi.org/10.1007/978-0-387-35973-1_1012
- Su, D., Alshehri, K., Ern, J., Chen, B., Chen, L., Chen, Z., Han, X., King, K. M., Li, H., & Li, J. (2022). Racism experience among American adults during COVID-19: a mixed-methods study. *Health Equity*, 6(1), 554–563.
- Tessler, H., Choi, M., & Kao, G. (2020). The anxiety of being Asian American: Hate crimes and negative biases during the COVID-19 pandemic. *American Journal of Criminal Justice*, 45, 636–646.
- The United States Department of Justice. (2020). *Hate Crime Statistics*. <https://www.justice.gov/crs/highlights/2020-hate-crimes-statistics>
- Tran, L. T., Knight, C. G., O'Neill, R. V., Smith, E. R., Riitters, K. H., & Wickham, J. (2002). Fuzzy decision analysis for integrated environmental vulnerability assessment of the Mid-Atlantic region. *Environmental Management*, 29, 845–859.
- United States Census Bureau. (2020). *2020 TIGER/Line Shapefiles (machine readable data files)*. U.S. Department of Commerce. www.census.gov/cgi-bin/geo/shapefiles/index.php
- United States Department of Justice. Federal Bureau of Investigation. (2023). *Uniform Crime Reporting Program Data: Hate Crime Data (Record-Type Files)*, United States, 2020. Inter-University Consortium for Political and Social Research [Distributor], 2023-12-11. <https://doi.org/https://doi.org/10.3886/ICPSR38790.v1>
- Wu, A. (2023). Perpetual Foreigner Stereotype: Third Class Americans. In *Asian American Educators and Microaggressions: More Than Just Work (ers)* (pp. 41–60). Springer.

Appendix

Source	Hohl et al. 2022	hatebase.org	Committee of 100
Keywords	Asian Invasion AsianInvasion Bamboo coon bambooocoon Blame China BlameChina Bomb China BombChina CCP is Virus CCPis Virus CCP isVirus CCPisVirus CCP Virus CCPVirus china did this China Didthis ChinaDid this ChinaDidthis China Is Terrorist ChinalsTerrorist China IsTerrorist Chinals Terrorist China lied people died ChinaLiedPeopleDied China LiedPeopleDied ChinaLied PeopleDied ChinaLiedPeople Died ChinaLied People Died China LiedPeople Died China Lied PeopleDied chinaman china man china men chinamen china virus chinavirus Chinazi ChineseBioterrorism Chinese Bioterrorism Chinese Eat Bats ChineseEatBats ChineseEat Bats Chinese EatBats Chinese Is Virus ChinesIsVirus Chinese IsVirus ChineseIs Virus chinese virus	dog eater dog eaters dogeater dogeaters gink ginks goloid goloids gook eyed gookeyed gookette gookettes gookie gookies gooklet gooklets gooky eyes gookyeyes mongoloid mongoloids rice nigger rice niggers ricensigger ricensiggers whoriental whorientals yellow invader yellow invaders yellowinvader yellowinvaders	asian driver asiandriver chinese driver chinesedriver banana twinkie chinaman chinesefire drill chinesefire drill chinesefiredrill chinesefiredrill chinese restaurant syndrome chineserestaurant syndrome chinese restaurantsyndrome chineserestaurantsyndrome dragon lady dragonlady fu manchu fumanchu lotus blossom lotusblossom china doll chinadoll oriental rice burner riceburner ricer rice machine ricemachine sick man of asia sickmanofasia tiger mom tigermom zipper head zipperhead

	chinese virus chinesevirus Ching Chong Ching Chongs ChingChong ChingChongs chinig chink chinks chinky chonky Chop Fluey ChopFluey churka cina cokin communist virus communistvirus coolie dink Fuck China FuckChina Hold China Accountable HoldChinaAccountable Hold ChinaAccountable HoldChina Accountable kung flu Kung Fu Virus kungflu kung-flu KungFu Virus Kung-Fu Virus KungFuVirus Kung-FuVirus Make China Pay MakeChinaPay Make ChinaPay MakeChina Pay niakoue oriental pastelde flango pastel deflango pasteldeflango pastel de flango ricerabies rice rabies sidewayscooter sideways cooter sidewayspussies sidewayspussy sideways pussies		
--	---	--	--

	sideways pussy sidewaysvagina sidewaysvaginas sideways vagina sideways vaginas slant slanteye slant eye slope head slopehead spink spinks ting tong tingtong wuhan flu wuhan virus wuhanflu wuhan-flu wuhanvirus		
--	--	--	--

Table A1: Anti-Asian hate terms.

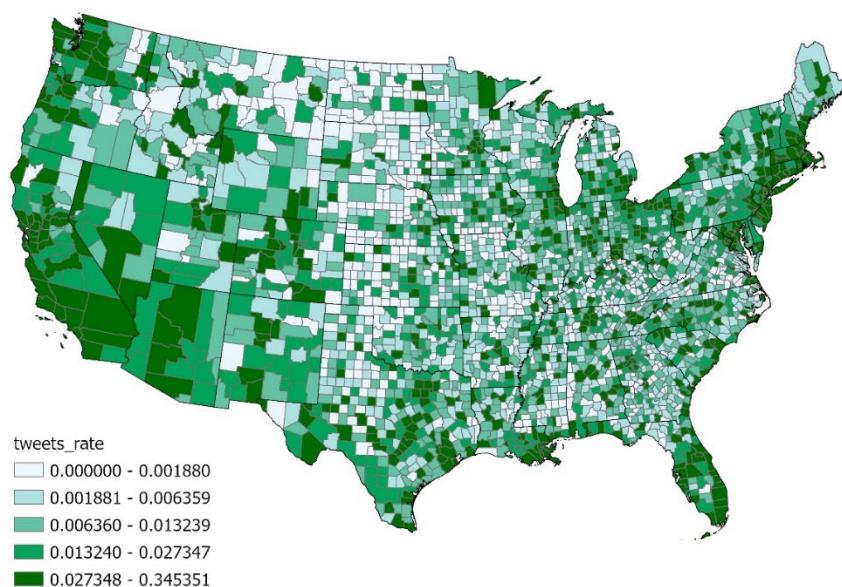


Figure A1: The spatial distribution of the *tweets_rate* variable.

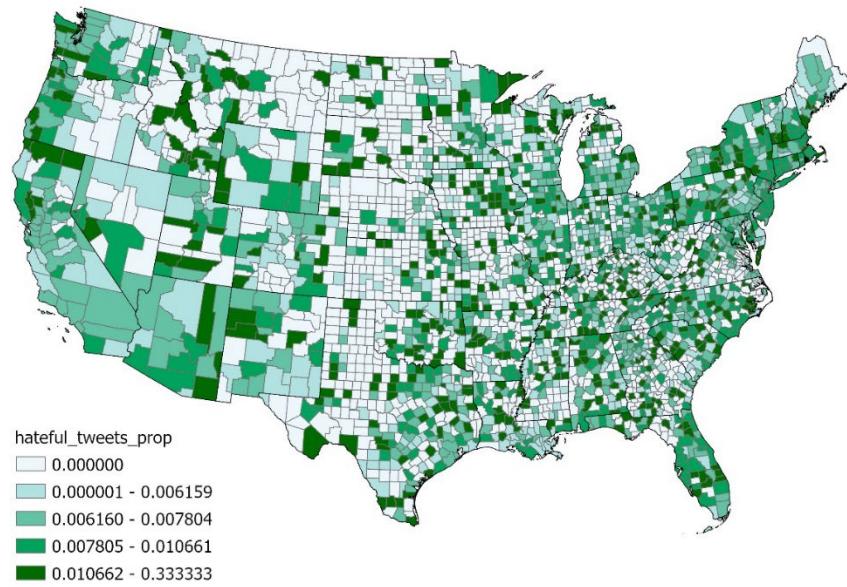


Figure A2: The spatial distribution of the *hateful_tweets_prop* variable.

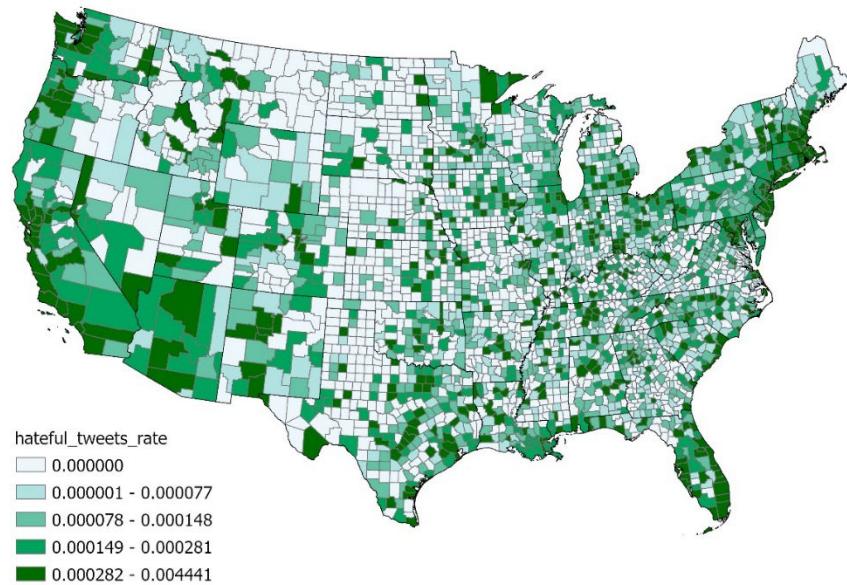


Figure A3: The spatial distribution of the *hateful_tweets_rate* variable.

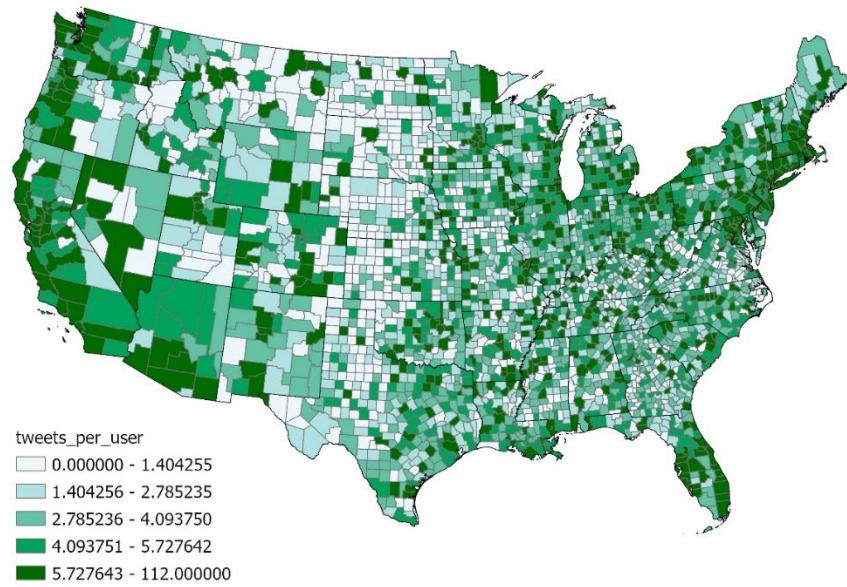


Figure A4: The spatial distribution of the *tweets_per_user* variable.

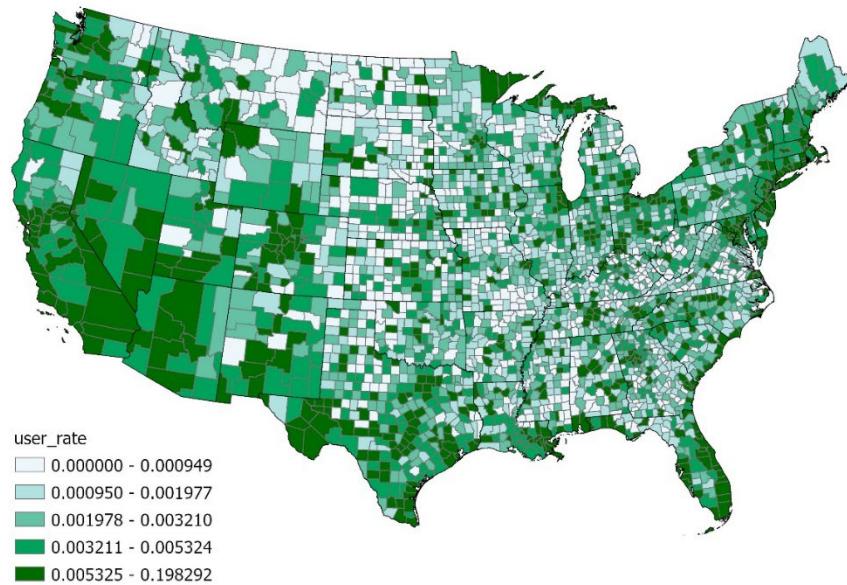


Figure A5: The spatial distribution of the *user_rate* variable.

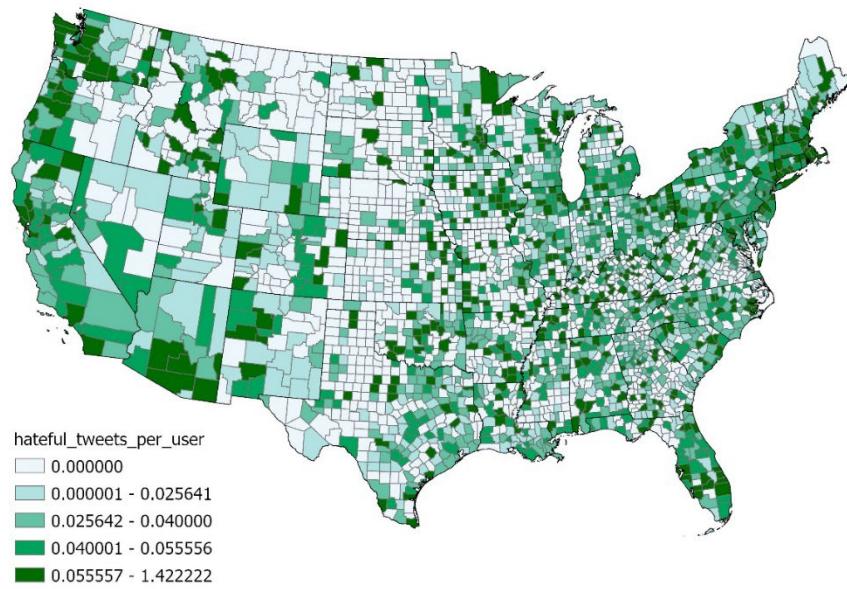


Figure A6: The spatial distribution of the *hateful_tweets_per_user* variable.

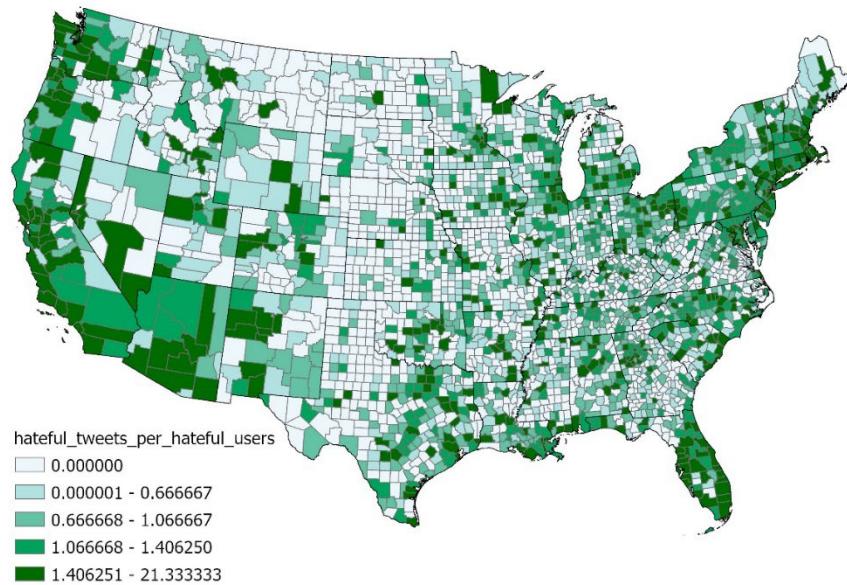


Figure A7: The spatial distribution of the *hateful_tweets_per_hateful_user* variable.

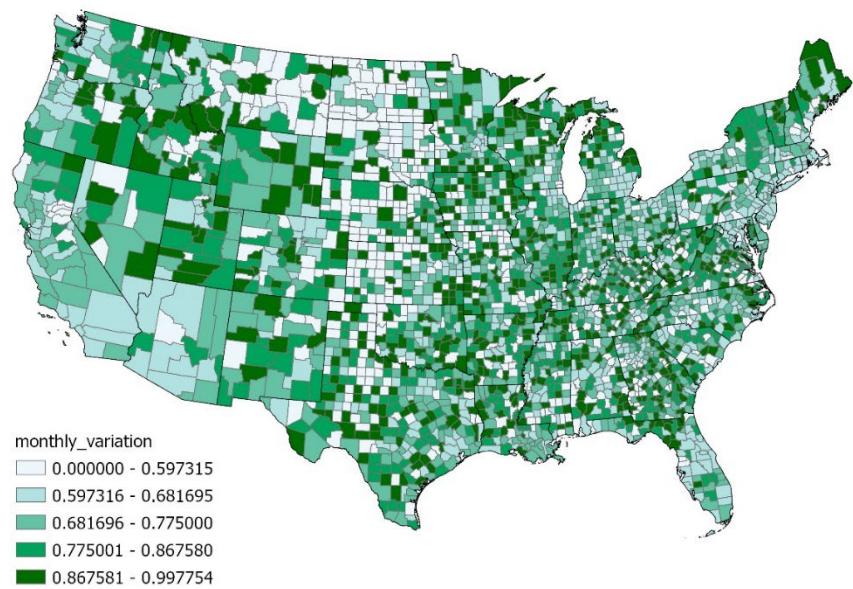


Figure A8: The spatial distribution of the *monthly_variation* variable.

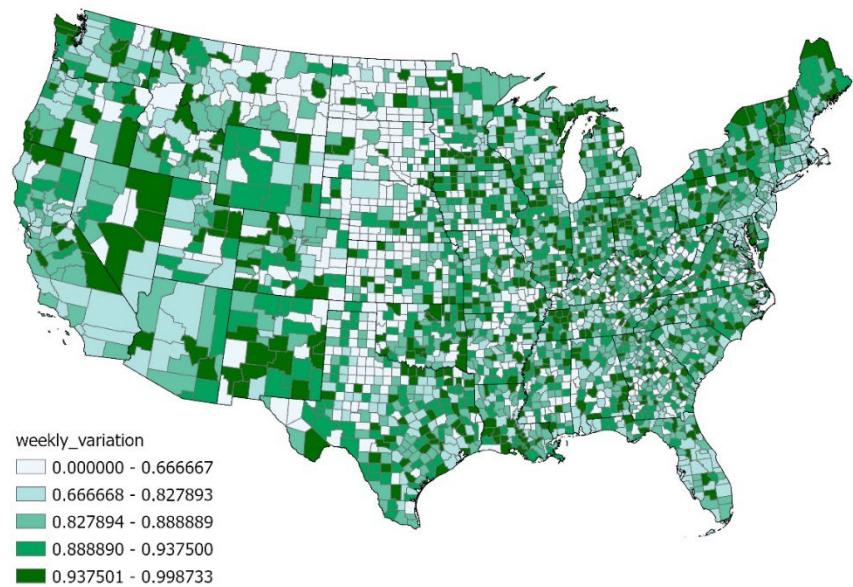


Figure A9: The spatial distribution of the *weekly_variation* variable.

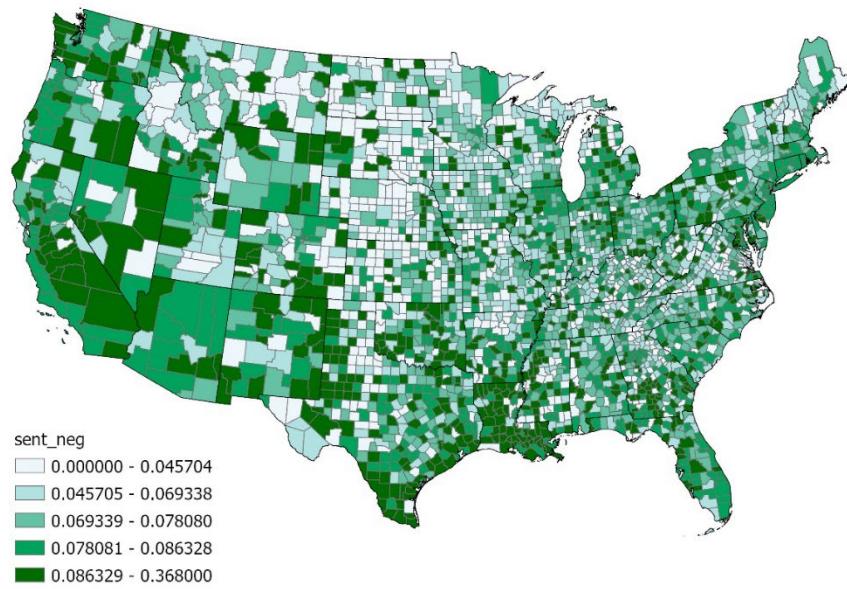


Figure A10: The spatial distribution of the *sent_neg* variable.

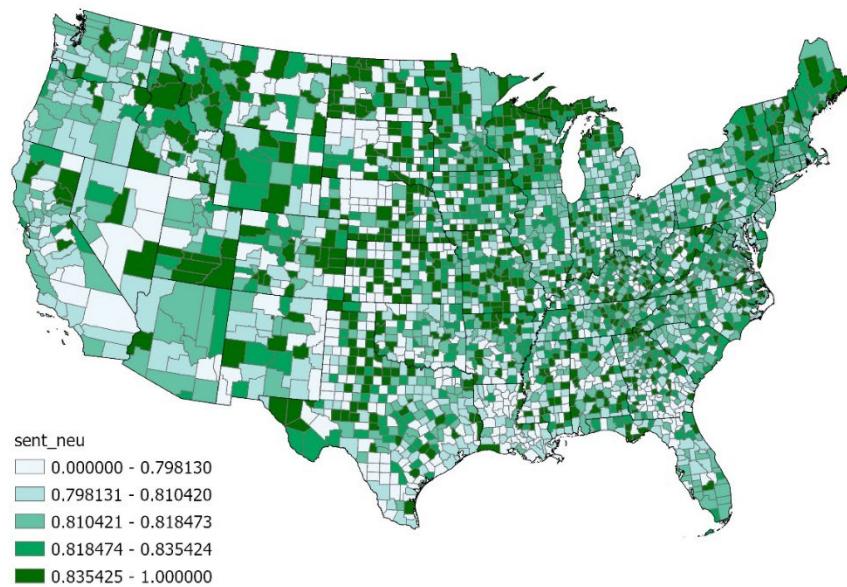


Figure A11: The spatial distribution of the *sent_neu* variable.

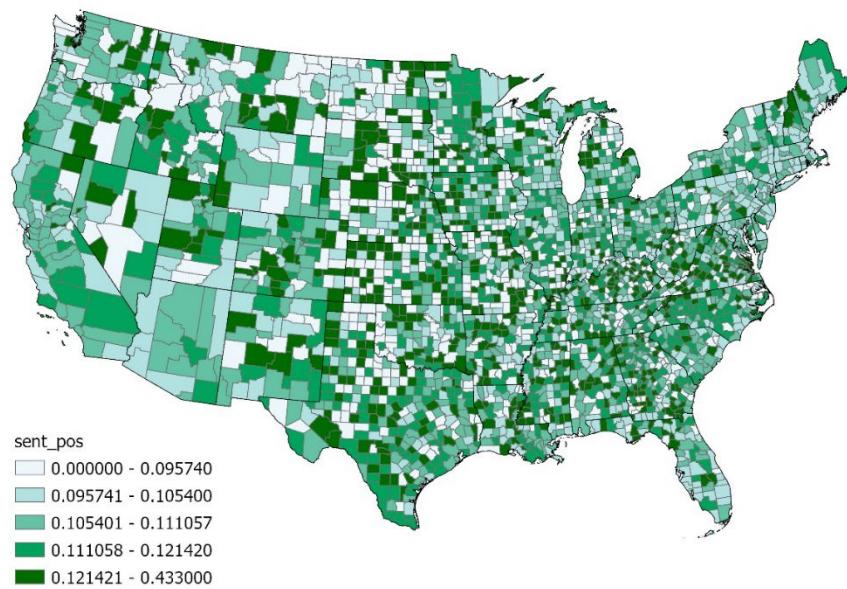


Figure A12: The spatial distribution of the *sent_pos* variable.

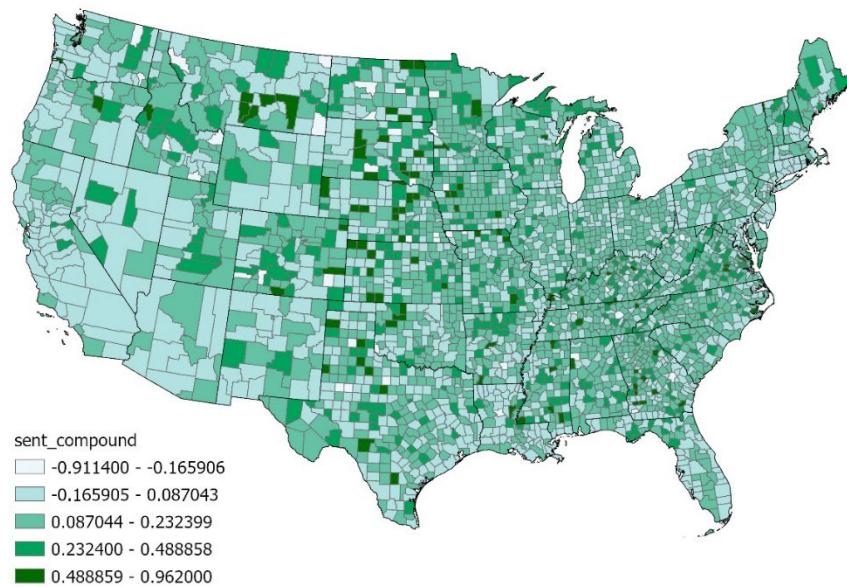


Figure A13: The spatial distribution of the *sent_compound* variable.

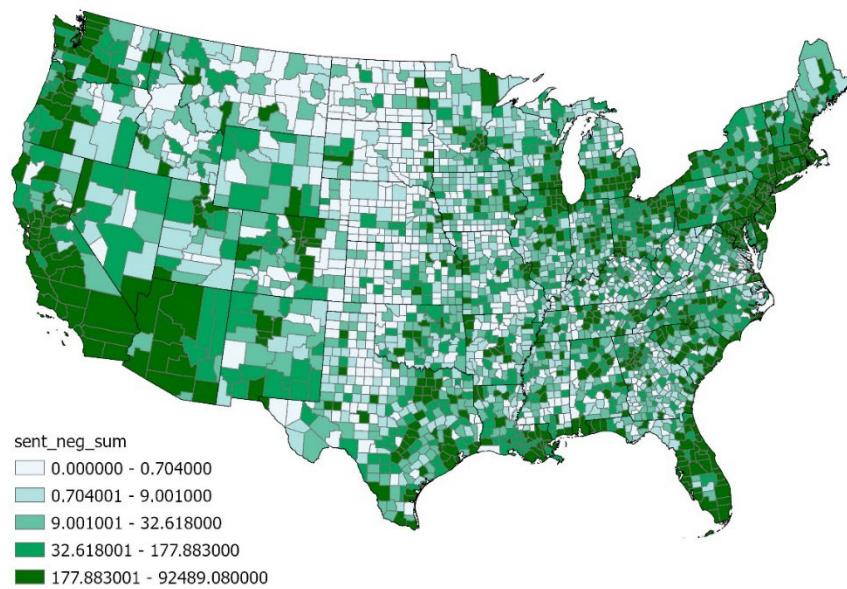


Figure A14: The spatial distribution of the *sent_neg_sum* variable.

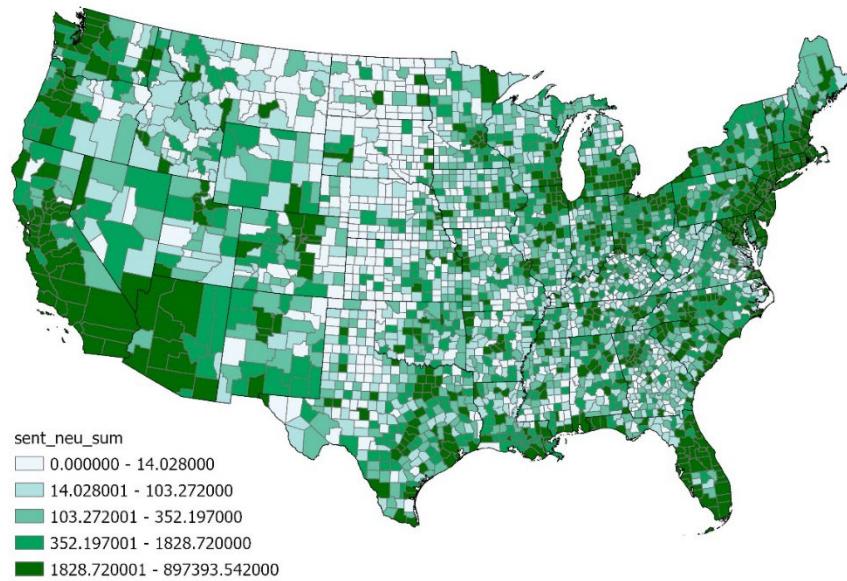


Figure A15: The spatial distribution of the *sent_neu_sum* variable.

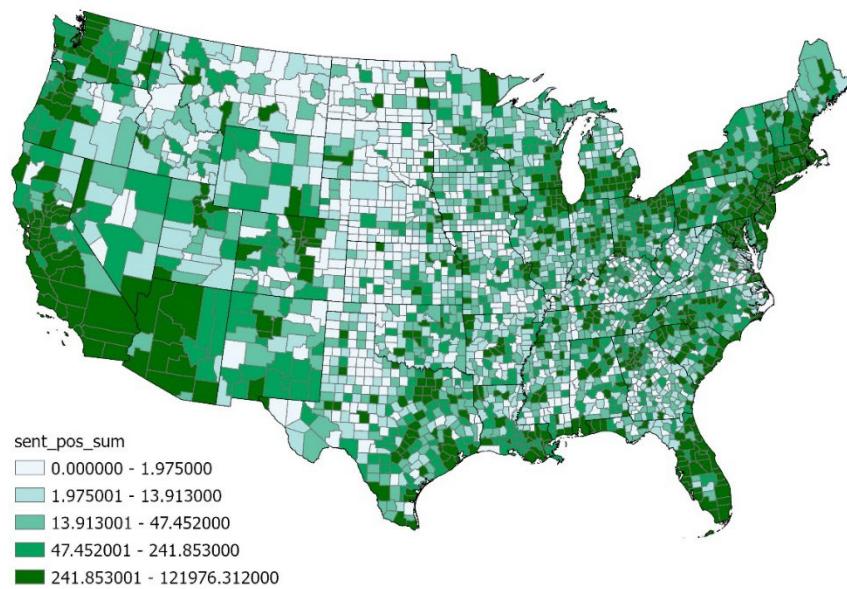


Figure A16: The spatial distribution of the *sent_pos_sum* variable.

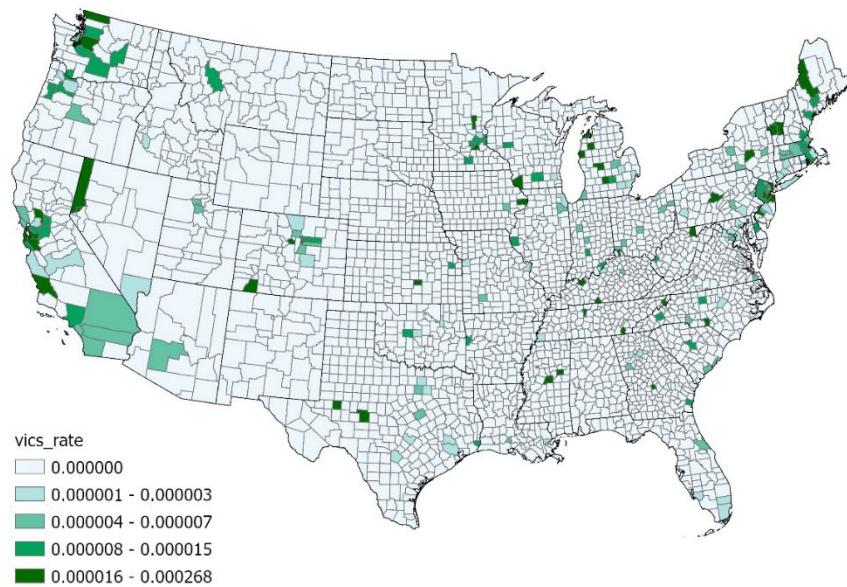


Figure A17: The spatial distribution of the *vics_rate* variable.

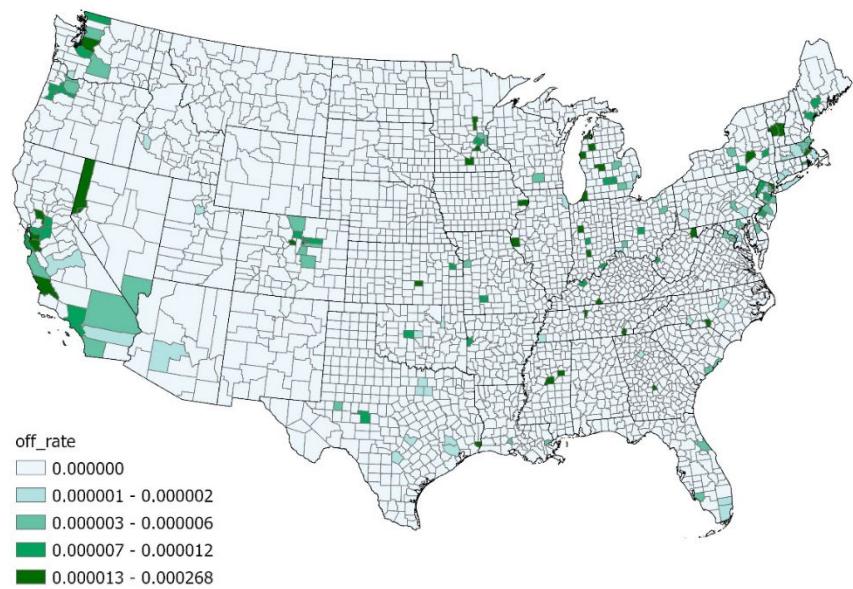


Figure A18: The spatial distribution of the *off_rate* variable.

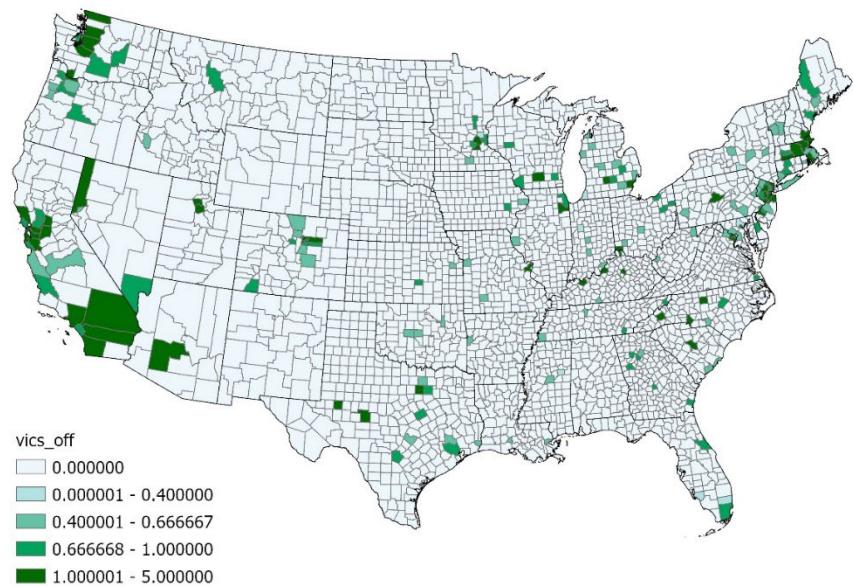


Figure A19: The spatial distribution of the *vics_off* variable.

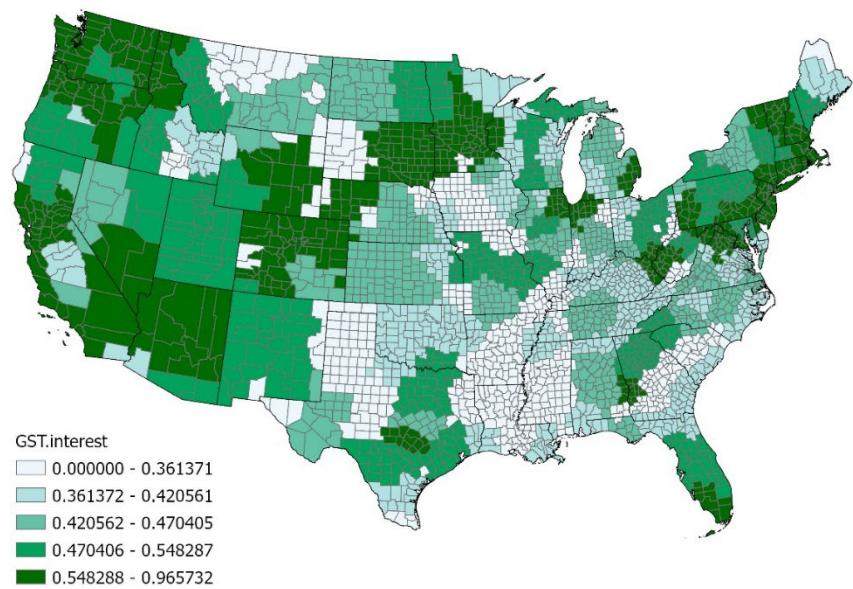


Figure A20: The spatial distribution of the *GST.interest* variable.

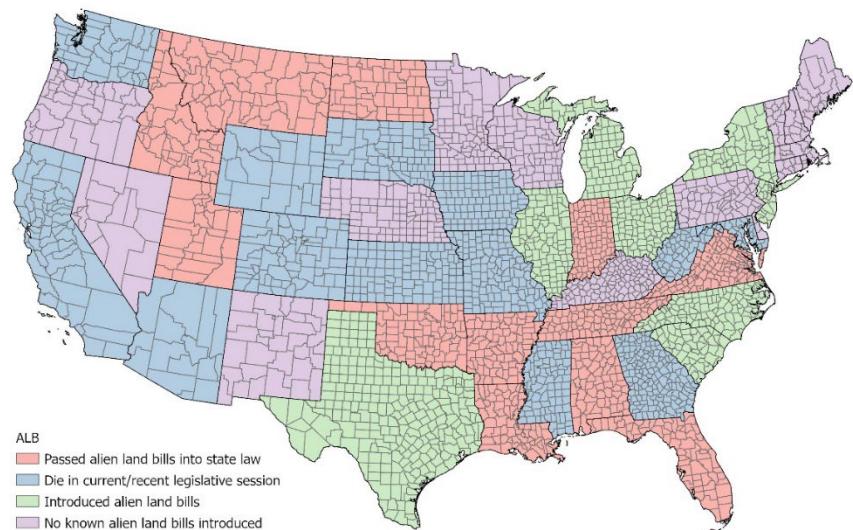


Figure A21: The spatial distribution of the *ALB* variable.

