

# **Project Report - Mental Health on the Map**

Shefali Verma, Tessa Joseph, Anvi Rakshe



## **Abstract**

This project seeks to examine the underlying factors contributing to mental health distress within the American population, focusing on the data from 2022 covering all 50 states. Through the analysis of state-wide health metrics provided by America's Health Rankings, we aim to find patterns shared among states that report high levels of mental health issues. This project also aims to find a way to identify at-risk communities that might be underreporting mental health concerns due to cultural or stigma-related barriers. By grouping communities based on mental health outcomes and correlating factors, such as location and environmental conditions, we aspire to understand the broader influences on mental health.

## **Introduction**

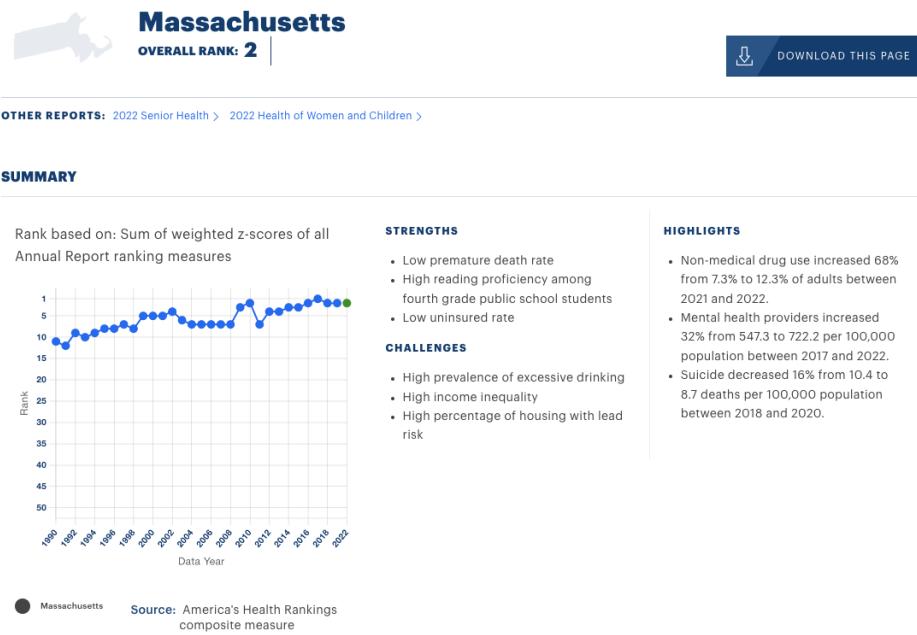
Mental well-being often takes a back seat in health discussions, primarily due to the discomfort and stigma associated with it. Nevertheless, it is necessary to prioritize mental health issues as a critical public health concern, given their status as leading causes of morbidity and mortality globally. While existing interventions such as medications and therapy offer relief to those already grappling with mental health problems, this reactive approach only addresses individual cases.

As medical professionals and researchers pursue cures, there is a notable lack of emphasis on the proactive prevention of mental health challenges at a broader level. Adding to this issue, certain communities may harbor stigmas surrounding mental health, leading to potential underreporting or even completely disregarding cases of mental distress within their population.

Through a comprehensive examination of health and socioeconomic characteristics across American states, coupled with Frequent Mental Distress data, there may be an opportunity to predict the likelihood of mental health distress in communities where underreporting of mental health issues may be a concern.

We ask ourselves, what are the indicators of mental distress within a population? And how can we identify other communities with low positive mental-health outcomes?

We scraped data from [America's Health Rankings.org](#). This site offers in-depth health-related data of each state in the United States. It not only describes the States' mental health data, like their average number of poor mental health days per month, or percentage of adults who experience frequent mental distress, but also contains comprehensive data about the State's health behaviors, access to clinical care, social and economic data, and physical environment/infrastructural data.



Measure	Rating	Value	Rank
<b>Social and Economic Factors *</b>			
Community and Family Safety - Annual *	++++	0.542	4
Firearm Deaths †	++++	0.651	8
Occupational Fatalities	++++	3.8	2
Public Health Funding	++++	3.3	6
Violent Crime	+++	\$185	7
Economic Resources - Annual *	++	309	18
Economic Hardship Index	++++	0.335	23
Crowded Housing †	++	25	4
Dependency (Ages <18 or >64) †	++++	2.2%	24
Less Than High School Education †	+++	36.9%	2
Per Capita Income †	++++	8.9%	26
Poverty †	++++	\$49,746	1
Unemployment †	+++	11.5%	17
Food Insecurity	++++	6.6%	38
Homeownership †	++	8.4%	9
Homeownership Racial Disparity †	++	63.2%	44
...	...	39.2	33

## Data Description

Our analysis revolves around a comprehensive dataset that looks at various factors potentially affecting mental health across all 50 U.S. states. The data, sourced from America's Health Rankings Website for the year 2022, includes a wide spectrum of indicators that may have implications on the prevalence of mental health issues at the state level.

For scraping the America's Health Rankings Website:

We use the function: `scrape_state_data(state_code)`, which builds a dataframe for each state based on the state code put through the function. This is the only function we have for scraping the website. The dataframes for each state will then be concatenated and cleaned:

We first use a dictionary with all the states and respective state codes to generate 50 dataframes using the `scrape_state_data` function. Then we combined all 50 dataframes together to make one comprehensive one. Then we clean the dataframe, getting rid of rows with missing values, making all values numerical, etc.

This is what the a part of the final data frame looks like:

Category	AL	AK	AZ	AR	CA	CO	CT	DE	FL	GA	HI	ID	IL	IN	IA	KS	KY
Social and Economic Factors	-0.377	-0.286	-0.177	-0.742	-0.126	0.299	0.438	0.265	0.146	-0.051	0.397	-0.035	0.059	-0.067	0.491	-0.006	-0.294
Community and Family Safety - Annual	-0.398	-0.909	-0.426	-1.070	0.077	0.027	0.729	-0.008	-0.205	-0.299	0.924	0.631	-0.028	-0.534	0.196	-0.538	0.063
Firearm Deaths	23.400	24.000	17.100	22.500	8.700	15.700	6.100	14.200	13.900	17.700	3.500	18.100	13.900	17.300	11.200	17.100	20.300
Occupational Fatalities	5.700	8.600	4.000	5.900	3.200	3.600	3.600	4.300	4.500	5.500	4.100	4.300	3.400	6.900	5.600	6.000	5.900
Public Health Funding (\$)	129.000	449.000	79.000	128.000	138.000	127.000	126.000	152.000	79.000	107.000	241.000	172.000	109.000	76.000	161.000	87.000	110.000
Violent Crime	454.000	838.000	485.000	672.000	442.000	423.000	182.000	432.000	384.000	400.000	254.000	243.000	426.000	358.000	304.000	425.000	259.000
																	6

The components of our dataset are:

- Frequent Mental Distress Score: This score is represented as the percentage of adults in the state who reported their mental health was not good for more than 14 days in the past month. We will use this as the primary indicator of mental health in a population.
- Health Distress Indicators: We look at data that reflects the frequency and intensity of mental distress experienced by the population.
- Healthcare Access: Parameters evaluating the availability and quality of mental healthcare services, including the number of mental health providers per capita and the percentage of residents with health insurance coverage.
- Socioeconomic Factors: Data points considering the impact of employment rates, education levels, and income inequality on mental health.
- Lifestyle and Community Metrics: Variables assessing the role of physical activity, diet, substance abuse, and community support in mental health.
- Environmental Influences: Information on factors such as air quality and urbanization that correlate with mental health outcomes.
-

# Pipeline Overview

Data retrieval from the America's Health Rankings Website:

- Our `scrape\_state\_data(state\_code)` function builds a data frame for each state based on the two-letter state code put through the function.

The Big DataFrame: Generation and Concatenation

- Then, we use a dictionary with all the states and respective state codes to generate 50 data frames using the `scrape\_state\_data` function. Then we combine the data frames together to make one comprehensive one.
- The third code cleans the data frame, getting rid of rows with missing values, making all values numerical, etc.

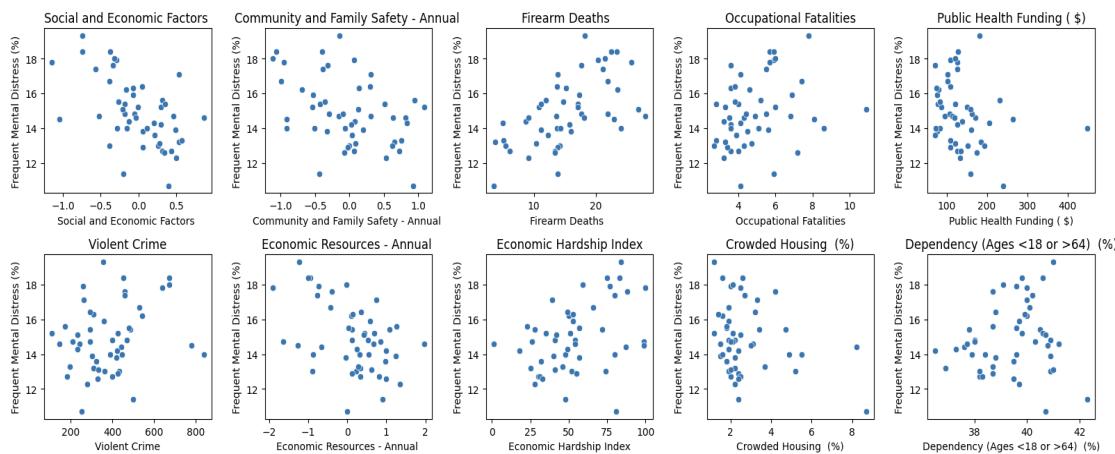
Cleaning the Dataframe:

- We changed missing values to NaN, get rid of rows with too many missing values, made all values numerical by removing "\$" and "%" symbols, and removed unnecessary symbols like "†" and "\*" from feature names.

## Method

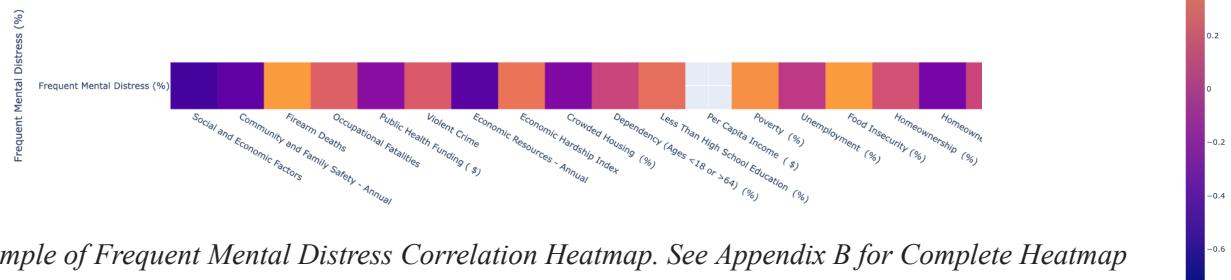
We are left with an extremely large data frame with over 100 health factors which may or may not correlate with frequent mental distress. Our goal is to determine which health factors actually do have a correlation with mental distress and isolate them. Then, run those factors through a machine learning model to be able to predict a frequent mental distress value from the other factors. Finally, we wanted to be able to run example values throughout our model, and receive a fairly accurate predicted Frequent mental distress Score.

Our first step was visualizing the relationship between each factor and frequent mental distress. This would give us a general idea of whether there are factors which are correlated.



*Fig. 1 Sample of Scatter-plots Displaying the Relationship Between all Given Health Metrics and Frequent Mental Distress. See Appendix A for Complete Scatterplot Matrix.*

We then constructed a correlation matrix using *scaled* values to determine which health metrics actually have any relationship with Frequent Mental Distress. This prevents any single variable from disproportionately affecting outcomes. We only considered those with a correlation of  $>|0.50|$ .



*Fig 2. Sample of Frequent Mental Distress Correlation Heatmap. See Appendix B for Complete Heatmap Visualization.*

Finally, we used the K-nearest neighbors method to model the relationship between our chosen features and predict mental distress outcomes from their values. Our data has already been scaled, reducing error and ensuring that each metric will contribute proportionally to the prediction model. The K-nearest neighbors method is suitable for our data because K-nearest neighbors takes into account the proximity of data points and by identifying and analyzing patterns among various health metrics, the method will effectively capture the factors that may play a role in mental distress. Our final function `predict_mental_distress(selected_feature_values)` will take any values for our chosen health metrics, and return a predicted frequent mental distress score. This function explores various scenarios by inputting different combinations of health metric values, enabling an adaptable analysis.

	True Values	Predicted Values
Index		
IA	13.9	13.80
WA	15.4	13.96
WI	13.6	14.38
VA	14.7	13.76
ME	15.2	15.06
NH	14.6	13.64
NV	17.6	15.72
NE	13.1	13.06
NC	13.8	14.80
MA	13.2	13.12

*Table 1. True values of State Frequent Mental Distress Scores vs. our Model's Predicted Values*

## Discussion

### Individual Categories

We found that the health and social metrics which had the highest correlations with Frequent mental distress were:

Social and Economic Factors, Firearm Deaths, percentage of people in poverty, Food insecurity, Social support, Adverse Childhood experiences, Access to clinical care and dental visits, flu vaccinations, at-risk behaviors, teen births, Smoking and Tobacco Usage, e-cigarette usage, general health outcomes, drug usage, physical health and distress, chronic conditions, arthritis, cardiovascular, kidney, and pulmonary diseases, depression, diabetes, and high blood pressure.

While this list may seem fairly obvious, it is necessary to note that the features were included as long they had a correlation with Frequent Mental distress of over 0.5. And while they may have a correlation with mental distress, there is no conclusion for causation.

### Machine Learning Model

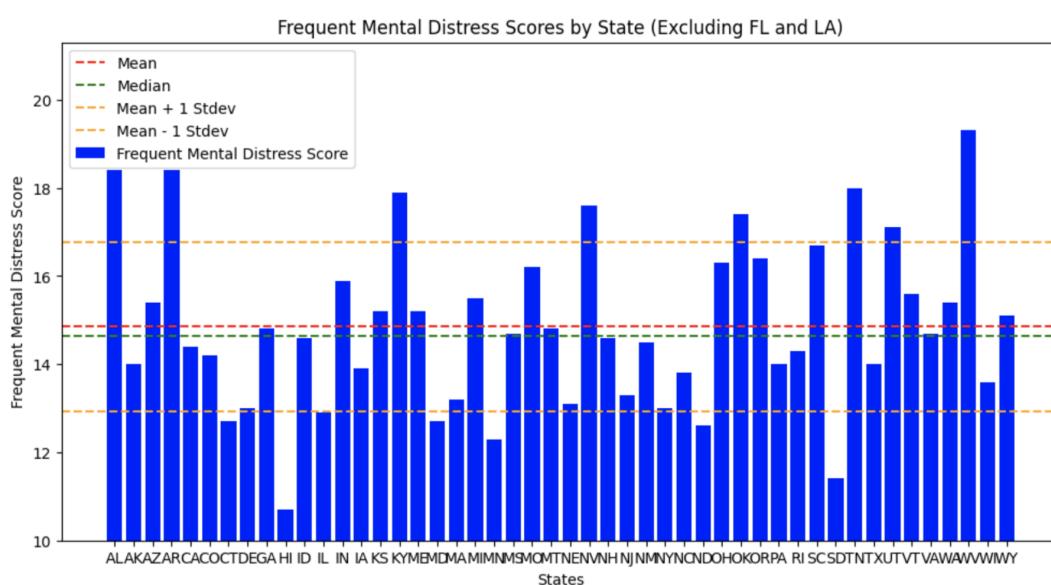
Our model had an R-squared value of 0.71, and a Mean Squared Error of 0.8445. The mean squared is fairly low, indicating a good fit of the model. The  $R^2$  value being 0.71 indicates a fair fit, however considering the number of variables factored in, the  $R^2$  score may be acceptable.

We are also able to see the scores model predicted for known values of mental distress, which are fairly similar to the real values.

### Testing the model

We input arbitrary feature values and received an appropriate predicted Frequent mental distress score, which answers our original question: how can we identify other communities with low positive mental-health outcomes?

By comparing the model's predicted output to the measures of central tendency and standard deviations of the states' Frequent mental distress scores, we can gain an understanding of whether the tested population has relatively better or worse mental health outcomes than the American average and by how much.



*Fig. 3. Chart displaying Frequent Mental Distress scores of all states excluding Florida and Louisiana (no reported Frequent Mental Distress scores). The mean, median and standard deviations of the states Frequent Mental Distress Scores are shown with red, green and yellow dashed lines, respectively.*

It is important to keep in mind however, that this data is based on American populations. It may be true that people in other countries and cultures may experience mental health issues due to completely different factors (ongoing wars, better weather, cultural behavior, etc). This predictive model may only be fit for those communities who experience mental distress in a similar way to Americans.

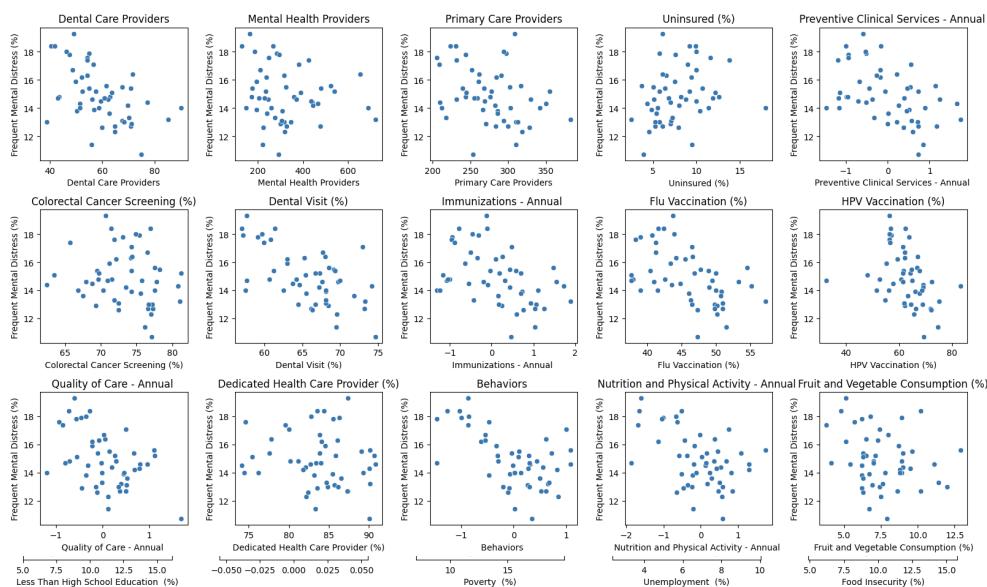
## Takeaway

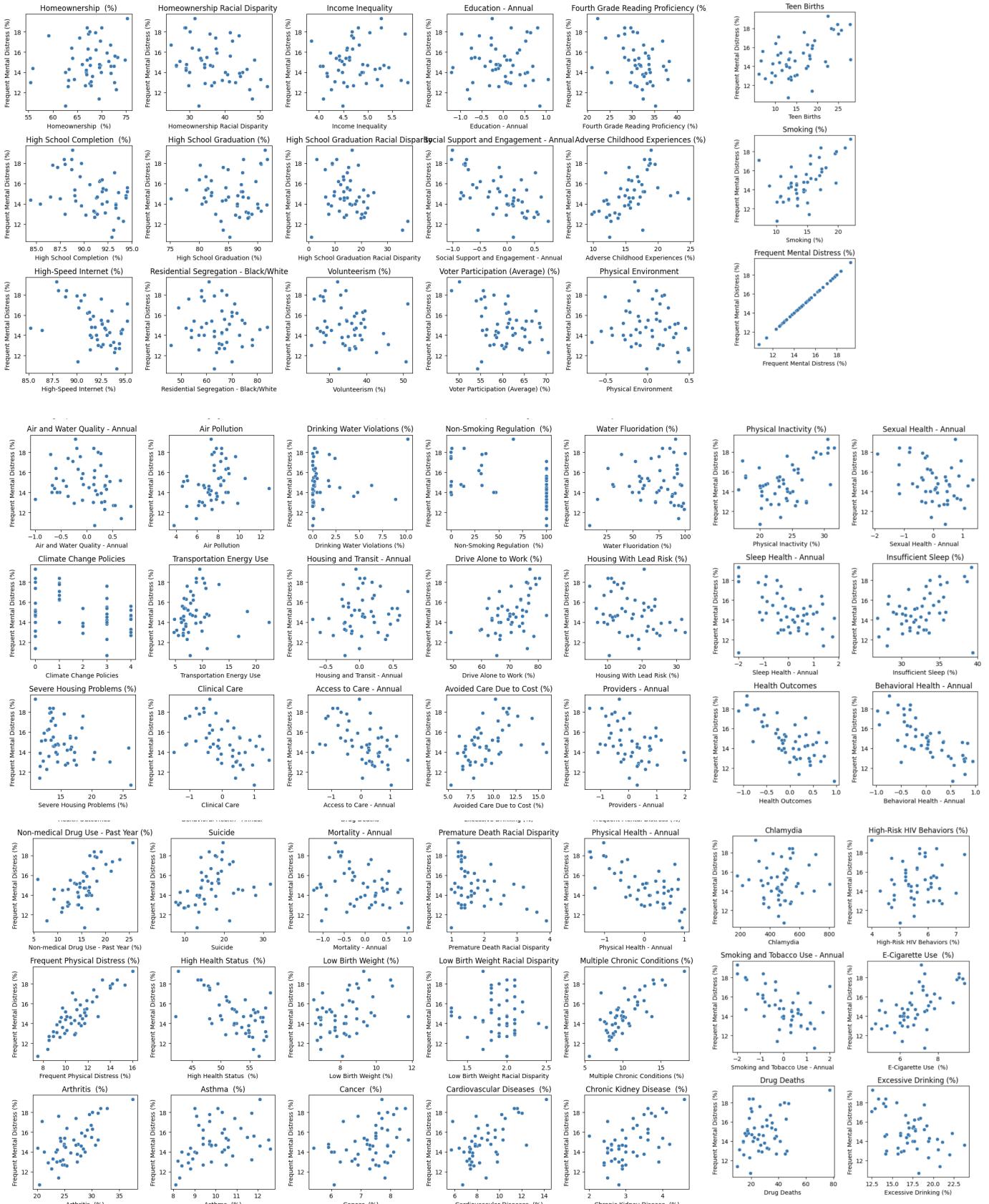
Our team remains optimistic about the potential insights that could be gained from our analysis. The visualizations, including the correlation heatmap and the scatter plots, suggest that while some relationships between mental health distress and various socioeconomic and health-related factors are evident, the complexity of mental health issues demands an approach that considers many elements. Since we have many factors that are in play, more advanced techniques, such as Random Forests, could provide a deeper understanding of the most influential factors so that we may further facilitate the identification of commonalities among affected groups and foster the development of targeted public health strategies.

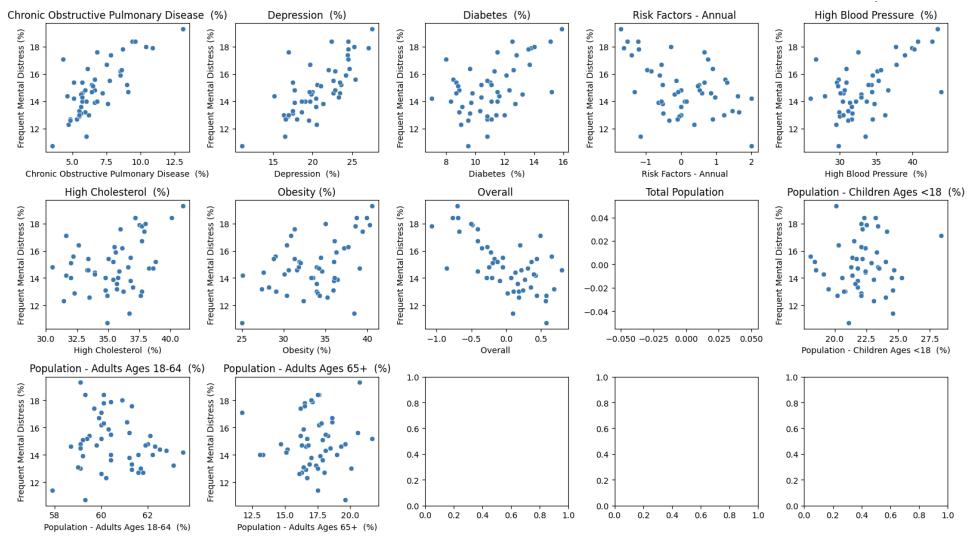
## References

U.S. Census Bureau. (n.d.). 2022 Annual Report. America's Health Rankings.  
<https://www.americashealthrankings.org/>

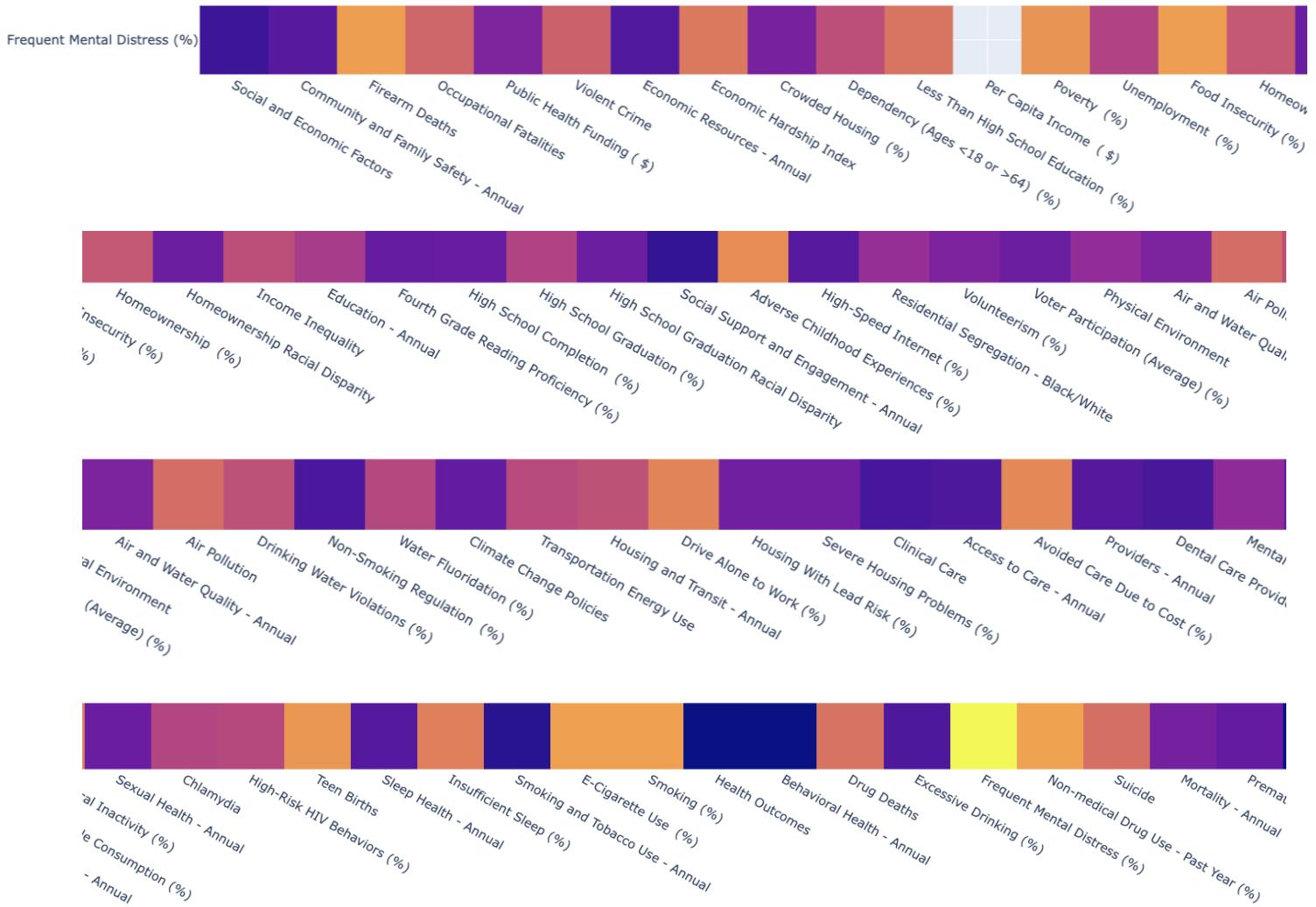
## Appendix A.







## Appendix B.



## APPENDIX B. (cont'd)

