
SPATIAL MODELING OF CARDIOVASCULAR DISEASE INCIDENCE POSITIVELY ASSOCIATED WITH PM_{2.5}

A PREPRINT

Johan Booc

Department of Statistics
Texas A&M University

jbooc24@tamu.edu

Christina Kim

Department of Statistics
Texas A&M University

christinaykim3@tamu.edu

Shombit Roy

Department of Statistics
Texas A&M University

shombit123@tamu.edu

April 23, 2024

ABSTRACT

1 Cardiovascular disease (CVD) is the leading cause of death in the United States. By using a spa-
2 tial modeling technique (geographically weighted regression) and even after adjusting for median
3 household income and unemployment rates (finding similar correlations studied in the past), we
4 observed that the concentration of PM 2.5 is centered in the Stroke Belt region. Policymakers and

health practitioners can use these results to identify targeted interventions to curb the increasing rates of CVD, aiming to halt one of the world's deadliest diseases.

Keywords Fine Particulate Matter (PM2_5) • Cardiovascular Disease (CVD) • Cardiovascular Mortality (CVM)

1 Introduction

As cardiovascular disease is the leading cause of death in the US, numerous past studies have been done on this disease, specifically on the older population the 65+ year age group. Our study aims to further investigate the impact, focusing on the 18-44-year-old age population, since generally, there has been less interest in the effects of CVD on them.

Our approach involves analyzing this relationship to produce risk prediction estimates for our specified age group. There are several factors involved in influencing the incidence rates of CVD from genetics, lifestyle, diet, and smoking/alcohol habits. We look into the covariates - PM 2.5 concentration, median household income, and unemployment rates - to analyze the spatial variation for our specified age population, thus decreasing future CVD incidence rates. Specifically, we focus on PM 2.5 concentration as the world continues to depend on fossil fuels/natural gas and climate change is an ongoing concern. Earlier studies such as (Warsito et al. 2018) and (Liu et al. 2020) highlight the modern-day threats from air pollution using a Robust GWR model and Bayesian-temporal model, respectively. From studying our covariates, we can examine the effects of certain CVD risk factors for the year 2015.

The geographically weighted regression model uses local variables and weights to produce statistical visualizations of the US's regional variation between our response variable (CVD deaths) and its covariates. Other methods, such as the traditional linear regression model and clustering, have been used in past CVD studies. However, we found the GWR model allows us to see each covariate's local significance and magnitude. The spatial distribution aspect of the GWR model highlights regions that are concentrated with CVD rates, whether incidence percentage rates are increasing or decreasing, and any significant correlations. This makes the GWR model efficient in analyzing the dynamic relationship of CVD rates and the factors that contribute to it.

Using a geographically weighted regression approach, the overall outcome of our study aids policymakers and health practitioners in implementing the necessary interventions for targeted regions that are being influenced the most by our covariates. For public health experts, it's ideal to tailor aid to each specific region of the US first to eventually reach a national decrease in CVD incidence rates.

2 Related Works

There's a growing popularity of using spatial models in the epidemiological domain to analyze the distribution and factors of a disease. The techniques often used differ between studies but all have the same goal: to reduce the disease incidence and mortality rates. Statistical regression models are popular methods for CVD studies. The underlying causes for CVD are dynamic and can be seen as an interconnected web. (Zelko et al. 2023) examines the relationship

between CVD and covariates similar to our study (air pollution, social determinants, and county-level data). From using a GWR model, their results found that counties in the South had the highest exposure to PM 2.5 concentrations whereas counties in the Northeast had the lowest. A strong correlation was also found between household income, race, and healthcare access. Overall, there are several causes to consider with CVD, and where the covariates are concentrated is another important consideration.

Compared to traditional regression models, the geographically weighted model (GWR) explores spatial data by considering varied coefficients for a certain spatial unit (Gebreab and Diez Roux 2012). The GWR model contains the parameter, neighborhood (also known as bandwidth) and builds on the weighted least squares method to estimate the regression coefficients. We want to see values closer to the point of interest since they carry more weight, thus having a greater influence.

The Ordinary Least Squares (OLS) model is a popular choice of method for public health studies because it generates the regression coefficients on a global scale and captures the average differences between covariates. However, this means spatial variability between units can be easily hidden for the OLS model. Our paper acknowledges that PM 2.5 concentration and socioeconomic factors are not spatially constant with cardiovascular disease. The GWR model is more suitable for our study because the goal of public health is to improve the population as a whole. To get to the root cause of CVD mortality rates and see future improvement, the OLS model should be treated as the ‘null’ model, with the GWR model used to test and verify that it is a statistically significant better fit.

The GWR model spatially displays a relationship between CVD deaths and their covariates and analyzes disparities at the local scale. Errors are also minimized between the actual model and any estimates. This makes the GWR model suitable for seeing the socioeconomic factors that affect different regions and narrows down our focus to the areas in need of improvement, which helps health practitioners implement policies for that region. Past studies (Zelko et al. 2023),(Terry et al. 2023), and(Singh et al. 2019) tend to focus on the trends of socioeconomic covariates and their spatial patterning. However, our study extends past studies by including PM 2.5 concentrations as one of our covariates. From this, we can fully understand the dynamic relationship between counties and CVD incidence/mortality rates.

Risk assessment and risk estimates uncover the key factors associated with CVD. Because the concentration of specific races varies by county, we studied the dose-response relationship of PM 2.5 concentrations and the socioeconomic covariates of CVD mortality at the county level. This helps researchers and health practitioners to develop the necessary risk-preventative measures and allocate resources to the areas that need them the most. By putting the focus on the county level (rather than on individual states/nationally), resources can be allocated accordingly to the regions that need the most assistance. This can lead to a reduction in CVD mortality rates, both nationally and globally.

3 Methods

3.1 Data Collection

A data set of Medicare services and claims from the Centers for Disease Control and Prevention (CDC) website was loaded to analyze Medicare claims data, specifically Cardiovascular death rates across different counties and States in the US. Racial and geographic data were retrieved from the Census and TIGER Bureau. The median income was extracted from the CENSUS API. The air quality index was extracted from the NASA PM 2.5 Concentration data. Finally, the unemployment rate was extracted from the CDC website. These data sources were collected and prepared for analysis to understand the relationship between these factors and Cardiovascular death rates. The code used to extract and filter the data is available in our GitHub repository.

3.2 Data Preprocessing

Several steps were taken to ensure the reliability and accuracy of the results when preparing the data for statistical analysis. The data from various sources was integrated into one file by year/county/latitude and longitude/geometries. This integration was achieved through coding in RStudio Version 4.3.2, specifically focusing on data from 2015, which allowed for a consistent time frame across all data sets. During the cleaning and processing phase, features with empty geometries were removed from the shapefile to ensure the removal of NA values in the dataset. Missing data entries were also omitted, which led to Nantucket County being removed as it had incomplete data and was removed from our analysis. Also, in this dataset, we are not analyzing Alaska and Hawaii due to their geographical isolation from the contiguous United States. Centroids of the multi-polygon geometries were calculated to provide a single point representing the location of these complex shapes. This step was essential for conducting precise location-based assessments. Lastly, the data frame was transformed into a SpatialPolygonsDataFrame, making it suitable for use in the Geographic Weighted Regression (GWR) model.

3.3 Statistical Analysis

A Geographically Weighted Regression (GWR) model was used in the statistical analysis, which was created using the GWmodel library. These packages were integral in setting up the model framework and executing the analysis, providing tools for spatial data manipulation, regression modeling, and data visualization. The optimal bandwidth for the GWR model was estimated using cross-validation with hyperparameters such as Gaussian kernel and fixed bandwidth. The Gaussian kernel was used because it assigns weights to the data points based on their distance from where the model is being calculated, with closer points receiving higher weights. This weighting falls off smoothly and symmetrically. The adaptiveness is FALSE because the goal is systematically comparing coefficients across different geographic regions. So, a fixed bandwidth can help ensure that each region is analyzed under the same spatial constraints. This function modeled the response variable, Cardiovascular death rate per 100,000 residents:

$$y = \beta_0 * \%White + \beta_1 * \%Black + \beta_2 * \%Hispanic + \beta_3 * \%Asian + \beta_4 * PM2.5 + \beta_5 * MedianIncome + \beta_6 * \%Unemployed$$

The variable “y” represents the CVM per 100,000 residents in each county. The GWR was performed using the identified optimal bandwidth, incorporating the same predictors as those used for bandwidth selection. This ensures the model’s findings are reliable and applicable to the predictors. Lastly, the p-values were found to determine significance by using the t-values obtained from the local regression coefficients from the GWR model. Then, the cumulative distribution function of the T-distribution is used to calculate the p-values. It was done as a two-tail test as we tried to find significant differences, and degrees of freedom were found from the GWR diagnostic output.

3.4 Mapping

We generated geographic plots to visualize the spatial distribution of the dependent variable (deaths from CVD per 100,000 people) across different counties in the US, highlighting the significance of the variables, and plotted the significance of each variable within different regions to assess the overall effect on the CVD death rate.

This approach allows for examining how various socioeconomic, environmental, and demographic factors influence health outcomes across different regions in the United States. We showed this in simulated data where we used a piecewise function to the coordinates, which divides the geographical space into four quadrants and assigns different coefficients to them to model for spatial heterogeneity. Also, a spatial autocorrelation test (Moran’s I) is performed on the residuals of the GWR model. From the table in our appendix and the graph in our appendix, we have a high standard deviation and low p-value, strongly suggesting that the residuals of the GWR model are not randomly distributed but show significant spatial autocorrelation. The graph shows that the majority of the residuals are close to zero, which means the predicted values are pretty close to the observed values, indicating that the GWR model has accounted for spatial variation effectively and that it is a good fit.

This methodological outline ensures that each step of the data handling, analysis, and visualization process is documented, providing transparency and reproducibility of the research findings.

4 Results

4.1 Regression Summary

The GWR regression showed two models: global regression and GWR regression output. They show the relationships between socioeconomic, demographic, and environmental variables and death rates. The first is a global regression model that does not consider spatial correlation despite revealing that all predictors are statistically significant. Hence, while the model does suggest that our variables are indeed important, the global model may overlook local variations that are crucial in understanding the true nature of the data. This can be seen in Figure 2, which shows the residual

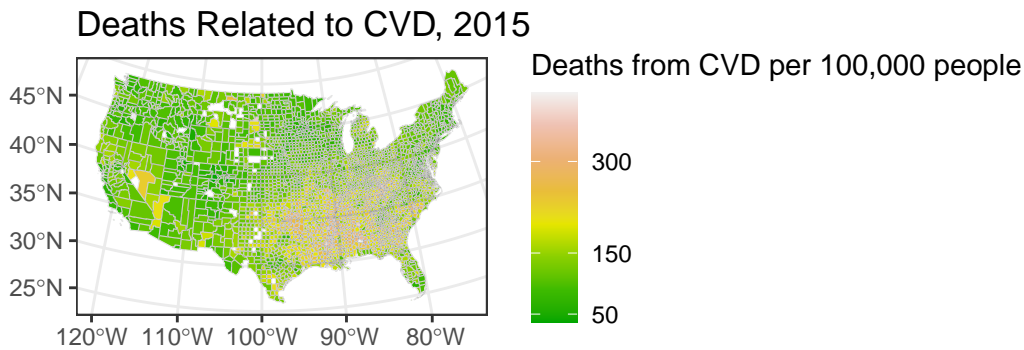


Figure 1

of the global model, and there appears to be a spatial pattern to the residuals, with specific areas showing clusters of higher residuals and the other regions showing clusters of lower residuals. This clustering of residuals suggests that the global regression model may not be capturing all the spatial variation in the data. This implies that the relationship between the independent and dependent variables might differ across different locations.

On the other hand, the geographically weighted regression (GWR) model incorporates spatial variation, which is a critical factor given the data context. We can see that the AIC and BIC values are lower for the GWR regression than the Global regression, which makes it the most preferred model given the data. In summary, by accommodating the spatial component present in the data, the GWR model provides a more realistic interpretation of how various factors influence death rates across different regions.

4.2 Local Significance Plots

In Figure 1, the plot shows an exploratory data analysis (EDA) plot, specifically a choropleth map displaying the number of deaths from cardiovascular disease (CVD) per 100,000 people across the contiguous United States for the year 2015. The particulate concentrations include all the covariates mentioned in the methods section. The regions along the higher latitudes (nearing 45 N) and towards the eastern section (approaching 80 W) display darker shades, suggesting higher CVD death rates in these areas, specifically the midwest and southwest regions.

Figure-3 represents the local significance and magnitude of the ‘% White’ demographic parameter on cardiovascular disease outcomes. We plotted for the local significance and magnitude because it can help identify areas where the predictor variable has a stronger or weaker influence on the outcome, leading to targeted insights that would not be

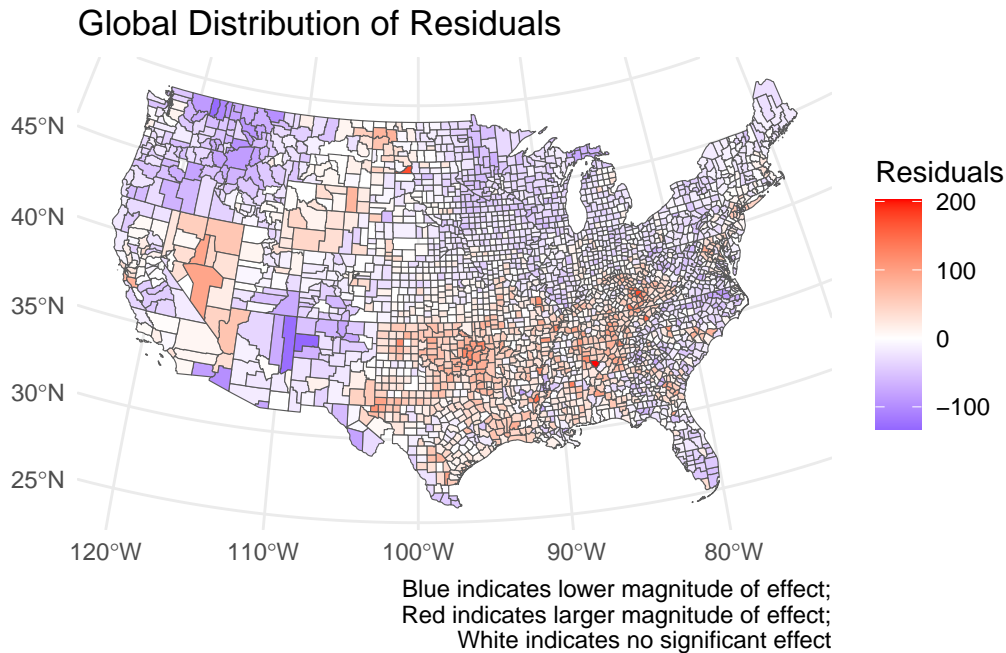


Figure 2

possible with a global model. From the plot, the prominent dark purple areas in the central part of the country, extending towards the southeastern regions, indicate a significant negative correlation between the percentage of the white population and CVD outcomes in these areas. Like the Northwest region, however, there are few areas, particularly in the southwest near New Mexico and Arizona, where there are positively correlated to CVD death rate, with the white regions of counties suggesting no significance and this is white region is similar for all the figures.

Fig-4 represents the local significance and magnitude of the ‘% African’ demographic parameter on cardiovascular disease outcomes. From the plot, areas located approximately in the northern central region indicate a significant positive correlation to CVD rate, and areas, notably in the central to southeastern areas of the map, suggest a significant negative correlation.

Fig-5 represents the local significance and magnitude of the ‘% Hispanic’ demographic parameter on cardiovascular disease outcomes. From the plot, areas across the central and southeastern regions are shaded in purple, suggesting a significant negative correlation between the percentage of the Hispanic population and CVD outcomes. The isolated red patches in the north-central region indicate areas where an increased Hispanic population correlates with higher CVD outcomes.

Fig-6 represents the local significance and magnitude of the ‘% Asian’ demographic parameter on cardiovascular disease outcomes. From the plot, a substantial portion of the map, particularly across the central to eastern regions, is colored in various shades of purple. This suggests that in these areas, an increased percentage of the Asian population correlates with lower CVD outcomes. However, there are some parts of the West where higher percentages of the Asian population are associated with an increase in CVD outcomes.

Figure 7 represents the local significance and magnitude of the “%p2.5 air quality” parameter on cardiovascular disease outcomes. From the plot, in the southern and northeastern regions, there increased levels of PM2.5 colored as red in the plot, meaning they are associated with higher rates of CVD. There are some regions in the West where it was negatively correlated with CVD.

?@fig-8 represents the local significance and magnitude of the “median income” parameter on cardiovascular disease outcomes. From the plot, most of the region suggests a significant negative correlation between the median income parameter and CVD outcomes. However, some parts of Texas indicate a positive correlation to CVD outcomes.

?@fig-9 represents the local significance and magnitude of the “Unemployment” parameter on cardiovascular disease outcomes. From the plot, there is a positive correlation to CVD rates in most of the central region, and in the southwest region, there is a negative correlation.

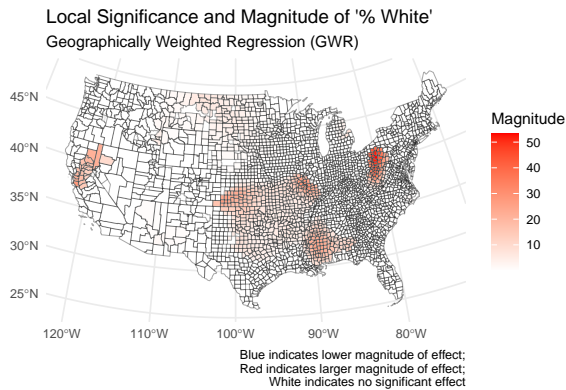


Figure 3

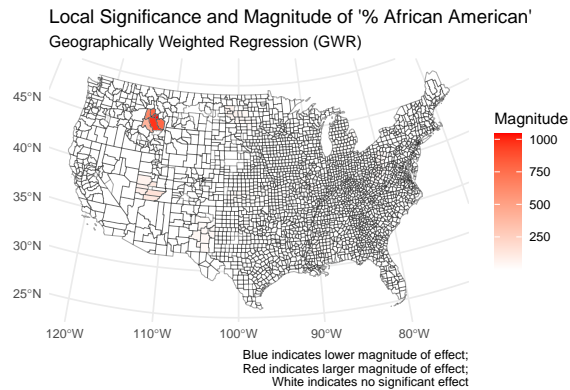


Figure 4

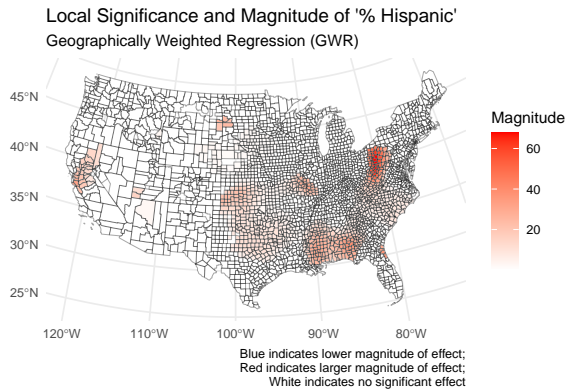


Figure 5

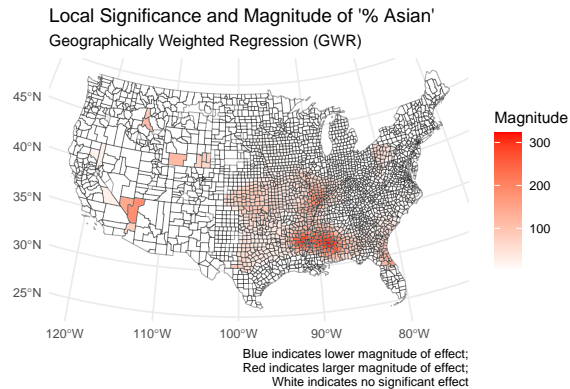


Figure 6

- Demographic Impact: The percentage of white and Hispanic populations shows significant regional variations in association with death rates. Higher proportions of white populations correlate with lower death rates in central areas, whereas higher percentages of Hispanic populations are linked to lower death rates in the West

and Southwest. Conversely, higher percentages of African American populations are associated with higher death rates in certain Midwestern and Southeastern regions.

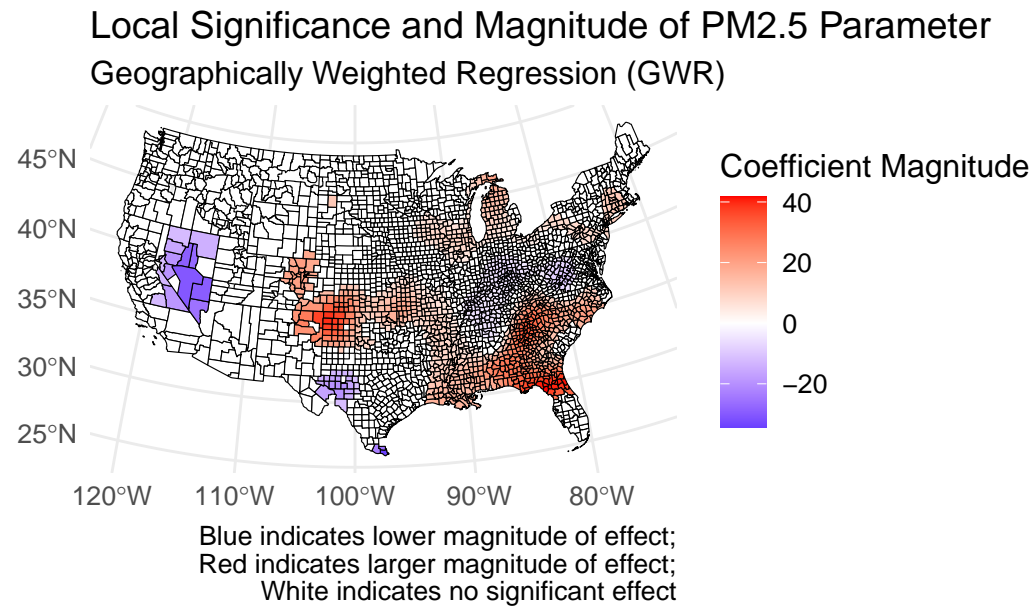


Figure 7

- Environmental Influence: Air quality, indicated by PM2.5 levels, demonstrates a significant positive relationship with death rates, particularly east of the Rockies, highlighting environmental health as a major concern.
- Socio-economic Correlation: Median income levels across many regions show a consistent negative association with death rates, suggesting that higher income areas generally experience fewer deaths.

5 Discussion

The purpose of this study was to fill a gap in prior studies where the impacts of CVD were primarily studied within the Stroke Belt, rather than the country at large.

Our goal was to put our findings towards answering the following question: what are the socioeconomic and environmental factors affecting CVD rates in the United States?

These maps reveal that the relationships between race, socio-economic factors, environmental quality, and death rates are complex and highly localized. The significance and strength of these relationships vary considerably across different parts of the United States. In contrast, some of the socio-economic factors such as median income show widespread, consistent significance, implying that the significance of the relationship with CVD outcomes is nearly constant across various locations. This differs from the heterogeneous local significance that we observed in the majority of our other variables.

Local Significance of Median Income Parameter
Geographically Weighted Regression (GWR)

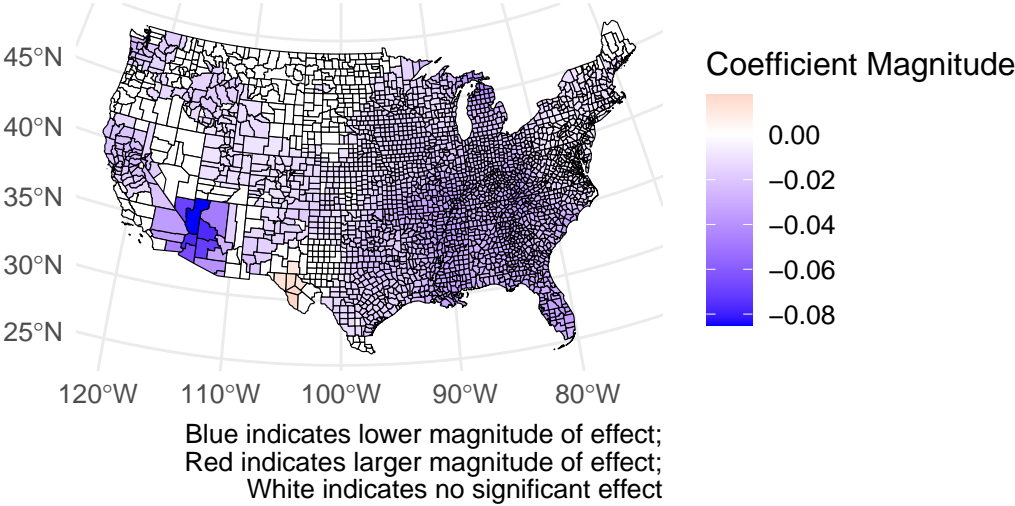


Figure 8

Local Significance of Unemployment Parameter
Geographically Weighted Regression (GWR)

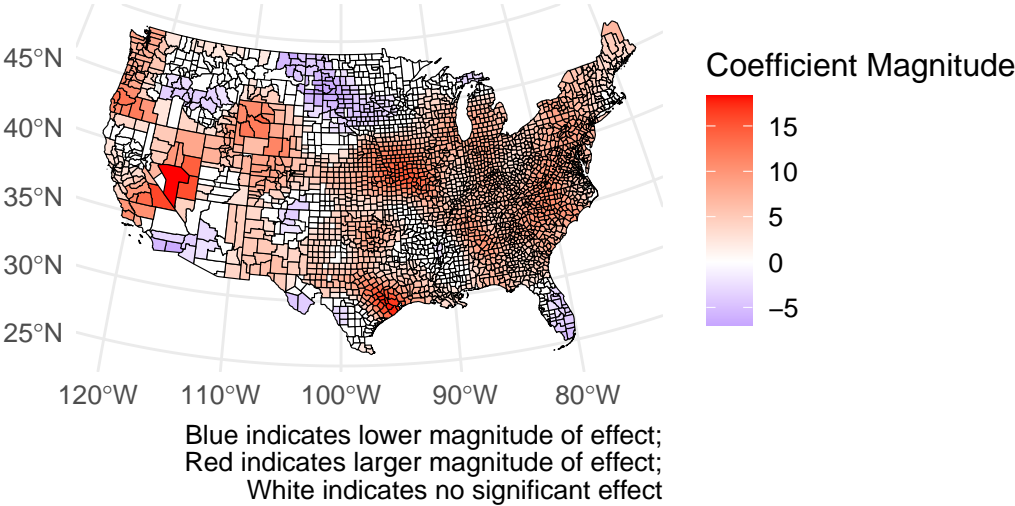


Figure 9

Table 1
 % latex table generated in R 4.3.2 by xtable 1.8-4 package % Tue Apr 23 00:11:25 2024

	1	2	3	4	5	6	7
	183.22	196.12	7.24	2.56	1.15	2.44	2.62
	81.32	93.80	4.14	2.53	1.77	3.39	2.65
[ht]	302.16	313.78	16.89	4.06	1.27	2.42	2.09
	140.83	151.46	5.25	2.13	1.10	2.35	2.73
	131.77	128.39	8.33	2.67	1.07	2.34	2.32
	282.66	294.23	14.14	3.77	1.24	2.44	2.22

Another factor to consider from the plots is the percentage of each race that inhabits each county. We can observe that for African Americans, Hispanics, and Asians that the impacts are significant in localized areas of the country, which may reflect underlying health disparities in access to medical care. One limitation of race percentages is the misreporting of medical records affecting minorities (Tabb et al.) However, this oversight lends further credence to the fact that intervention is needed in order to combat racial health disparities. We can look at the Variance Inflation Factor (VIF) to determine if multi-collinearity has an impact on our results:

Typically, when examining the VIF of a GWR model, we want the coefficients for each of the seven columns—which represent our seven independent variables—to be less than 15. We can see in Table 1 that this standard holds for most cases of the columns outside the % White column, which has extremely high values. This represents a limitation in our model, as it shows multicollinearity with respect to the white variable, however, this is to be expected given that the white population makes up a majority of the United States.

The concentration of PM2.5 and the consequential reduction in air quality has a significant impact localized in the central and southeastern regions of the United States, suggesting environmental health concerns that might require region-specific intervention.

We have shown that by using a GWR model to analyze the relationship between CVM and socioeconomic covariates, the factors that have the most significant impact on death rates vary by area of the country. This highlights the need to identify efficient methods of intervention in order to curb one of the deadliest groups of diseases on the planet (Marlow 1994).

6 Appendix

The R code used for this study can be found in our public GitHub repository: <https://github.com/jbooc117/STAT489-Project.git>.

References

- Gebreab, Samson Y., and Ana V. Diez Roux. 2012. "Exploring Racial Disparities in CHD Mortality Between Blacks and Whites Across the United States: A Geographically Weighted Regression Approach." *Health & Place* 18 (5): 1006–14. <https://doi.org/10.1016/j.healthplace.2012.06.006>.
- Liu, Yi, Jingjie Sun, Yannong Gou, Xiubin Sun, Dandan Zhang, and Fuzhong Xue. 2020. "Analysis of Short-Term Effects of Air Pollution on Cardiovascular Disease Using Bayesian Spatio-Temporal Models." *International Journal of Environmental Research and Public Health* 17 (3): 879. <https://doi.org/10.3390/ijerph17030879>.
- Marlow, Hilary F. 1994. "The Pharmaceutical Industry Viewpoint." *Cardiology* 85 (1): 102–12. <https://doi.org/10.1159/000176769>.
- Singh, Gitanjali M., Ninon Becquart, Melissa Cruz, Andrea Acevedo, Dariush Mozaffarian, and Elena N. Naumova. 2019. "Spatiotemporal and Demographic Trends and Disparities in Cardiovascular Disease Among Older Adults in the United States Based on 181 Million Hospitalization Records." *Journal of the American Heart Association* 8 (21): e012727. <https://doi.org/10.1161/JAHA.119.012727>.
- Terry, Katrina, Mohamed Makhoul, Salah E. Altarabsheh, Vaishali Deo, Fanny Petermann-Rocha, Yakov Elgudin, Khurram Nasir, Sanjay Rajagopalan, Sadeer Al-Kindi, and Salil Deo. 2023. "Trends in Cardiovascular Disease Mortality by County-Level Social Vulnerability Index in the United States." *Journal of the American Heart Association* 12 (20): e030290. <https://doi.org/10.1161/JAHA.123.030290>.
- Warsito, Budi, Hasbi Yasin, Dwi Ispriyanti, and Abdul Hoyyi. 2018. "Robust Geographically Weighted Regression of Modeling the Air Polluter Standard Index (APSI)." *Journal of Physics: Conference Series* 1025 (May): 012096. <https://doi.org/10.1088/1742-6596/1025/1/012096>.
- Zelko, Andrea, Pedro R. V. O. Salerno, Sadeer Al-Kindi, Fredrick Ho, Fanny Petermann Rocha, Khurram Nasir, Sanjay Rajagopalan, Salil Deo, and Naveed Sattar. 2023. "Geographically Weighted Modeling to Explore Social and Environmental Factors Affecting County-Level Cardiovascular Mortality in People With Diabetes in the United States: A Cross-Sectional Analysis." *The American Journal of Cardiology* 209 (December): 193–98. <https://doi.org/10.1016/j.amjcard.2023.09.084>.