
SPATIAL MODELING OF CARDIOVASCULAR DISEASE

INCIDENCE POSITIVELY ASSOCIATED WITH PM2.5

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

May 5, 2024

ABSTRACT

Cardiovascular disease (CVD) is the leading cause of death in the United States. After adjusting for median household income and unemployment rates (finding similar correlations studied in the past), we used the spatial modeling technique, geographically weighted regression (GWR), and observed that the coefficient for the concentration of PM 2.5 is centered in the Stroke Belt region.

5 Policymakers and health practitioners can use these results to identify targeted interventions to curb
 6 the increasing rates of CVD and help halt one of the world's deadliest diseases.

7 **Keywords** Fine Particulate Matter (PM2.5) • Cardiovascular Disease (CVD) • Cardiovascular Mortality (CVM)

8 1 Introduction

9 As cardiovascular disease is the leading cause of death in the US, numerous past studies have been done, specifically
 10 on the 65+ year age group. Our study aims to investigate the impact of CVD mortality rates further, focusing on the
 11 18-44-year-old age population. There has been less interest in the effects of CVD outcomes on them because they are
 12 not seen as susceptible to the disease compared to the older population.

13 Our approach analyzes the relationship between CVD rates and our covariates (socioeconomic factors and PM 2.5
 14 concentration) to produce risk prediction estimates for our specified age group. There are several factors involved in
 15 influencing the mortality rates of CVD however, we look at PM 2.5 concentration, median household income, and
 16 unemployment rates to analyze the spatial variation of CVD rates for our specified age group. Earlier studies such as
 17 (Warsito et al. 2018; Liu et al. 2020) highlight the modern-day air pollution threats by using a Robust GWR model
 18 and a Bayesian spatiotemporal model, respectively. By studying our covariates, we can examine their effects on CVD
 19 rates to produce risk estimates for the year 2015.

20 To produce statistical visualizations of the US's regional variation between our response variable (CVD deaths) and
 21 its covariates, we used the geographically weighted regression model through local variables and weights. Other
 22 methods, such as the traditional linear regression model and clustering have been used in past CVD studies. However,
 23 we found the GWR model allows us to see the local significance (compared to the Ordinary Least Squares model)
 24 and nonstationary effects (compared to most clustering methods) from each of our covariates. The spatial distribution
 25 aspect of the GWR model highlights CVD mortality rate percentage changes and correlations between CVD outcomes
 26 and our covariates. Overall, we found the GWR model useful for analyzing the dynamic relationship of CVD rates
 27 and the factors that contribute to it.

28 Using a geographically weighted regression approach, we found distinct geographic patterns for demographics and
 29 CVD outcomes. PM 2.5 levels and CVD rates showed a positive correlation in the south and northeast regions while
 30 western regions exhibited a negative correlation. For our socioeconomic covariates, a higher median income was asso-
 31 ciated with lower CVD levels and the correlation between unemployment rates and CVD outcomes varied. From this,
 32 the overall outcome of our study aims to assist policymakers and health practitioners implement necessary intervention
 33 for targeted regions.

34 2 Related Works

35 There's been a growing popularity of using spatial models in the epidemiological domain to analyze the distribution
 36 and factors of a disease. The techniques used often differ between studies but all have the same goal: to reduce the
 37 disease incidence and mortality rates. Statistical regression models are popular methods for CVD studies and the
 38 underlying causes are dynamic. (Zelko et al. 2023) examined the relationship between CVD and covariates similar
 39 to our study (air pollution, social determinants, and county-level data). From using a GWR model, their coefficients
 40 showed counties in the South had the highest exposure to PM 2.5 concentration whereas counties in the Northeast
 41 had the lowest. In addition, they found a strong correlation between household income, race, and healthcare access.
 42 However, their study focused on Type 2 diabetes mortality rates in addition to CVD outcomes. Our study emphasizes
 43 CVD rates for our specified age group with a focus on PM 2.5 levels. The study found that their coefficients had
 44 statistically significant spatial variation so it is important to see where the covariates are concentrated. Overall, there
 45 are several causes to consider with CVD.

46 Compared to traditional regression models, the geographically weighted model (GWR) is a local model with spatially
 47 varying coefficients, as opposed to the global model from Ordinary Least Squares (OLS) (Gebreab and Diez Roux
 48 2012). The GWR model includes the kernel density function, weights, and the parameters distance and the number of
 49 neighborhoods (also known as bandwidth). It builds on the weighted least squares method to estimate the regression
 50 coefficient, with the diagonals of the weighted matrix representing the location of each observation. We want to see
 51 values closer to the point of interest since they carry more weight, thus having a greater influence. We acknowledge
 52 that PM 2.5 concentration and socioeconomic factors are not spatially constant with cardiovascular disease. While a
 53 global model helps improve the population as a whole, we found the GWR model allows us to see the local significance
 54 of our covariates. This is an important feature to see future improvement in CVD mortality rates. We consider the
 55 OLS model as the ‘null’ model, with the GWR model used to test and verify that it is a statistically significant better
 56 fit.

57 The GWR model spatially displays a relationship between CVD deaths and their covariates to analyze disparities at the
 58 local scale and minimize errors. This makes the GWR model suitable for seeing the socioeconomic factors that affect
 59 different regions and narrows down our focus to the areas in need of improvement, which helps health practitioners
 60 implement policies for that region. Past studies (Zelko et al. 2023; Terry et al. 2023; Singh et al. 2019) have focused
 61 on the trends of socioeconomic covariates and their spatial patterning. However, our study extends the efforts of past
 62 studies by including PM 2.5 concentrations as one of our covariates. From this, we can fully understand the dynamic
 63 relationship between air quality and CVD mortality rates.

64 Risk assessment and risk estimates uncover the key factors associated with CVD. We studied the concentration of
 65 specific races, the relationship of PM 2.5 concentrations, and the socioeconomic covariates of CVD rates at the county
 66 level. This helps researchers and health practitioners to develop the necessary risk-preventative measures and allocate
 67 resources to the areas that need them the most. By putting the focus on the county level (rather than individual

68 states/nationally), resources can be allocated accordingly to the regions that need the most assistance, leading to a
 69 reduction in CVD mortality outcomes.

70 **3 Methods**

71 **3.1 Data Collection**

72 A data set of Medicare services and claims from the Centers for Disease Control and Prevention (CDC) website was
 73 loaded to analyze Medicare claims data, specifically Cardiovascular death rates across different counties in the US.
 74 Racial and geographic data were retrieved from the Census and TIGER Bureau. The median income was extracted
 75 from the CENSUS API. The air quality index was extracted from the NASA PM 2.5 Concentration dataset. Finally,
 76 the unemployment rate was extracted from the CDC website. These data sources were collected and prepared for
 77 analysis to understand the relationship between these factors and Cardiovascular death rates. The code used to extract
 78 and filter the data is available in our GitHub repository: <https://github.com/jbooc117/STAT489-Project.git>.

79 **3.2 Data Preprocessing**

80 We took several steps to ensure the reliability and accuracy of the results when preparing the data for statistical
 81 analysis. We integrated the data from various sources into one file and grouped it by year, county, and geometries.
 82 This integration was achieved through coding in RStudio Version 4.3.2, specifically focusing on data from 2015,
 83 which allowed for a consistent time frame across all data sets. During the cleaning and processing phase, we removed
 84 features with empty geometries from the shapefile to ensure the removal of missing values in the dataset, which led to
 85 Nantucket County being removed as it had incomplete data and was removed from our analysis. Also, we removed
 86 Alaska and Hawaii in this dataset as they are geographically separate from the US. Centroids of the multi-polygon
 87 geometries were calculated to provide a single point representing the location of these complex shapes.

88 **3.3 Statistical Analysis**

89 A geographically weighted regression (GWR) model was used in the statistical analysis, and it was run using the
 90 GWmodel library (Gollini et al. 2015). The functions within the GWmodel library were integral in setting up the model
 91 framework and executing the analysis, providing tools for spatial data manipulation, regression modeling, and data
 92 visualization. The optimal bandwidth for the GWR model was estimated using cross-validation with hyperparameters
 93 we chose to use, such as Gaussian kernel and fixed bandwidth. The reason we chose the Gaussian Kernel was its
 94 smoothness. We set adaptiveness to false because the goal is to systematically compare coefficients across different
 95 geographic regions. A fixed bandwidth can help ensure that each region is analyzed under the same spatial constraints,
 96 and its cheaper. This function modeled the response variable, Cardiovascular death rate, per 100,000 county residents:

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

97 To calculate t-values from the GWR model, we introduced a correction for the t-values. This ensures accurate sig-
 98 nificance testing as the t-values from the GWR model do not follow a regular t-distribution. Then, the cumulative
 99 distribution function of the adjusted T-distribution is used to calculate the p-values. We conducted a two-tailed test to
 100 find significant differences, and effective degrees of freedom were found from the GWR diagnostic output.

101 **3.4 Mapping**

102 We generated geographic plots to visualize the spatial distribution of the dependent variable (deaths from CVD per
 103 100,000 people) across different counties in the US, highlighting the significance of the variables, and plotted the
 104 significance of each variable within different regions to assess the overall effect on the CVD death rate.

105 This approach allows for examining how various socioeconomic, environmental, and demographic factors influence
 106 health outcomes across different regions in the United States. We showed this in our simulated data where we used a
 107 piecewise function on the coordinates, which divides the geographical space into four quadrants and assigns different
 108 coefficients to them to model for spatial heterogeneity, which can be viewed in Figure 6. Also, a spatial autocorrelation
 109 test (Moran's I) was performed on the residuals of the GWR model. Based on the table and graph in our appendix,
 110 we have a high standard deviation, low p-value, and a moderate Moran's I, suggesting that the residuals of the GWR
 111 model are not randomly distributed but instead show significant spatial autocorrelation. The graph shows that the
 112 majority of the residuals are close to zero, which means the predicted values are close to the observed values. This
 113 indicates that the GWR model has accounted for spatial variation effectively and that it is a good fit.

114 This methodological outline ensures that each step of the data handling, analysis, and visualization process is docu-
 115 mented, providing transparency and reproducibility of the research findings.

116 **4 Results**

117 **4.1 Regression Summary**

118 We used two models: global regression and GWR regression output. They show the relationships between socioeco-
 119 nomic, demographic, and environmental variables and death rates. The first is a global regression model that does not
 120 consider spatial correlation despite revealing that all predictors are statistically significant. Hence, while the model
 121 does suggest that our variables are indeed important, the global model may overlook local variations that are crucial
 122 in understanding the true nature of the data. This can be seen in Figure 2, which shows the residual of the global
 123 model, and there appears to be a spatial pattern to the residuals, with specific areas showing clusters of higher resid-
 124 uals and the other regions showing clusters of lower residuals. This clustering of residuals suggests that the global

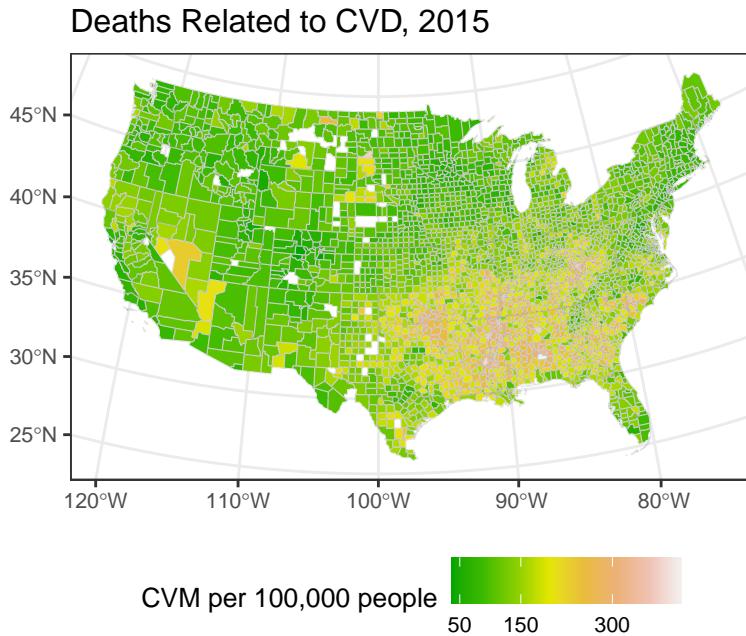


Figure 1: Higher death rates can be seen in Stroke Belt region for 2015

125 regression model may not be capturing all the spatial variation in the data. This implies that the relationship between
126 the independent and dependent variables might differ across different locations.

127 On the other hand, the geographically weighted regression (GWR) model incorporates spatial variation, which is a
128 critical factor given the data context. In summary, by accommodating the spatial component present in the data,
129 the GWR model provides a more realistic interpretation of how various factors influence death rates across different
130 regions.

131 **4.2 Local Significance Plots**

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

137 Figure 3a represents the local significance and magnitude of the ‘% White’ demographic parameter on cardiovascular
138 disease outcomes. We plotted for the local significance and magnitude because it can help identify areas where the
139 predictor variable has a stronger or weaker influence on the outcome, leading to targeted insights that would not be
140 possible with a global model. From the plot, the prominent dark red areas in the central part of the country, extend-
141 ing towards the southeastern regions, indicate a significant negative correlation between the percentage of the white
142 population and CVD outcomes in these areas. Like the Northwest region, however, there are few areas, particularly in

Global Distribution of Residuals

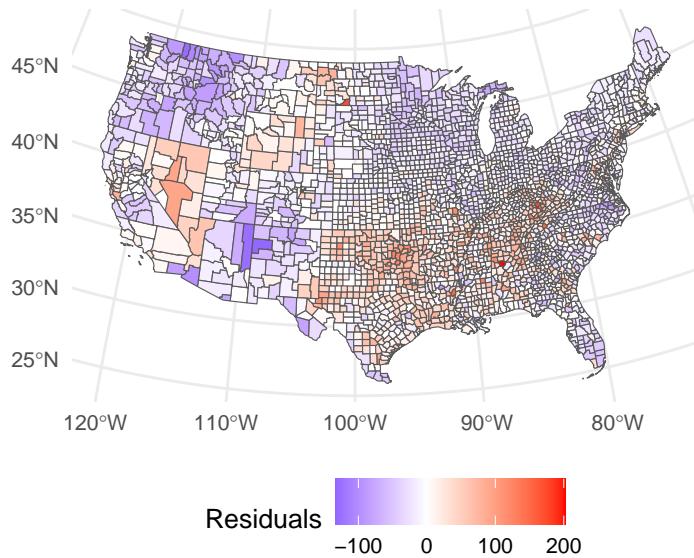


Figure 2

143 the southwest near New Mexico and Arizona, where there are positively correlated to CVD death rate, with the white
 144 regions of counties suggesting no significance and this is white region is similar for all the figures.

145 Figure 3b represents the local significance and magnitude of the ‘% African’ demographic parameter on cardiovascular
 146 disease outcomes. From the plot, areas located approximately in the northern central region indicate a significant
 147 positive correlation to CVD rate, and areas, notably in the central to southeastern areas of the map, suggest a significant
 148 negative correlation.

149 Figure 3c represents the local significance and magnitude of the ‘% Hispanic’ demographic parameter on cardiovas-
 150 cular disease outcomes. From the plot, areas across the central and southeastern regions are shaded in red, suggesting
 151 a significant negative correlation between the percentage of the Hispanic population and CVD outcomes. The isolated
 152 red patches in the north-central region indicate areas where an increased Hispanic population correlates with higher
 153 CVD outcomes.

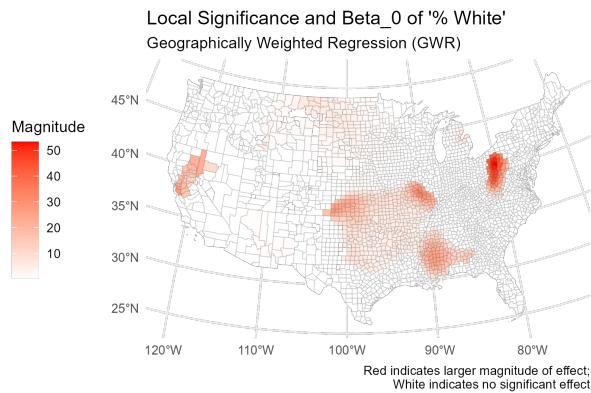
154 Figure 3d represents the local significance and magnitude of the ‘% Asian’ demographic parameter on cardiovascular
 155 disease outcomes. From the plot, a substantial portion of the map, particularly across the central to eastern regions,
 156 is colored in various shades of red. This suggests that in these areas, an increased percentage of the Asian population
 157 correlates with lower CVD outcomes. However, there are some parts of the West where higher percentages of the
 158 Asian population are associated with an increase in CVD outcomes.

159 Figure 4 represents the local significance and magnitude of the “%p2.5 air quality” parameter on cardiovascular disease
 160 outcomes. From the plot, in the southern and northeastern regions, there increased levels of PM2.5 colored as red in

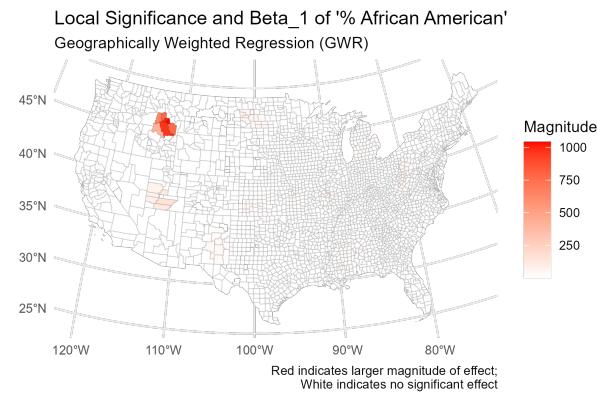
161 the plot, meaning they are associated with higher rates of CVD. There are some regions in the West where it was
 162 negatively correlated with CVD, highlighted in blue.

163 Figure 5a represents the local significance and magnitude of the “median income” parameter on cardiovascular disease
 164 outcomes. From the plot, most of the region suggests a significant negative correlation between the median income
 165 parameter and CVD outcomes. However, some parts of Texas indicate a positive correlation to CVD outcomes.

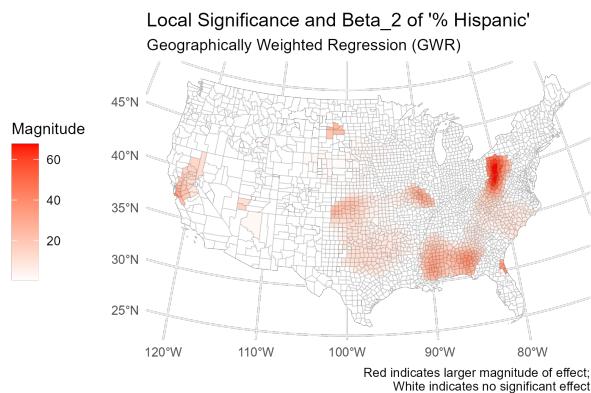
166 Figure 5b represents the local significance and magnitude of the “Unemployment” parameter on cardiovascular disease
 167 outcomes. From the plot, there is a positive correlation to CVD rates in most of the central region, and in the southwest
 168 region, there is a negative correlation.



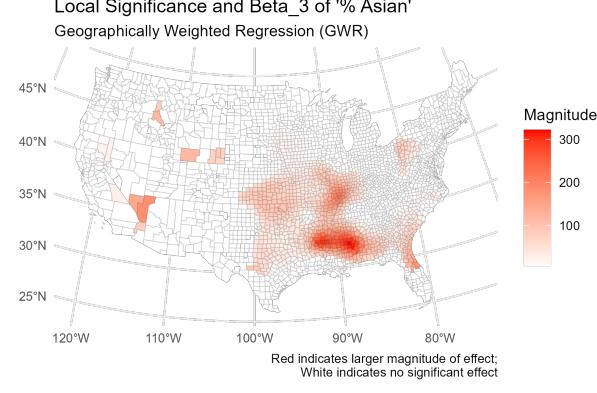
(a) Figure 3a: Strong effect in southern New England area



(b) Figure 3b: Significant positive effect in Idaho region



(c) Figure 3c: Positive effect focused in New England



(d) Figure 3d: Strong positive effects focused in western Stroke Belt, minor positive effects through region

169 5 Discussion

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

172 Our goal was to put our findings towards answering the following question: where are the socioeconomic and envi-
 173 ronmental factors affecting CVD rates in the United States?

Local Significance and Beta_4 of PM2.5 Parameter Geographically Weighted Regression (GWR)

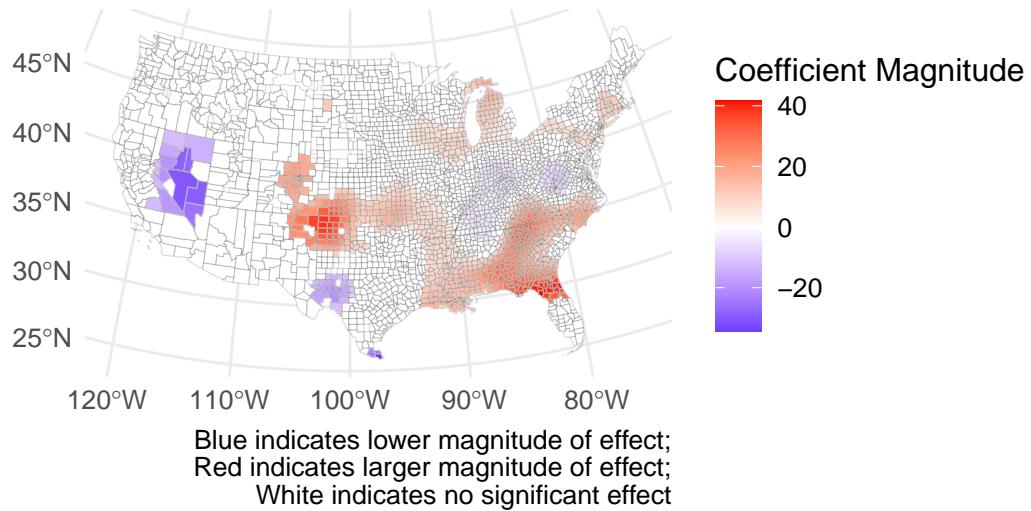
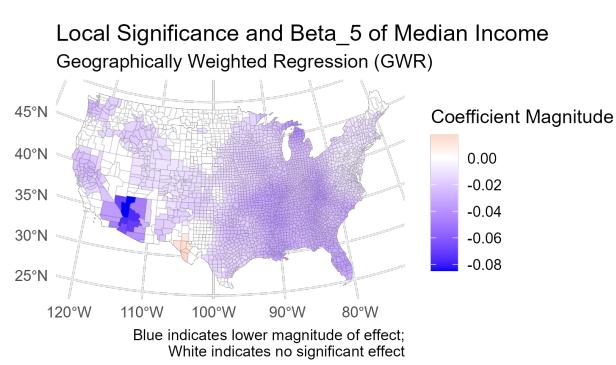
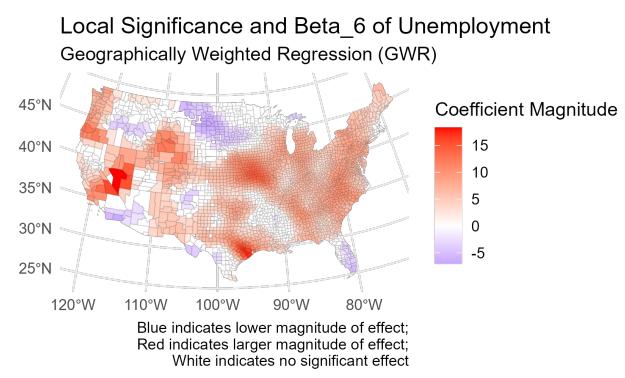


Figure 4: Higher effect in south and northeast. Negative effect in the west.



(a) Figure 5a: As expected, median income has a consistent negative relationship with CVM



(b) Figure 5b: Consistently positive relationship in central regions, however negative patches exist in Florida and northern states

Table 1

%White	%Black	%Hispanic	%Asian	PM2.5	MedIncome	Unemployment
183.00	196.00	7.00	3.00	1.00	2.00	3.00
81.00	94.00	4.00	3.00	2.00	3.00	3.00
302.00	314.00	17.00	4.00	1.00	2.00	2.00
141.00	151.00	5.00	2.00	1.00	2.00	3.00
132.00	128.00	8.00	3.00	1.00	2.00	2.00
283.00	294.00	14.00	4.00	1.00	2.00	2.00

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

180 Another factor to consider from the plots is the percentage of each race that inhabits each county. We can observe
 181 that for African Americans, Hispanics, and Asians that the impacts are significant in localized areas of the country,
 182 which may reflect underlying health disparities in access to medical care. One limitation of race percentages is the
 183 misreporting of medical records affecting minorities (Tabb et al. 2020). However, this oversight lends further credence
 184 to the fact that intervention is needed in order to combat racial health disparities. We can look at the Variance Inflation
 185 Factor (VIF) to determine if multi-collinearity has an impact on our results:

186 Typically, when examining the VIF of a GWR model, we want the coefficients for each of the seven columns—which
 187 represent our seven independent variables—to be less than 15. We can see in Table 1 that this standard holds for most
 188 cases of the columns outside the % White column, which has extremely high values. This represents a limitation in
 189 our model, as it shows multicollinearity with respect to the white and black variables, however, this is to be expected
 190 given that the racial variables we are using sum up to one. An adjustment to the model could be made in a future study
 191 to eliminate the % White variable and aggregate the other three into a single % Non-White variable. .

192 The concentration of PM2.5 and the consequential reduction in air quality has a significant impact localized in the
 193 central and southeastern regions of the United States, suggesting environmental health concerns that might require
 194 region-specific intervention. This aligns with previous research done on the stroke belt, where factors such as smoking,
 195 limited access to healthcare, and food insecurity (Zelko et al. 2023) combined with the high concentrations of PM2.5
 196 have created a hotspot of CVM far surpassing the rates of the rest of the country.

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

200 **6 Appendix**

201 The R code used for this study can be found in our public GitHub repository: [https://github.com/jbooc117/STAT489-
202 Project.git.](https://github.com/jbooc117/STAT489-Project.git)

203 **6.1 Data Sets**

204 The CDC dataset containing Medicare claims: <https://www.cdc.gov/dhdsp/maps/hd-stroke-mortality-dashboard.htm>
205 PM2.5 concentrations from 2015 dataset: [https://sedac.ciesin.columbia.edu/data/set/sdei-global-annual-gwr-pm2-5-
206 modis-misr-seawifs-aod-v4-gl-03/data-download](https://sedac.ciesin.columbia.edu/data/set/sdei-global-annual-gwr-pm2-5-modis-misr-seawifs-aod-v4-gl-03/data-download)

207 **6.2 Simulation Study Figures**

Table 2: Moran's I test demonstrates spatial correlation in our model

Moran I Statistic	Expectation	Variance
0.3010445156	0.0003274394	0.0001178233

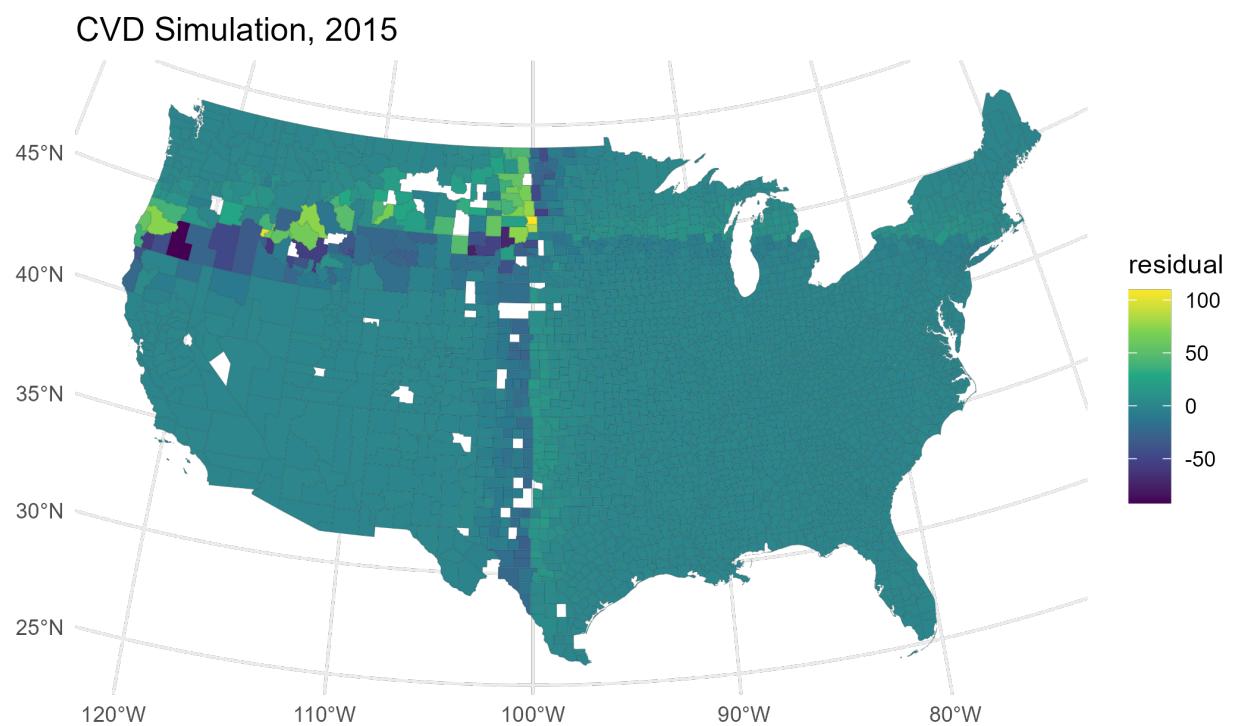


Figure 6: Division of country into four quadrants shows the effectiveness of the GWR model

208 **References**

- 209 Gebreab, Samson Y., and Ana V. Diez Roux. 2012. "Exploring Racial Disparities in CHD Mortality Between Blacks
 210 and Whites Across the United States: A Geographically Weighted Regression Approach." *Health & Place* 18 (5):
 211 1006–14. <https://doi.org/10.1016/j.healthplace.2012.06.006>.
- 212 Gollini, Isabella, Binbin Lu, Christopher Brunsdon, and Paul Harris. 2015. "{GWmodel}: An {r} Package for
 213 Exploring Spatial Heterogeneity Using Geographically Weighted Models" 63. <https://doi.org/10.18637/jss.v063.i17>.
- 215 Liu, Yi, Jingjie Sun, Yannong Gou, Xiubin Sun, Dandan Zhang, and Fuzhong Xue. 2020. "Analysis of Short-
 216 Term Effects of Air Pollution on Cardiovascular Disease Using Bayesian Spatio-Temporal Models." *International
 217 Journal of Environmental Research and Public Health* 17 (3): 879. <https://doi.org/10.3390/ijerph17030879>.
- 218 Marlow, Hilary F. 1994. "The Pharmaceutical Industry Viewpoint." *Cardiology* 85 (1): 102–12. <https://doi.org/10.1159/000176769>.
- 220 Singh, Gitanjali M., Ninon Becquart, Melissa Cruz, Andrea Acevedo, Dariush Mozaffarian, and Elena N. Naumova.
 221 2019. "Spatiotemporal and Demographic Trends and Disparities in Cardiovascular Disease Among Older Adults
 222 in the United States Based on 181 Million Hospitalization Records." *Journal of the American Heart Association* 8
 223 (21): e012727. <https://doi.org/10.1161/JAHA.119.012727>.
- 224 Tabb, Loni Philip, Angel Ortiz, Suzanne Judd, Mary Cushman, and Leslie A. McClure. 2020. "Exploring the Spatial
 225 Patterning in Racial Differences in Cardiovascular Health Between Blacks and Whites Across the United States:
 226 The REGARDS Study." *Journal of the American Heart Association* 9 (9): e016556. <https://doi.org/10.1161/JAHA.120.016556>.
- 228 Terry, Katrina, Mohamed Makhlouf, Salah E. Altarabsheh, Vaishali Deo, Fanny Petermann-Rocha, Yakov Elgudin,
 229 Khurram Nasir, Sanjay Rajagopalan, Sadeer Al-Kindi, and Salil Deo. 2023. "Trends in Cardiovascular Disease
 230 Mortality by County-Level Social Vulnerability Index in the United States." *Journal of the American Heart Asso-*
231 ciation 12 (20): e030290. <https://doi.org/10.1161/JAHA.123.030290>.
- 232 Warsito, Budi, Hasbi Yasin, Dwi Ispriyanti, and Abdul Hoyyi. 2018. "Robust Geographically Weighted Regression of
 233 Modeling the Air Polluter Standard Index (APSI)." *Journal of Physics: Conference Series* 1025 (May): 012096.
 234 <https://doi.org/10.1088/1742-6596/1025/1/012096>.
- 235 Zelko, Andrea, Pedro R. V. O. Salerno, Sadeer Al-Kindi, Fredrick Ho, Fanny Petermann Rocha, Khurram Nasir, Sanjay
 236 Rajagopalan, Salil Deo, and Naveed Sattar. 2023. "Geographically Weighted Modeling to Explore Social and
 237 Environmental Factors Affecting County-Level Cardiovascular Mortality in People With Diabetes in the United
 238 States: A Cross-Sectional Analysis." *The American Journal of Cardiology* 209 (December): 193–98. <https://doi.org/10.1016/j.amjcard.2023.09.084>.