

Distribution of the environmental and socioeconomic risk factors on COVID-19 death rate across continental United States: A spatial nonlinear analysis

Yaowen Luo

Wuhan University

Jianguo Yan (✉ jgyan@whu.edu.cn)

Wuhan University

Stephen McClure

Wuhan University

Research Article

Keywords: COVID-19 death rate; environment; socioeconomic; health; local nonlinear model; spatial variation

Posted Date: August 20th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-61369/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Distribution of the environmental and socioeconomic risk factors on COVID-19**
2 **death rate across continental United States: A spatial nonlinear analysis**

3 Yaowen Luo^{a,b,*}, Jianguo Yan^b and Stephen McClure^b

4 ^a *Electronic Information School, Wuhan University, Wuhan, China*

5 ^b *State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan*
6 *University, Wuhan, China*

7 *Corresponding author email: luoyw@lreis.ac.cn

8

9 Address:

10 Electronic Information School Wuhan University
11 127 Luoyu Road
12 Wuhan, Hubei 430070

13

14 **Abstract**

15 The COVID-19 outbreak has become a global pandemic. Spatial variation in the environmental, health,
16 socioeconomic, and demographic risk factors of COVID-19 death rate is not well understood. Global
17 models and local linear models were used to estimate the impact of risk factors of the COVID-19, but
18 these do not account for the nonlinear relationships between the risk factors and the COVID-19 death rate
19 at various geographical locations. We proposed a local nonlinear nonparametric regression model named
20 geographically weighted random forest (GW-RF) to estimate the nonlinear relationship between COVID-
21 19 death rate and 47 risk factors derived from US Environmental Protection Agency, National Center for
22 Environmental Information, Centers for Disease Control and the US census. The COVID-19 data were
23 employed to a global regression model random forest (RF) and a local model GW-RF. The adjusted R^2 of
24 the RF is 0.69. The adjusted R^2 of the proposed GW-RF is 0.78. The result of GW-RF showed that the
25 risk factors (i.e. going to work by walking, airborne benzene concentration, householder with a mortgage,
26 unemployment, airborne $PM_{2.5}$ concentration and percent of the black or African American) have a high
27 correlation with the spatial distribution of the COVID-19 death rate and these key factors driven from the
28 GW-RF were mapped, which could provide useful implications for controlling the spread of COVID-19
29 pandemic.

30 **Keywords:** COVID-19 death rate; environment; socioeconomic; health; local nonlinear model; spatial
31 variation

32 **1. Instruction**

33 The 2019 novel coronavirus disease (COVID-19) caused by SARS-CoV-2, is a rapidly spreading
34 infectious disease that mainly affects the respiratory system (Landi et al. 2020). Because the disease is
35 highly contagious with rapid transmission between human (Huang et al. 2020), World Health Organization
36 (WHO) declared on March 11 2020 that the COVID-19 outbreak has been a global pandemic (World
37 Health Organization 2020b). As of July 6, 2020, a total of 11 520 953 COVID-19 confirmed cases and
38 532 633 deaths were recorded worldwide. The current epicentre of the COVID-19 is the United States
39 with 2 982 928 confirmed cases and 132 569 deaths as of July 6, 2020. The economic impact from the
40 COVID-19 crisis are unprecedented in United States with the substantial stock market shifting and
41 unemployment rate reaching the peak (O'Connor et al. 2020). The health care system is also
42 overwhelmed across the world, which are already operating at full capacity struggling to meet the demand
43 for ventilators, intensive care beds and personal protective equipment.

44 Some researches about the COVID-19 have found that various factors including environment
45 (Ahmadi et al. 2020; Bashir et al. 2020; Xu et al. 2020), socioeconomic (de León-Martínez et al. 2020;
46 Zheng et al. 2020), demographic (Serge et al. 2020) and underlying disease (Dariya and Nagaraju 2020;
47 Malik et al. 2020; Marhl et al. 2020; Ruthberg et al. 2020) may influence the transmission of COVID-19.
48 Bashir et al. (2020) found that air pollution including PM₁₀, PM_{2.5}, SO₂, NO₂, and CO are significant risk
49 factors to the COVID-19 epidemic. Tosepu et al. (2020) analysed the correlation between weather and the
50 COVID-19, and found that the average temperature was highly correlated with the COVID-19. Virus
51 carrier via public transportation played an important role in the transmission of COVID-19 (Zheng et al.
52 2020). Serge et al. (2020) found that males are about 60% more likely than females to suffer severe
53 illness or death from the COVID-19 complications. Targher et al. (2020) found that patients with diabetes
54 were at an approximately 4 times risk of having severe COVID-19.

55 With the increased availability of health care data online and the development of the spatial
56 analysis techniques, multiple analysis by GIS tool(Guliyev 2020; Rosenkrantz et al. 2020) found that the
57 distribution of COVID-19 cases (Desjardins et al. 2020; Lau et al. 2020; Shim et al. 2020) and its risk
58 factors (Mollalo et al. 2020) exhibit patterns of spatial heterogeneity. A study by Lau et al. (2020) showed

59 that the number of flight routes was a highly relevant factor of the COVID-19 spread. Their study
60 showed that regions in Asia, North America and Europe were at a serious risk of constant exposure to
61 highly infected countries, while the exposure risk to the COVID-19 was relatively low in South America
62 and Africa. Liu et al. (2020) employed a contact model to reconstruct the contact and air spread to
63 simulate the outbreak of COVID-19 on the “Diamond Princess”. They suggested rigorous prevention
64 measure should be followed by the high-risk susceptible people. Mollalo et al. (2020) mapped the spatial
65 variability of the relationships between COVID-19 incidence rate and income inequality, median
66 household, the proportion of black females and proportion of nurse practitioners using multiscale
67 geographically weighted regression (MGWR).

68 Many mathematical models have been employed to explore the risk factors of COVID-19.
69 Typical global models such as Partial correlation coefficient (PCC) (Ahmadi et al. 2020), Ordinary least
70 squares (OLS), Poisson regression model (Xu et al. 2020), Bayesian hierarchical model (Millett et al.
71 2020) and geographical local model such as geographically weighted regression model (GWR) (Imran et
72 al. 2015; Mollalo et al. 2020) were used to model the correlations between COVID-19 data and other
73 impacting factors. However, the global model assumes the relationship between risk factors do not vary
74 over space and is inconsistent with the imbalanced distribution of COVID-19. Although Spatial Error
75 Model (SEM) and Spatial Lag Model (SLM) do consider spatial factors, they focus more on the analysis
76 of spatial correlation and do not analyze the spatial variation of the relationships between variables in
77 different regions from the perspective of spatial heterogeneity (Ahmadi et al. 2020). The GWR (Brunsdon
78 et al. 2010; Fotheringham et al. 2002; Lu et al. 2017) as a local regression model can obtain the linear
79 relationship between variables in different locations. However, the GWR is constructed based on multiple
80 linear regression models, thus it is not suitable to estimate the nonlinear relationships between
81 independent and dependent variables, and local multicollinearity exists when dealing with correlated
82 variables (Wheeler and Tiefelsdorf 2005). The real relationship between risk factors and COVID-19 is
83 complex and is not always linear. In order to explore the spatial variation of the nonlinear relationship
84 between multiple risk factors and COVID-19, it is necessary to deal with the nonlinear situation in a local
85 regression model.

86 In this study, we proposed a local nonlinear nonparametric regression method, Geographically
87 Weighted Random Forest (GW-RF), to evaluate the relationship between COVID-19 death rate and
88 multiple risk factors including air pollution, climate, landcover, disaster, health status, commuting to

89 work, socioeconomic and demographic indicators at county level across the continental United States.
 90 This paper first tries to explore the variation in the nonlinear relationships between multiple risk factors
 91 and COVID-19 death rate in different locations by using the GW-RF. We expect this study can provide
 92 scientific evidence for implementing control and prevention measure in COVID-19.

93 **2. Materials and methods**

94 **2.1 Data and preparation**

95 The county level daily COVID-19 death cases data and population data of 3108 counties of the
 96 continental United State from Jan 22, 2020 to June 26, 2020 were downloaded from the website of USA
 97 FACTS (<https://usafacts.org/>). The death rate at county level was calculated based on the daily COVID-
 98 19 death cases and population data. We selected 47 indicators including atmosphere, climate, landcover,
 99 disaster, health status, commuting to work, socioeconomic and demographic factors as independent
 100 variables to evaluate their correlation with the COVID-19 death rate. The indicators we selected and their
 101 meanings and sources are presented in Table 1. The shapefile of the selected 3108 counties was
 102 downloaded from geographical program of United States Census Bureau
 103 (<https://www.census.gov/programs-surveys/geography.html>).

104 Table 1. Definitions of indicators and sources

Theme	Indicators	Indicator meaning	Source
Atmosphere	Airborne PM _{2.5} concentration	Annual average ambient concentrations of PM _{2.5} in micrograms per cubic meter	United States Environmental Protection Agency (https://www.epa.gov/) and Centers for Diseases Control and Prevention (https://www.cdc.gov/)
	Airborne benzene concentration	Annual average concentration of benzene estimates in microgram per cubic meter	
	Airborne formaldehyde concentration	Annual average air concentration of formaldehyde estimates in microgram per cubic meter	
	Airborne acetaldehyde concentration	Annual average air concentration of acetaldehyde estimates in microgram per cubic meter	

	Airborne carbon tetrachloride concentration	Annual average air concentration of carbon tetrachloride estimates in microgram per cubic meter	
Climate	Air temperature	Average Daily Max Air Temperature (F)	National Center for Environmental Information (https://www.ncei.noaa.gov/)
	Precipitation	Average Daily Precipitation (mm)	
	Sunlight exposure	Annual average sunlight exposure measured by solar irradiance (kJ/m ²)	Centers for Diseases Control and Prevention (https://www.cdc.gov/)
	UV radiation exposure	Annual average daily dose of UV irradiance (J/m ²)	
Landcover	Landcover with water	Percent of land covered by water	
	Land cover with forest	Percent of land covered by forest	
Disaster	Drought	Number of weeks of moderate drought or worse per year	
	Flood	Percentage of people within fema designated flood hazard area	
Health status	Disability	Percentage of population aged 5 years and over with a disability	
	Asthma	Percent of adults diagnosed with asthma	
	Obese	Percentage of adults aged 18 years and over who were obese	
	Overweight	Percentage of adults aged 18 years and over who were overweight	
	Cancer	Number of people with lung and brouchus cancer per 1000000 population	
Commuting to work	Go to work by private transportation	Percentage of workers 16 years and over who drove alone (car, truck, or van)	United States Census Bureau (https://www.census.gov/en.html)
	Go to work by public transportation	Percentage of workers 16 years and over who go to work by public transportation (excluding taxicab)	

	Go to work by walking	Percentage of workers 16 years and over who go to work by walking	
	Work at home	Percentage of workers 16 years and over who worked at home	
	Mean travel time to work	Mean travel time to work (minutes) of the workers 16 years and over	
Socioeconomic	Health insurance	Percentage of population without health insurance	
	Householder with a mortgage	Percentage of household with a mortgage	
	Poverty	Percentage of population whose income is below the poverty level	
	Service occupations	Percentage of employed population 16 years and over with service occupations	
	Unemployment	Percentage of population 16 years and over unemployed	
	Hospital	Number of hospitals	Centers for Diseases Control and Prevention (https://www.cdc.gov/)
	Hospital beds	Number of hospital beds per 10000 population	
	People living in group quarter	Percentage of population living in group quarter	United States Census Bureau (https://www.census.gov/en.html)
	People living near a park	Percentage of population living within a half mile of a park	
	Householder with no internet access	Percentage of households with no internet access	
	Median household income		
	Mean household retirement income		
	Mean household cash public assistance income		

Mean household
Supplemental
Security Income

Demographic Percent of males

Median age

Percent of people
under 18 years

Percent of people 65
years and over

Percent of the white
race

Percent of the black
or African American

Percent of American
Indian and Alaska
Native

Percent of Asian

Percent of native
Hawaiian and other
Pacific islander

Percent of Hispanic
or Latino

105 Due to the units of these 47 indicators are different, the indicators should be normalized before

106 regression. The method is as follows:

107
$$X_{ki} = \frac{X_{ki} - \bar{X}_k}{\sigma_k} (i \in 1, 2, \dots, 2056; k \in 1, 2, \dots, 28) \quad (1)$$

108 where X_{ki} represents the normalized value of the k th indicator in the i th county, X_{ki} represents
109 the original value of the k th indicator in the i th county; \bar{X}_k represents the average value of the k th
110 indicator; σ_k represents the standard deviation of k th index. The COVID-19 death rate and 47 indicators
111 were joined to the county level shapefile for further processing.

112 2.2 Nonlinear nonparametric model

113 2.2.1 Random forest (RF)

114 We selected the Random forest (RF) machine learning method (Breiman 2001) because it is
115 nonparametric, it can easily learn nonlinear relationships and interactions from data without explicitly
116 modelling them. RF is an ensemble of multiple decision trees. The decision tree is a non-parametric
117 model which does not have a fixed structure. The decision tree grows according to the complexity of the
118 input data in the learning process. The RF works well for high-dimensional variables with a relatively
119 small number of samples, and can access variable importance (Grömping 2009). The algorithm flow of
120 the RF is as follows.

- 121 (1) The n data sets D_1, D_2, \dots, D_n are extracted by repeatedly using the bootstrap method to
122 randomly extract the whole data set D ; and the corresponding n decision trees H_1, H_2, \dots, H_n are
123 generated.
- 124 (2) At each node of the decision tree, randomly select m ($m < k$) variables from all the k variables
125 of the decision tree, and each node is split using the selected m variables by the optimal
126 segmentation method determined by a segmentation criterion.
- 127 (3) The value of m remains unchanged while the forest grows. Each tree grows to its largest extent
128 without pruning until it cannot be split.

129 Thus the correlation between the decision trees in the forest decreases through a random
130 selection of variables at each node of the tree and the optimal split of each node is determined by the
131 selected variables only, instead of all variables. Each tree can grow to its largest extent without pruning.
132 Therefore, the algorithm can deal with excessive redundant features and avoid over fitting.

133 In the first step in constructing the RF, whether with or without replacement, approximately
134 36.8% of the data samples are not used to grow the tree; these samples are the Out-Of-Bag (OOB) for the
135 tree. The accuracy of the RF model can be estimated from the OOB data as presented by Equation (1):

$$136 \quad MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y}_i)^2 \quad (1)$$

137 where N is the number of samples which from OOB data, y_i is the actual value of the i th sample, \bar{y}_i is the
138 average prediction for the i th sample from all trees.

139 The overall sum of squares (SST) and coefficient of determination (R^2) are respectively defined
 140 in Equations (2) and (3):

$$141 \quad SST = \sum_{i=1}^N (y_i - \bar{y})^2 \quad (2)$$

$$142 \quad R^2 = 1 - N \frac{MSE}{SST} \quad (3)$$

143 where $R^2 \in (0,1)$. The closer of the value of R^2 to 1, the better the regression performance of the GW-RF
 144 will be.

145 Variable importance can sort the independent (predictor) variables according to their degree of
 146 correlation to the dependent (response) variable. There are two popular methods to measure the variable
 147 importance in the RF, which are average impurity reduction (Gini importance) and Mean Square Error
 148 (MSE) reduction. Because the result of variable importance by impurity reduction is biased (Strobl et al.
 149 2007), many researchers have verified and suggested choosing the MSE reduction method when
 150 permuting the variables (Ishwaran 2007; Strobl et al. 2008). The MSE reduction method uses the MSE
 151 value of the Out-Of-Bag (OOB) data to evaluate the variable importance (Cai et al. 2018). It is
 152 determined as follows:

153 (1) Calculate the MSE of the OOB data for each tree. For tree t , the MSE of OOB data is calculated
 154 by Equation (4):

$$155 \quad MSE_t = \frac{1}{N_t} \sum_{i=1}^{N_t} (y_i - \hat{y}_{i,t})^2 \quad (4)$$

156 where N_t is the number of samples from OOB data in the tree t , $\hat{y}_{i,t}$ is the prediction for the i th
 157 sample of the tree t .

158 (2) Randomly replace the target variable j , and then the new value of the MSE of tree t is
 159 calculated by Equation (5):

$$160 \quad MSE_t(j) = \frac{1}{N_t} \sum_{i=1}^{N_t} (y_i - \hat{y}_{i,t}(j))^2 \quad (5)$$

161 where $\hat{y}_{i,t}(j)$ is the prediction for the i th sample of the new tree t when randomly replace the
 162 target variable j .

163 (3) Calculate the difference between MSE_t and $MSE_t(j)$, and the MSE reduction is the variable
 164 importance for variable j of tree t . The MSE reduction of variable j of the whole forest is
 165 obtained as the average over MSE reduction of all n trees. The variable importance of variable j
 166 is expressed as in Equation (6):

$$167 \quad VI(j) = MSE(j) = \frac{1}{n} \sum_{t=1}^n (MSE_t - MSE_t(j)) \quad (6)$$

168 2.2.2 Geographically weighted random forest (GW-RF)

169 In this section, a local nonlinear machine learning method, denoted as GW-RF, is proposed. The GW-RF
 170 is designed by integrating spatial weight matrix (SWM) and RF into a local regression analysis
 171 framework. The GW-RF inherits the merits of the RF, making the RF from being applicable from a
 172 global system to a local system. Thus, it can handle high-dimensional variables with nonlinear
 173 relationships and multicollinearity. The variable importance for each spatial unit can be obtained from the
 174 GW-RF. The process of constructing the GW-RF model is designed as follows:

175 (1) The SWM for each spatial unit of the study area should first be made according to the specified
 176 spatial weight rule. The SWM for the whole study area with p spatial units can be expressed as
 177 in Equation (7):

$$178 \quad W = \begin{bmatrix} W(1) \\ W(2) \\ \vdots \\ W(i) \\ \vdots \\ W(p) \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1p} \\ w_{21} & w_{22} & \cdots & w_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ w_{i1} & w_{i2} & \cdots & w_{ip} \\ \vdots & \vdots & \vdots & \vdots \\ w_{p1} & w_{p2} & \cdots & w_{pp} \end{bmatrix}, \quad i \in (1, 2, \dots, p) \quad (7)$$

179 As the local random forest of an individual unit need to consider the unit itself, the value of w_{ii} is
 180 set to 1 ($w_{ii} = 1$). According to the spatial weight rule, for spatial unit i , if sample j ($j \in$
 181 $(1, 2, \dots, p) \wedge i \neq j$) is a “neighbour” of unit i , the value of spatial weight between them is set to 1,
 182 that is, $w_{ij} = 1$. While spatial unit j is far away from spatial unit i , not a “neighbour” of spatial
 183 unit i , $w_{ij} = 0$.

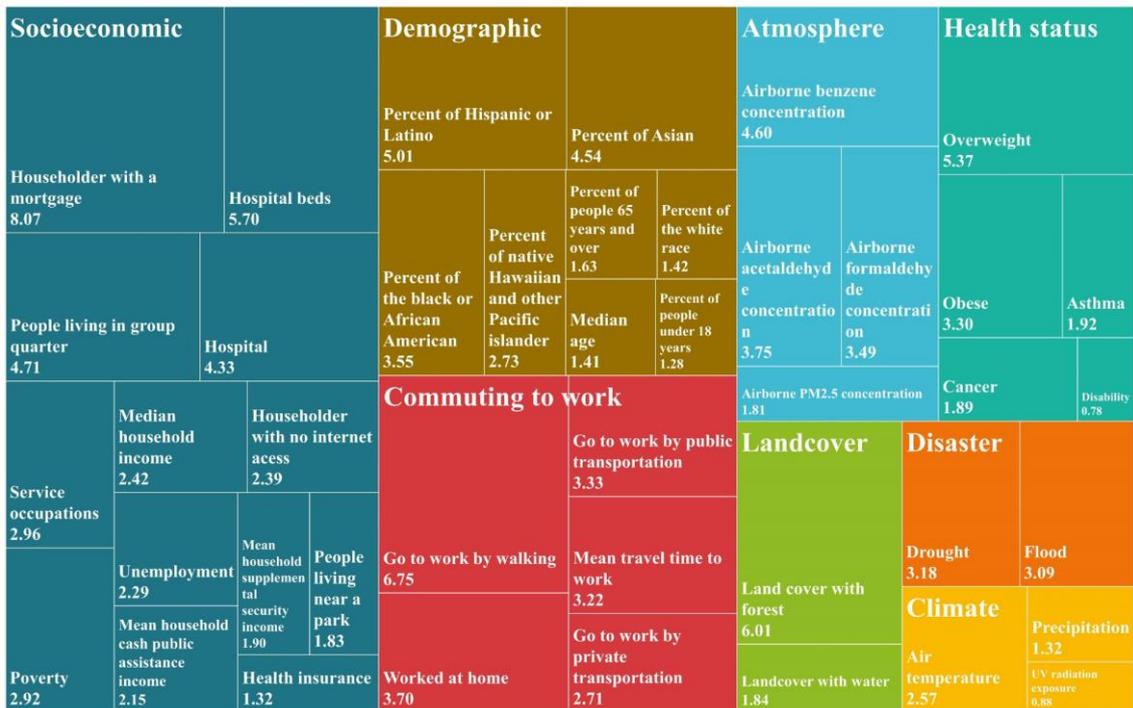
184 (2) Select all the “neighbours” of each spatial unit according to the spatial weight matrix. For unit i ,
 185 the “neighbors” of it can be selected from the special weight matrix W where $w_{ij} \neq 0$, ($j \in$
 186 $(1, 2, \dots, p) \wedge i \neq j$).

187 (3) The spatial unit i and its “neighbours” are as the inputs to construct a local RF for unit i (RF (i)).
 188 By executing RF (i), the variable importance for spatial unit i can be computed.
 189 (4) Repeat steps (2) and (3) to construct a local RF for each spatial unit in the study area and
 190 estimate the local variable importance for each spatial unit.

191 The nonlinear nonparametric models (RF, GW-RF) does not need to consider multicollinearity,
 192 and can analyse all independent variables without screening. R software (version 3.5.3, [http://cran.r-](http://cran.r-project.org)
 193 [project.org](http://cran.r-project.org)) was used to perform the regression analysis.

194 **3. Results**

195 All 47 indicators were employed to the nonlinear nonparametric models (RF, GW-RF). The adjusted
 196 fitting coefficient (R^2) of the RF was 0.69, while the adjusted R^2 of the GW-RF was 0.78, indicating that
 197 the regression result of the GW-RF was more accurate than that of the RF. The variable importance of a
 198 independent variable represents the correlation between the independent variable and the dependent
 199 variable and the higher value of the variable importance is, the stronger the correlation will be. The
 200 variable importance of 47 independent variables in modelling COVID-19 death rate using the RF is
 201 showed in Figure 1. The risk factors referring to socioeconomic are most correlated with COVID-19
 202 death rate, followed by risk factors referring to demographic, commuting to work, atmosphere, health
 203 status, landcover, disaster and climate. The variables including householder with a mortgage, going to
 204 work by walking, land cover with forest, hospital beds, overweight, percent of Hispanic or Latino, people
 205 living in group quarter and Airborne benzene concentration have a high correlation with the COVID-19
 206 death rate.



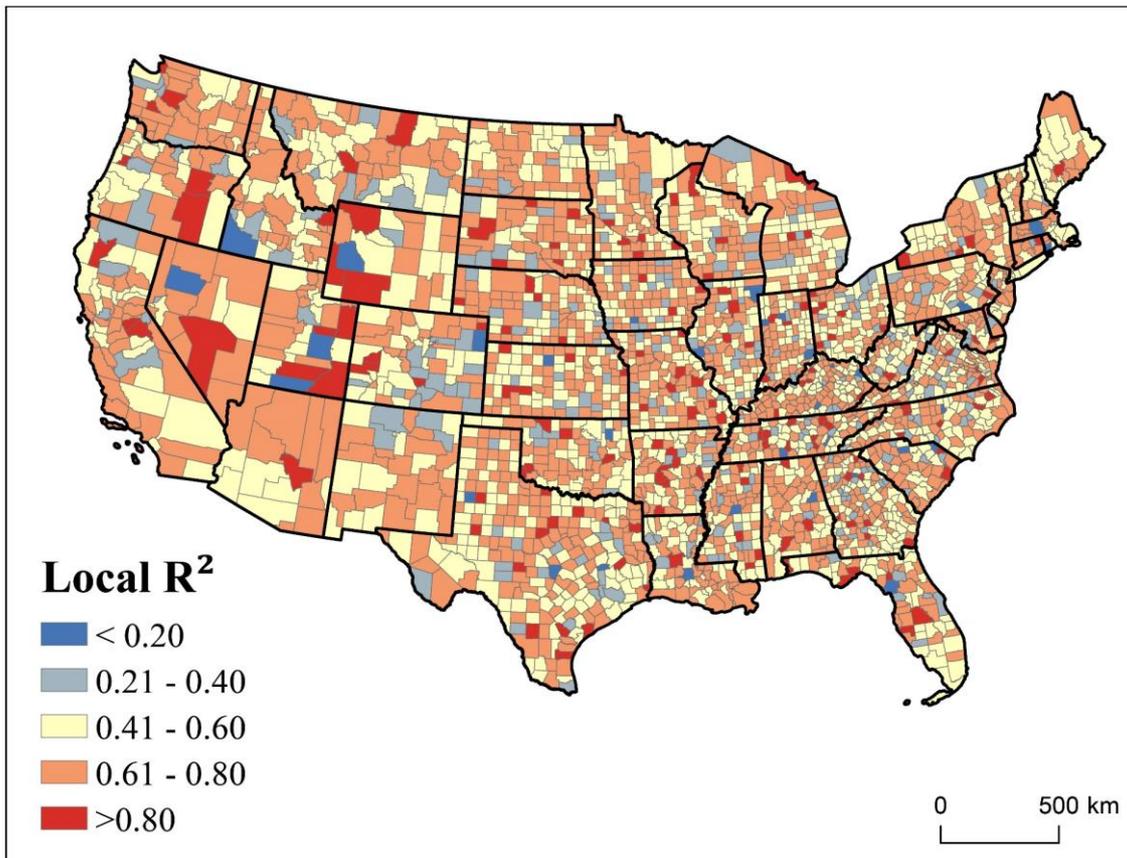
207 Figure 1. The variable importance of the independent variables of the RF model in modelling COVID-19
 208 death rate

209 We used the local R^2 to estimate the performance of the GW-RF. Table 2 describes the statistic
 210 of local R^2 of the GW-RF. The average value of local R^2 was 0.59. The value of local R^2 was higher than
 211 0.4 in 89.4% of the counties and higher than 0.6 in 50.5% of the counties. This shows that the GW-RF
 212 can accurately evaluate the correlation between the risk factors and the COVID-19 death rate in most of
 213 the study areas.

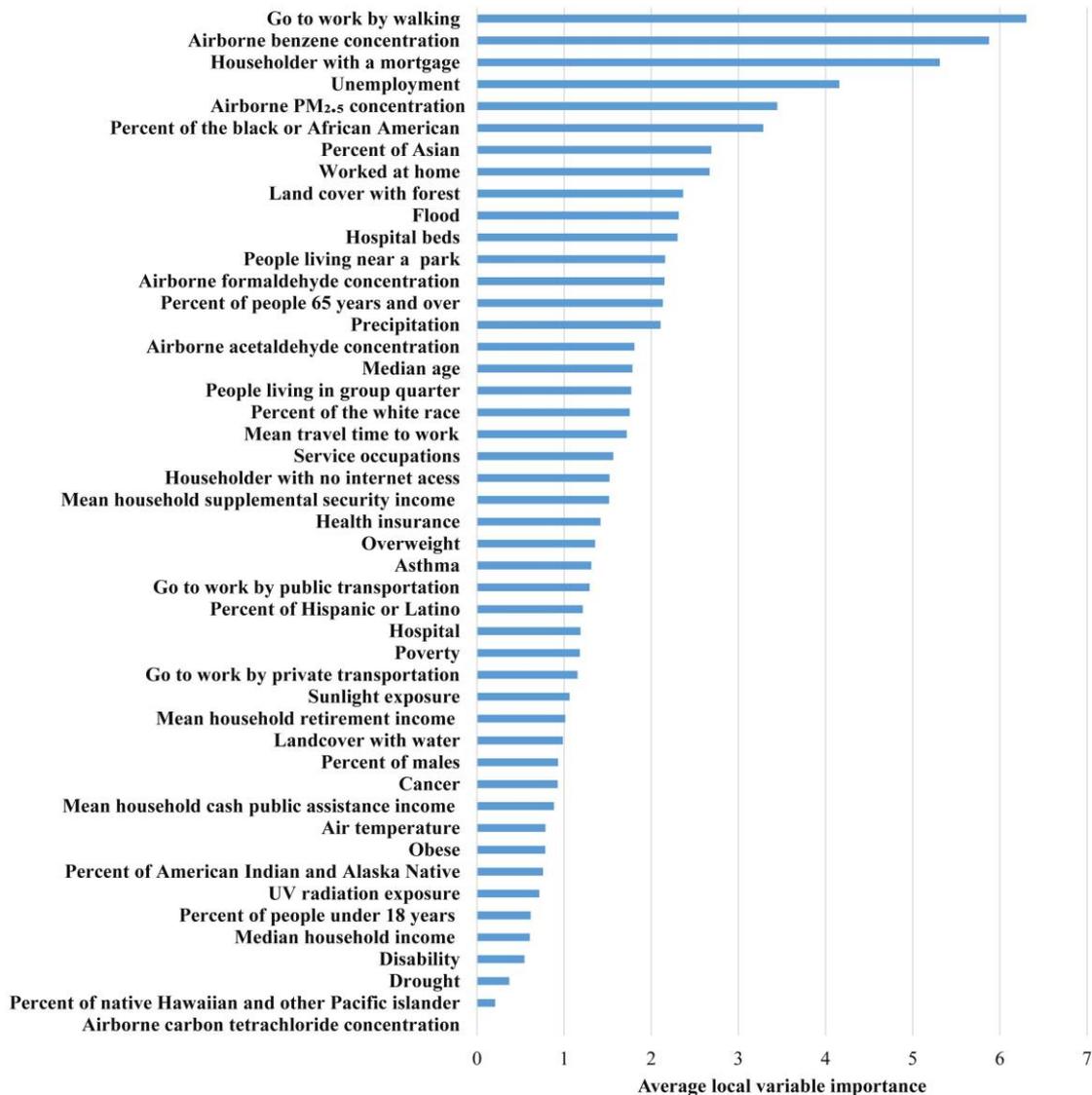
214 Table 2. The statistic of local R^2 of the GW-RF in modelling COVID-19 death rate, we calculated the
 215 average value of local R^2 and the percentage of counties in five local R^2 range (≤ 0.2 , (0.2, 04], (0.4, 06],
 216 (0.6, 08], >0.8)

The value of local R^2	GW-RF
Average value	0.59
≤ 0.2	1.1%
(0.2, 04]	9.5%
(0.4, 06]	38.9%
(0.6, 08]	44.8%
>0.8	5.7%

217 Figure 2 shows the distribution of the local R^2 of the GW-RF across the study area. As can be
 218 seen from Figure 2, the distribution of local R^2 was imbalanced in the whole study area. The local R^2
 219 value of the GW-RF was high in most of the counties across the whole continental United States,
 220 indicating that the GW-RF worked well in the prediction of the local COVID-19 death rate in most
 221 regions across the study area, especially in Nevada, Arizona, Washington and some counties in the East-
 222 central region.



223 Figure 2. The distribution of local R^2 of the GW-RF
 224 We computed the average local effect of each independent variable on COVID-19 death rate in
 225 the GW-RF model (see Figure 3). The effect of going to work by walking had the highest correlation with
 226 the COVID-19 death rate, followed by airborne benzene concentration, householder with a mortgage,
 227 unemployment, airborne $PM_{2.5}$ concentration and percent of the black or African American.



228

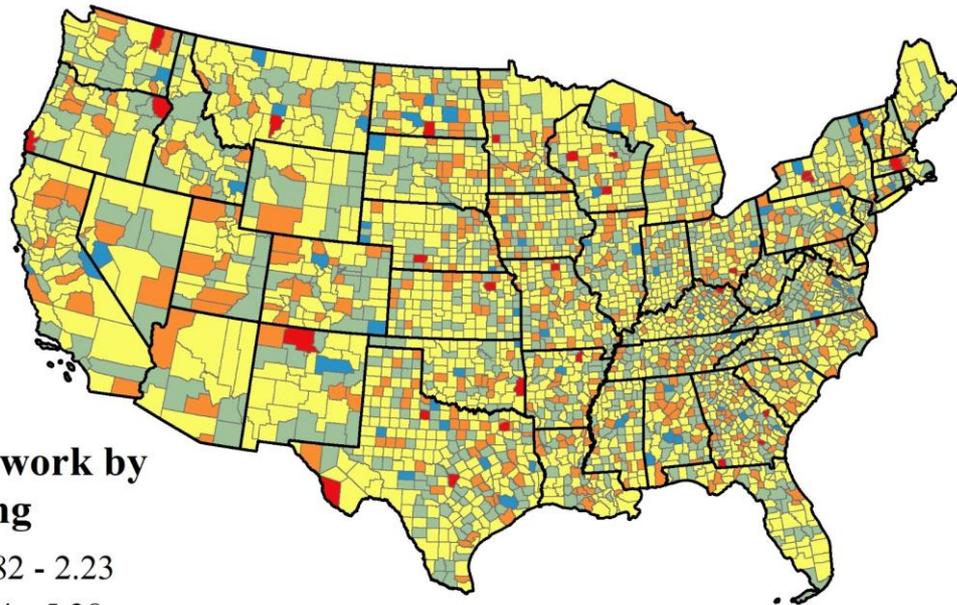
229 Figure 3. The average local variable importance of 47 potential risk factors on COVID-19 death rate in
 230 the GW-RF model

231 The proportion of counties with local primary risk factor (the risk factor with the highest value
 232 of local variable importance) at county level in the GW-RF was calculated (see Table 3). Going to work
 233 by walking was the most influential risk factor in 35% of the counties. The airborne benzene
 234 concentration was the leading risk factor in 24% of the counties. 13% percent of counties were most
 235 affected by householder with a mortgage and 12% percent of counties were most affected by
 236 unemployment. Figure 4, Figure 5 and Figure 6 provide a detailed spatial distribution of the local
 237 variable importance of first six factors with the highest value of average variable importance on the
 238 COVID-19 death rate using the GW-RF.

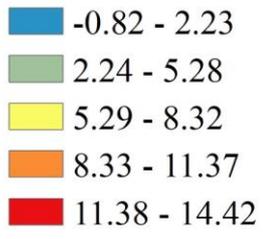
239 Table 3. The proportion of counties with local primary risk factor (the risk factor with the highest value of
240 local variable importance) on COVID-19 death rate at county level in the GW-RF.

Local primary risk factor	Proportion of counties
Go to work by walking	35%
Airborne benzene concentration	25%
Householder with a mortgage	13%
Unemployment	12%
Other risk factors	16%

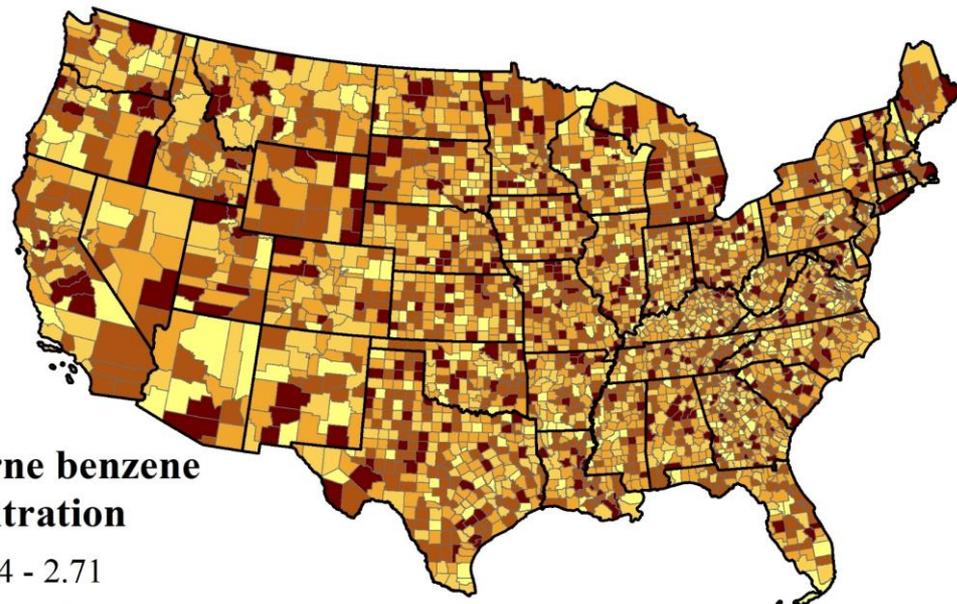
A



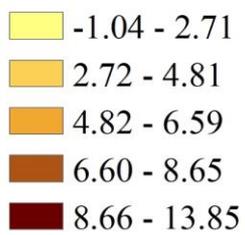
Go to work by walking



B



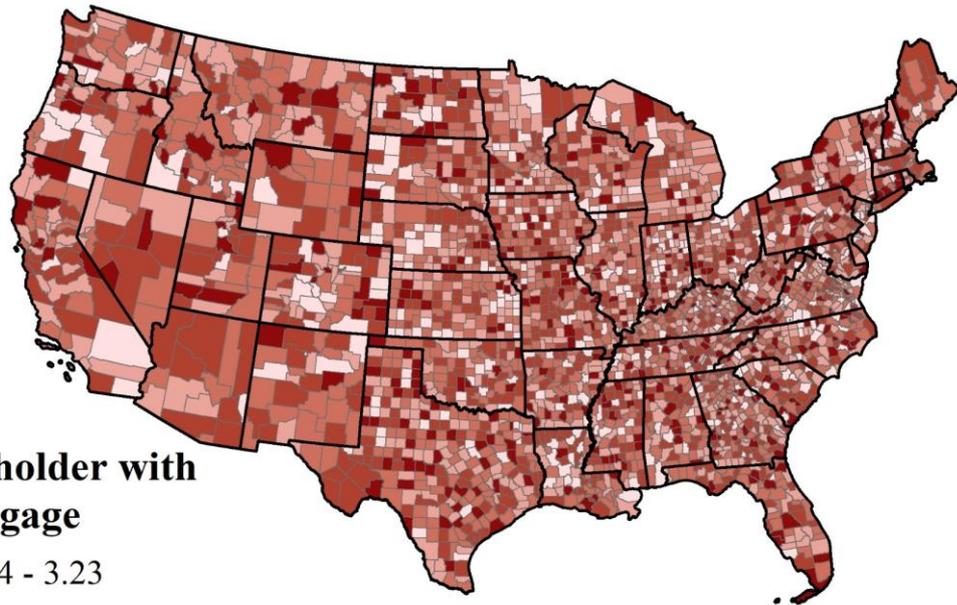
Airborne benzene concentration



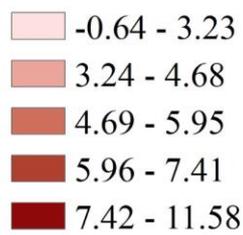
241

242 Figure 4. The spatial distribution of the local variable importance of (A) go to work by walking, (B)
243 airborne benzene concentration on COVID-19 death rate in GW-RF model

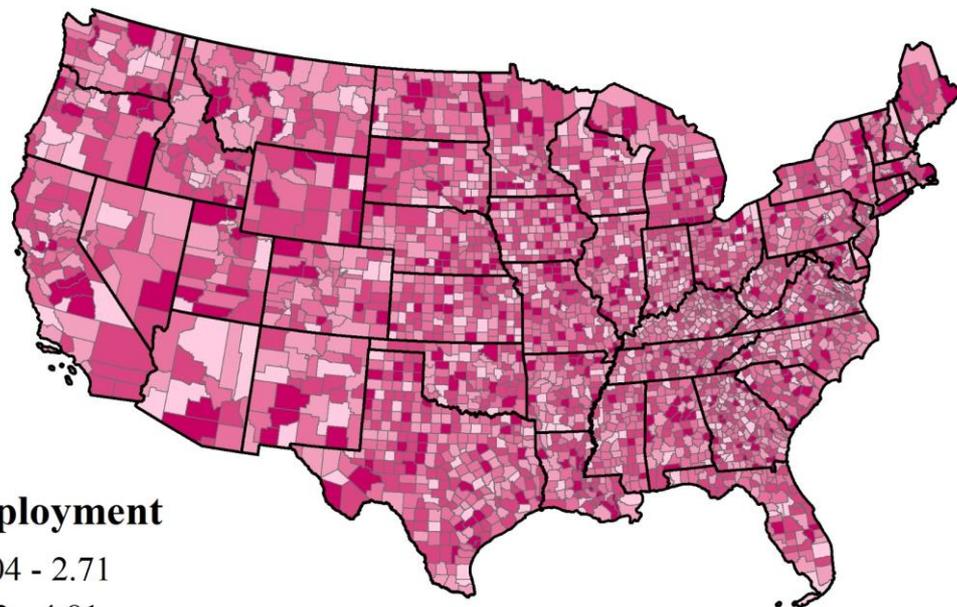
A



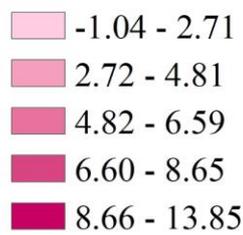
Householder with a mortgage



B



Unemployment



244

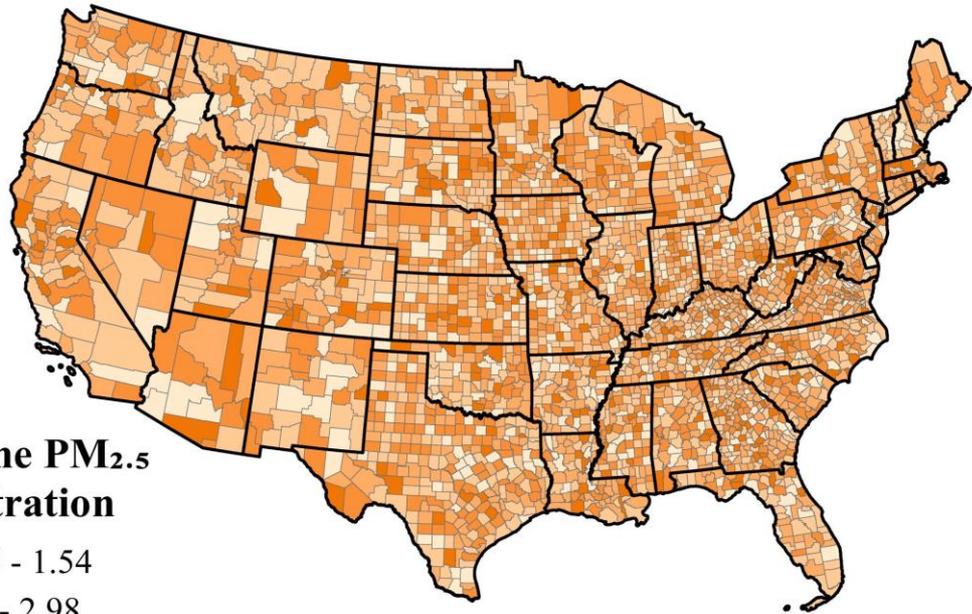
245

246

247

Figure 5 The spatial distribution of the local variable importance of (A) householder with a mortgage, (B) unemployment on COVID-19 death rate in GW-RF model

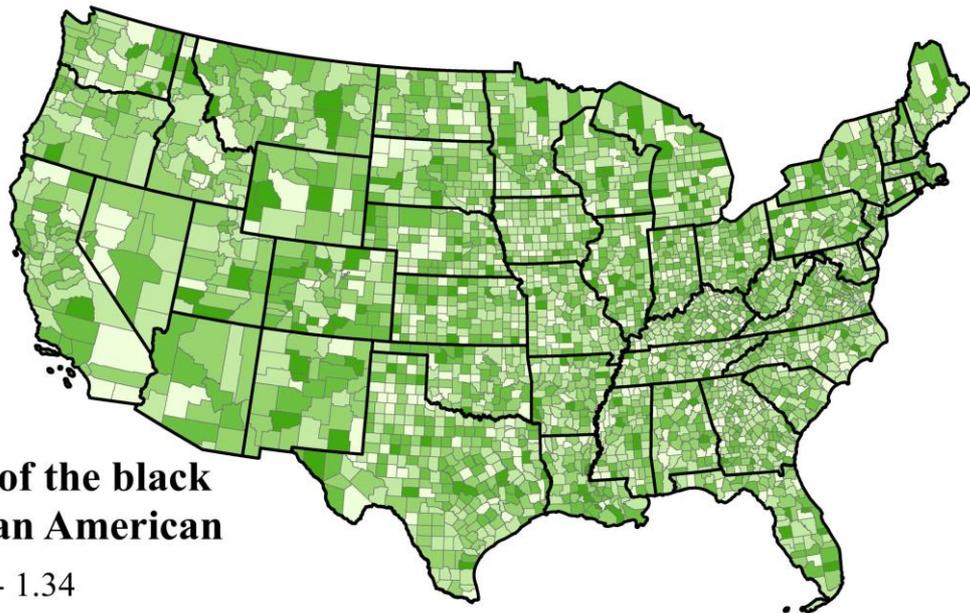
A



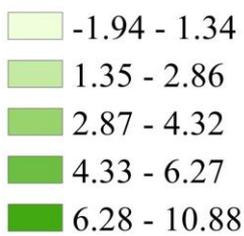
**Airborne PM_{2.5}
concentration**



B



**Percent of the black
or African American**



248

249

250

251

252

Figure 6 The spatial distribution of the local variable importance of (A) airborne PM_{2.5} concentration, (B) percent of the black or African American on COVID-19 death rate in GW-RF model

From Figure 4, Figure 5 and Figure 6, the distribution of the variable importance of each variable on COVID-19 death rate in GW-RF model was imbalanced in different counties even the

253 counties in a same state. For example, in the southern part of Arizona, the COVID-19 death rate was
254 mainly affected by the airborne benzene concentration and unemployment, and the northern part was
255 mainly affected by going to work by walking and airborne PM_{2.5} concentration. The regions obviously
256 effected by going to work by walking were distributed in California, Arizona, the west of Utah, South
257 Carolina and Massachusetts. The areas influenced by airborne benzene concentration were scattered
258 throughout the study area. New Mexico, Florida, Texas, Missouri, the south of Nevada, the north of
259 Arizona, Massachusetts and Connecticut were sensitive to householder with a mortgage. The regions
260 obviously effected by airborne PM_{2.5} concentration and percent of the black or African American are
261 similar, mainly located in the north of Nevada, the north of Arizona, the southeast of Oregon, the east of
262 Wyoming and the central part of the continental United States. In addition, the same area was affected by
263 several risk factors. For example, airborne benzene concentration, householder with a mortgage,
264 unemployment and the percent of black of African American were influential factors in the southeast of
265 Arizona.

266 **4. Discussion and conclusion**

267 Identifying the risk factors that highly correlated with the transmission will provide guidance in
268 containing the spread of the COVID-19 disease. In this study, we selected 47 potential risk factors from
269 atmosphere, climate, land cover, disaster, health status, commuting to work, socioeconomic and
270 demographic categories as independent variables to estimate their impact on the distribution of the
271 COVID-19 death rate at county level across continental United States. Due to the imbalanced distribution
272 of COVID-19 death rate and the complex relationship between the COVID-19 death rate and its risk
273 factors, the linear models could not accurately identify the key risk factors in different locations. To solve
274 this problem, we applied GW-RF, a local regression model capable of identifying nonlinear relationships
275 between variables at various geographical locations and suitable for dealing with high-dimensional
276 variables even for correlated variables.

277 In this study, we used two nonlinear regression models (RF, GW-RF) to identify the key risk
278 factors to the COVID-19 death rate. The result showed that the nonlinear models effectively modelled the
279 relationship between the risk factors and COVID-19 death rate both in global and local regressions. The
280 adjusted R^2 of the GW-RF was 0.78, higher than that of the RF, indicating the GW-RF is more suitable to
281 estimate the local risk factors of the COVID-19 death rate compared with the global model RF. The

282 average value of local R^2 of the GW-RF is 0.59. In GW-RF, the value of local R^2 higher than 0.4 in
283 89.4% of the counties and higher than 0.6 in 50.5% of the counties, indicating that the GW-RF performed
284 well in most of the study area. This shows that that the local nonlinear nonparametric model GW-RF can
285 accurately estimate the relationship between the risk factors and COVID-19 death rate at various
286 geographical locations.

287 Our result shows that several risk factors from environment, socioeconomic, demographic and
288 commuting to work are associated with the COVID-19 death rate. Finding of the global model RF
289 showed that householders with a mortgage had a highest correlation with the number of COVID-19 death
290 rate, followed by going to work by walking, land cover with forest, hospital beds, and overweight.
291 Findings of the geographical local model GW-RF is similar to that of the RF, but a little different. The the
292 GW-RF results show that going to work by walking, airborne benzene concentration, householder with a
293 mortgage, unemployment, airborne $PM_{2.5}$ concentration and percent of the black or African American
294 played an important role in the distribution of the COVID-19 death rate. Most of our findings are
295 consistent with previous research on COVID-19. Zheng et al. (2020) found that the frequency of public
296 transportation including flights, trains, and buses from the epicentre are important determinants of
297 transmission risks of COVID-19. They suggested preventive measures should be taken in public
298 transportation in order to contain the COVID-19 epidemic. Several studies found that the air pollution
299 have a significant correlation with the COVID-19 confirmed cases (Bashir et al. 2020; Xu et al. 2020).
300 Viruses are usually not spread as independent individuals in air, they are more likely to attach to other
301 suspended particles (Yang et al. 2011). Therefore, the concentration of air pollutants may affect the
302 aerosol transmission of SARS-CoV-2. These studies encouraged the formulation of environmental
303 policies to control pollution sources, which can reduce the harmful effects of air pollutants. Mollalo et al.
304 (2020) found that the proportion of black females and median household income had significant influence
305 on the spatial distribution of the COVID-19 incidence rate.

306 By exploring the spatial distribution of risk factors of the COVID-19 death rate, we found that
307 COVID-19 death rate in each region was affected by various factors and the association between each
308 risk factor and the COVID-19 death rate was not consistent in different spatial locations. About 35% of
309 the counties are most affected by Going to work by walking, so it is necessary to call on people to pay
310 attention to the social distancing and to wear medical masks. The western and central east region were
311 affected by the airborne benzene concentration; toxic particles in the air affect the spread of viruses.

312 Therefore, these regions should pay attention to the impact of air pollution on human health and take
313 measures to protect the environment. The southern part of the continental United States was heavily
314 affected by the proportion of the black or African American and householder with a mortgage, so there
315 are some assistance probably can be taken in these regions to provide people with financial help such as
316 food and medical supplies.

317 The current research, despite showing the spatial variability of the correlation between multiple
318 risk factors and the COVID-19 death rate at a county level, has the following limitations. First, the current
319 study only focused on the spatial dimension of the data based on a period, but the data about the COVID-
320 19 death rate is constantly changing over time. Future study can research in the spatiotemporal direction.
321 Secondly, we have not accounted for policy factors at local area. Policy factors would be an interesting
322 research contribution to the transmission of COVID-19. Thirdly, the GW-RF model only assess the
323 goodness-of-fit test of the regression, but does not assess the significance of the single variable. The test
324 method of this model needs to be improved in the future study.

325 At present, few geographic local models study the nonlinear relationship between variables. The
326 proposed GW-RF model could accurately estimate the spatial variability of nonlinear relationship
327 between the risk factors and COVID-19 death rate, thus this method is applicable in many use instances
328 where this is an issue about selecting significantly correlated variables at various geographical locations.
329 Our results confirmed the findings of existing work on COVID-19, but extends it by using a nonlinear
330 approach to quantify the impact of risk factors relevant in local areas. We expect this study could provide
331 a reference for the geographical local nonlinear modelling in the future epidemiological studies.

332 **Declaration of Competing Interest**

333 The authors declare that they have no known competing financial interests or personal relationships that
334 could have appeared to influence the work reported in this paper.

335 **Credit authorship contribution statement**

336 Yaowen Luo: Conceptualization, Methodology, Software, Formal analysis, Investigation, Writing -
337 Original Draft. Jianguo Yan: Conceptualization, Validation, Writing - Review & Editing, Supervision,
338 Funding acquisition. Stephen McClure: Validation, Writing - Review & Editing.

339 **Acknowledgements**

340 This research is supported by a grant provided by National Scientific Foundation of China (Grant No.

341 U1831132 and 41374024), Innovation Group of Natural Fund of Hubei province (Grant No.
342 2018CFA087) and the fundamental research funds for the central universities (2042018KF0231).

343 **Consent for publication**

344 All the co-authors consent the publication of this work.

345 **Declaration of Competing Interest**

346 The authors declare that they have no known competing financial interests or personal relationships that
347 could have appeared to influence the work reported in this paper.

348

349

350

351 **Reference**

- 352 Ahmadi M, Sharifi A, Dorosti S, Jafarzadeh Ghouschi S, Ghanbari N (2020)
353 Investigation of effective climatology parameters on COVID-19 outbreak in
354 Iran. *Science of The Total Environment* 729:138705
355 doi:<https://doi.org/10.1016/j.scitotenv.2020.138705>
- 356 Bashir MF, Bilal BMA, Komal B (2020) Correlation between environmental pollution
357 indicators and COVID-19 pandemic: A brief study in Californian context.
358 *Environmental Research*:109652
359 doi:<https://doi.org/10.1016/j.envres.2020.109652>
- 360 Breiman L (2001) *Random Forests* vol 45. doi:10.1023/A:1010933404324
- 361 Brunsdon C, Fotheringham AS, Charlton ME (2010) Geographically Weighted
362 Regression : A Method for Exploring Spatial Nonstationarity. *geographical*
363 *analysis* 28:281-298 doi:10.1111/J.1538-4632.1996.TB00936.X
- 364 Cai H, Lam NSN, Qiang Y, Zou L, Correll RM, Mihunov V (2018) A synthesis of
365 disaster resilience measurement methods and indices. *international journal of*
366 *disaster risk reduction* 31:844-855 doi:10.1016/J.IJDRR.2018.07.015
- 367 Dariya B, Nagaraju GP (2020) Understanding novel COVID-19: Its impact on organ
368 failure and risk assessment for diabetic and cancer patients. *Cytokine & Growth*
369 *Factor Reviews* doi:<https://doi.org/10.1016/j.cytogfr.2020.05.001>
- 370 de León-Martínez LD, de la Sierra-de la Vega L, Palacios-Ramírez A, Rodríguez-
371 Aguilar M, Flores-Ramírez R (2020) Critical review of social, environmental
372 and health risk factors in the Mexican indigenous population and their capacity
373 to respond to the COVID-19. *Science of The Total Environment*:139357
374 doi:<https://doi.org/10.1016/j.scitotenv.2020.139357>
- 375 Desjardins MR, Hohl A, Delmelle EM (2020) Rapid surveillance of COVID-19 in the
376 United States using a prospective space-time scan statistic: Detecting and
377 evaluating emerging clusters. *Applied Geography* 118:102202
378 doi:<https://doi.org/10.1016/j.apgeog.2020.102202>
- 379 Fotheringham AS, Brunsdon C, Charlton M (2002) Geographically Weighted
380 Regression: The Analysis of Spatially Varying Relationships.

- 381 Grömping U (2009) Variable Importance Assessment in Regression: Linear Regression
 382 versus Random Forest. *the american statistician* 63:308-319
 383 doi:[10.1198/TAST.2009.08199](https://doi.org/10.1198/TAST.2009.08199)
- 384 Guliyev H (2020) Determining the spatial effects of COVID-19 using the spatial panel
 385 data model. *Spatial Statistics* 38:100443
 386 doi:<https://doi.org/10.1016/j.spasta.2020.100443>
- 387 Huang C et al. (2020) Clinical features of patients infected with 2019 novel coronavirus
 388 in Wuhan, China. *The Lancet* 395:497-506 doi:[https://doi.org/10.1016/S0140-6736\(20\)30183-5](https://doi.org/10.1016/S0140-6736(20)30183-5)
- 389
- 390 Imran M, Stein A, Zurita-Milla R (2015) Using geographically weighted regression
 391 kriging for crop yield mapping in West Africa. *International Journal of*
 392 *Geographical Information Science* 29:234-257
 393 doi:[10.1080/13658816.2014.959522](https://doi.org/10.1080/13658816.2014.959522)
- 394 Ishwaran H (2007) Variable importance in binary regression trees and forests. *electronic*
 395 *journal of statistics* 1:519-537 doi:[10.1214/07-EJS039](https://doi.org/10.1214/07-EJS039)
- 396 Landi F et al. (2020) THE GERIATRICIAN: THE FRONTLINE SPECIALIST IN THE
 397 TREATMENT OF COVID-19 PATIENTS. *Journal of the American Medical*
 398 *Directors Association* doi:<https://doi.org/10.1016/j.jamda.2020.04.017>
- 399 Lau H et al. (2020) The association between international and domestic air traffic and
 400 the coronavirus (COVID-19) outbreak. *Journal of Microbiology, Immunology*
 401 *and Infection* doi:<https://doi.org/10.1016/j.jmii.2020.03.026>
- 402 Liu F, Li X, Zhu G (2020) Using the contact network model and Metropolis-Hastings
 403 sampling to reconstruct the COVID-19 spread on the “Diamond Princess”.
 404 *Science Bulletin* doi:<https://doi.org/10.1016/j.scib.2020.04.043>
- 405 Lu B, Brunsdon C, Charlton M, Harris P (2017) Geographically weighted regression
 406 with parameter-specific distance metrics. *International Journal of Geographical*
 407 *Information Science* 31:982-998 doi:[10.1080/13658816.2016.1263731](https://doi.org/10.1080/13658816.2016.1263731)
- 408 Malik VS, Ravindra K, Attri SV, Bhadada SK, Singh M (2020) Higher body mass index
 409 is an important risk factor in COVID-19 patients: a systematic review and meta-
 410 analysis. *Environmental Science and Pollution Research* doi:[10.1007/s11356-020-10132-4](https://doi.org/10.1007/s11356-020-10132-4)
- 411
- 412 Marhl M, Grubelnik V, Magdič M, Markovič R (2020) Diabetes and metabolic
 413 syndrome as risk factors for COVID-19. *Diabetes & Metabolic Syndrome:*
 414 *Clinical Research & Reviews* doi:<https://doi.org/10.1016/j.dsx.2020.05.013>
- 415 Millett GA et al. (2020) Assessing Differential Impacts of COVID-19 on Black
 416 Communities. *Annals of Epidemiology*
 417 doi:<https://doi.org/10.1016/j.annepidem.2020.05.003>
- 418 Mollalo A, Vahedi B, Rivera KM (2020) GIS-based spatial modeling of COVID-19
 419 incidence rate in the continental United States. *Science of The Total*
 420 *Environment* 728:138884 doi:<https://doi.org/10.1016/j.scitotenv.2020.138884>
- 421 O'Connor CM, Anoushiravani AA, DiCaprio MR, Healy WL, Iorio R (2020) Economic
 422 Recovery After the COVID-19 Pandemic: Resuming Elective Orthopedic
 423 Surgery and Total Joint Arthroplasty. *The Journal of Arthroplasty*
 424 doi:<https://doi.org/10.1016/j.arth.2020.04.038>
- 425 Rosenkrantz L, Schuurman N, Bell N, Amram O (2020) The need for GIScience in
 426 mapping COVID-19. *Health & Place*:102389
 427 doi:<https://doi.org/10.1016/j.healthplace.2020.102389>
- 428 Ruthberg JS, Quereshey HA, Jella TK, Kocharyan A, D'Anza B, Maronian N, Otteson
 429 TD (2020) Geospatial analysis of COVID-19 and otolaryngologists above age

430 60. American Journal of Otolaryngology:102514
431 doi:<https://doi.org/10.1016/j.amjoto.2020.102514>
432 Serge R, Vandromme J, Charlotte M (2020) Are we equal in adversity? Does Covid-19
433 affect women and men differently? Maturitas
434 doi:<https://doi.org/10.1016/j.maturitas.2020.05.009>
435 Shim E, Tariq A, Choi W, Lee Y, Chowell G (2020) Transmission potential and
436 severity of COVID-19 in South Korea. International Journal of Infectious
437 Diseases 93:339-344 doi:<https://doi.org/10.1016/j.ijid.2020.03.031>
438 Strobl C, Boulesteix A-L, Kneib T, Augustin T, Zeileis A (2008) Conditional Variable
439 Importance for Random Forests. bmc bioinformatics 9:307-307
440 doi:10.1186/1471-2105-9-307
441 Strobl C, Boulesteix A-L, Zeileis A, Hothorn T (2007) Bias in random forest variable
442 importance measures: Illustrations, sources and a solution. bmc bioinformatics
443 8:25-25 doi:10.1186/1471-2105-8-25
444 Targher G et al. (2020) Patients with diabetes are at higher risk for severe illness from
445 COVID-19. Diabetes & Metabolism
446 doi:<https://doi.org/10.1016/j.diabet.2020.05.001>
447 Tosepu R, Gunawan J, Effendy DS, Ahmad LOAI, Lestari H, Bahar H, Asfian P (2020)
448 Correlation between weather and Covid-19 pandemic in Jakarta, Indonesia.
449 Science of The Total Environment 725:138436
450 doi:<https://doi.org/10.1016/j.scitotenv.2020.138436>
451 Wheeler D, Tiefelsdorf M (2005) Multicollinearity and correlation among local
452 regression coefficients in geographically weighted regression. journal of
453 geographical systems 7:161-187 doi:10.1007/S10109-005-0155-6
454 World Health Organization b (2020b) Coronavirus disease (COVID-19) outbreak
455 situation. Accessed 21 March 2020
456 Xu H et al. (2020) Possible environmental effects on the spread of COVID-19 in China.
457 Science of The Total Environment 731:139211
458 doi:<https://doi.org/10.1016/j.scitotenv.2020.139211>
459 Yang W, Elankumaran S, Marr LC (2011) Concentrations and size distributions of
460 airborne influenza A viruses measured indoors at a health centre, a day-care
461 centre and on aeroplanes. journal of the royal society interface 8:1176-1184
462 doi:10.1098/RSIF.2010.0686
463 Zheng R, Xu Y, Wang W, Ning G, Bi Y (2020) Spatial transmission of COVID-19 via
464 public and private transportation in China. Travel Medicine and Infectious
465 Disease 34:101626 doi:<https://doi.org/10.1016/j.tmaid.2020.101626>
466

Figures

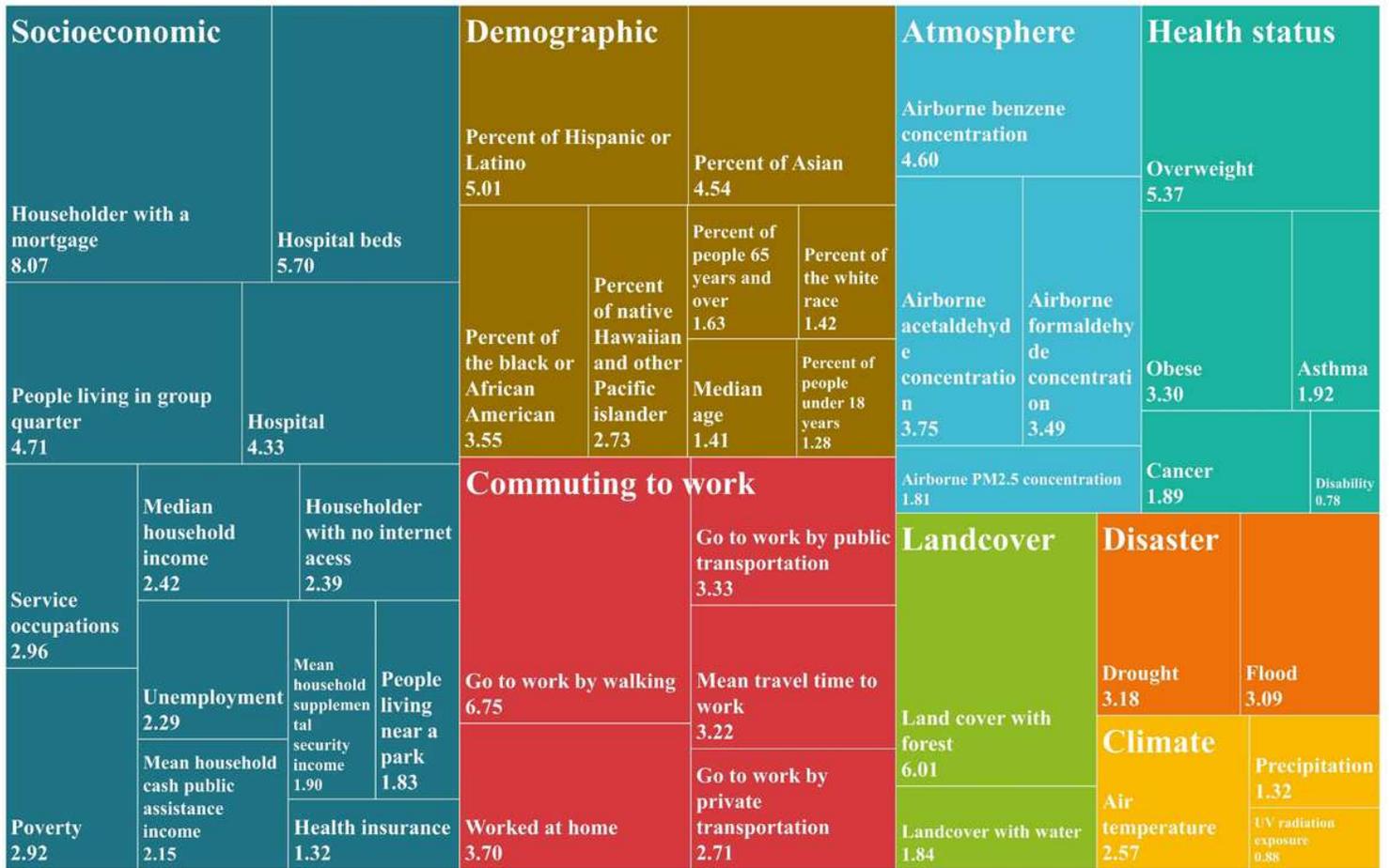


Figure 1

The variable importance of the independent variables of the RF model in modelling COVID-19 death rate

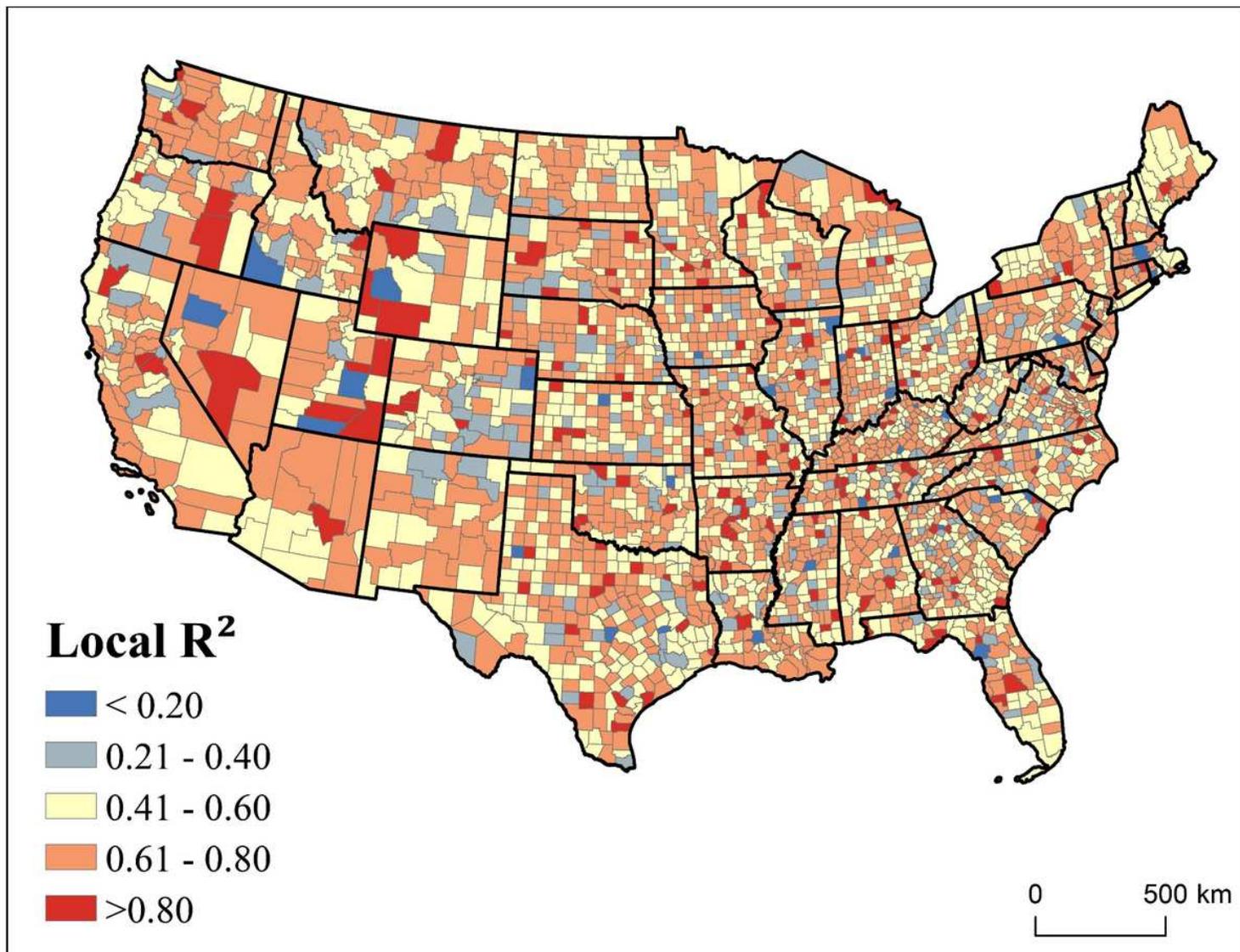


Figure 2

The distribution of local R^2 of the GW-RF

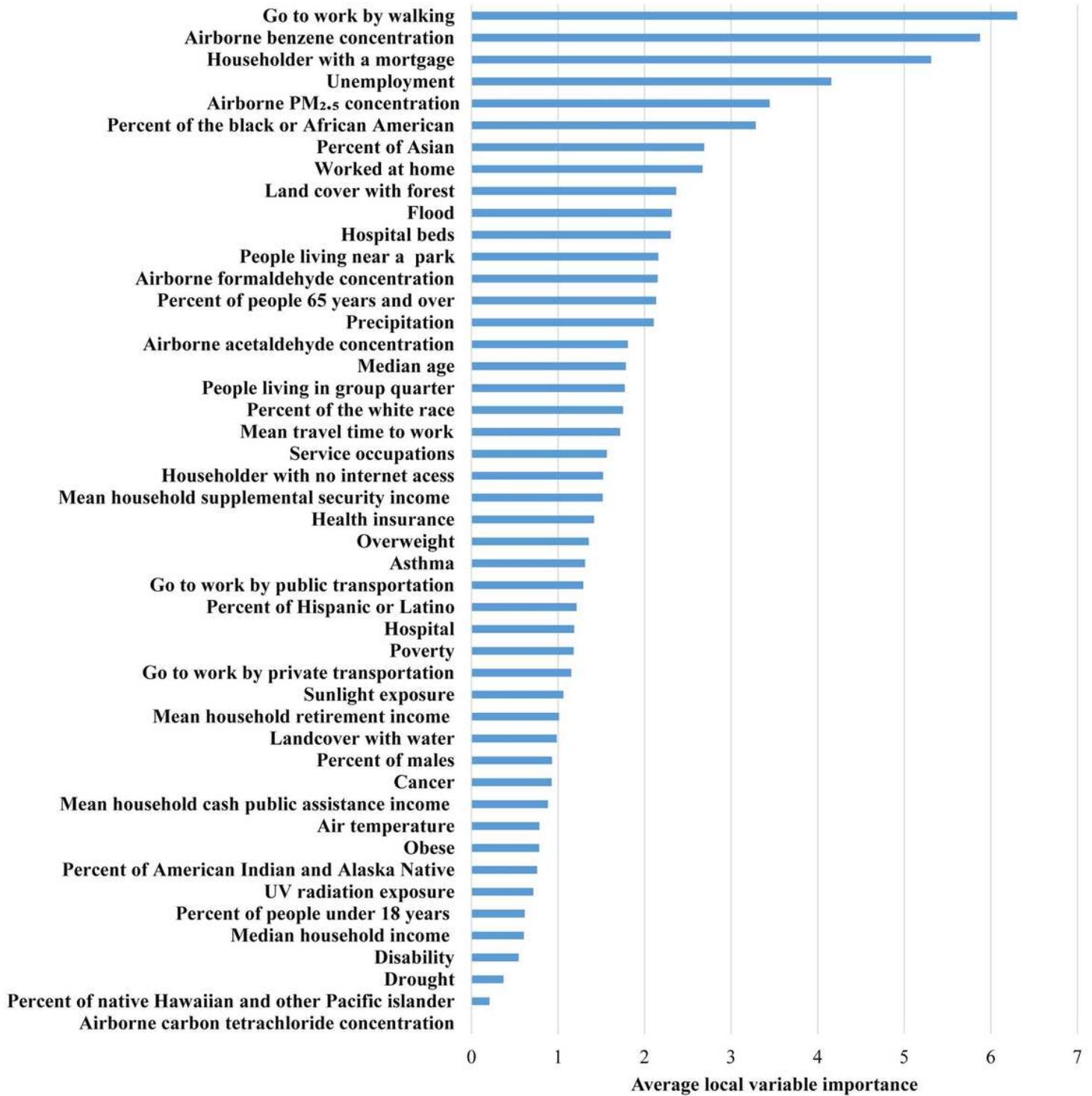
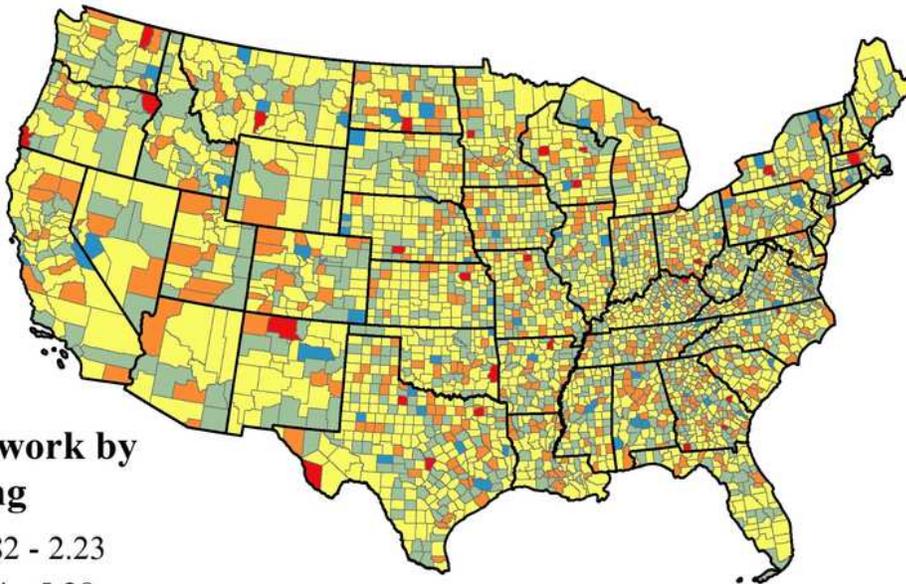


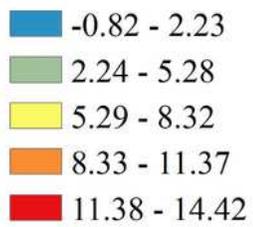
Figure 3

The average local variable importance of 47 potential risk factors on COVID-19 death rate in the GW-RF model

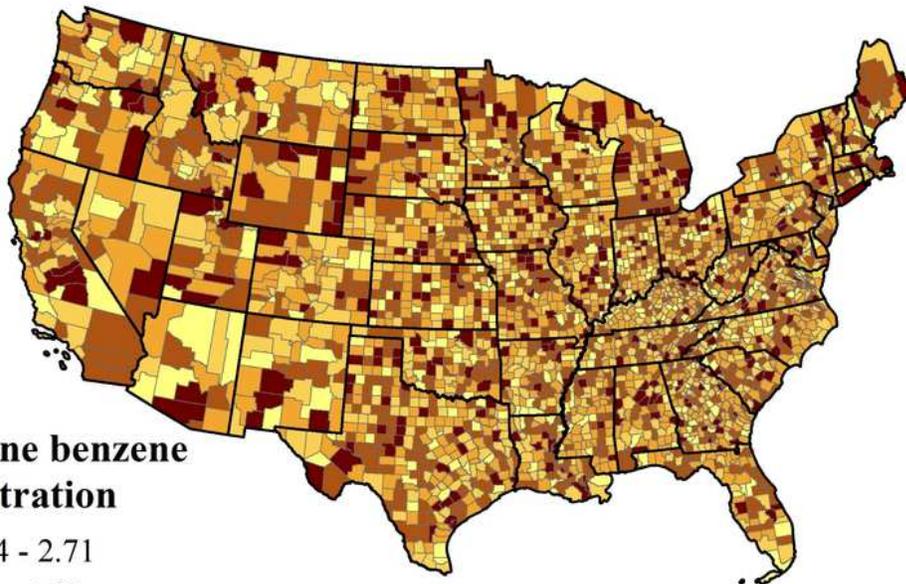
A



Go to work by walking



B



Airborne benzene concentration

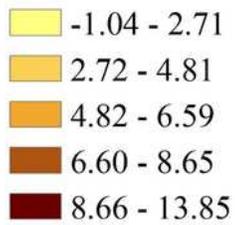
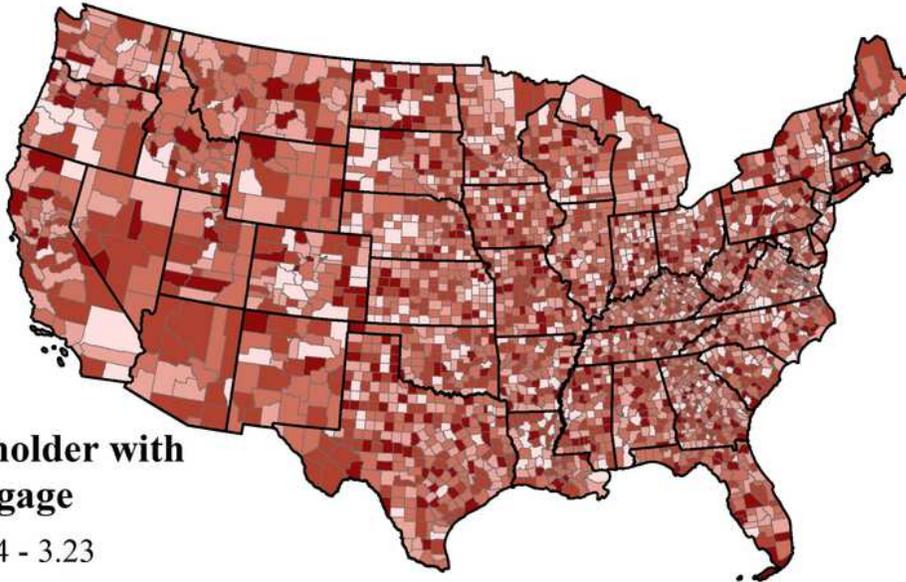


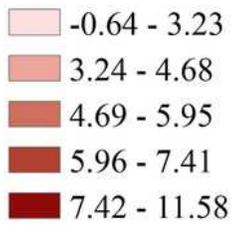
Figure 4

The spatial distribution of the local variable importance of (A) go to work by walking, (B) airborne benzene concentration on COVID-19 death rate in GW-RF mode

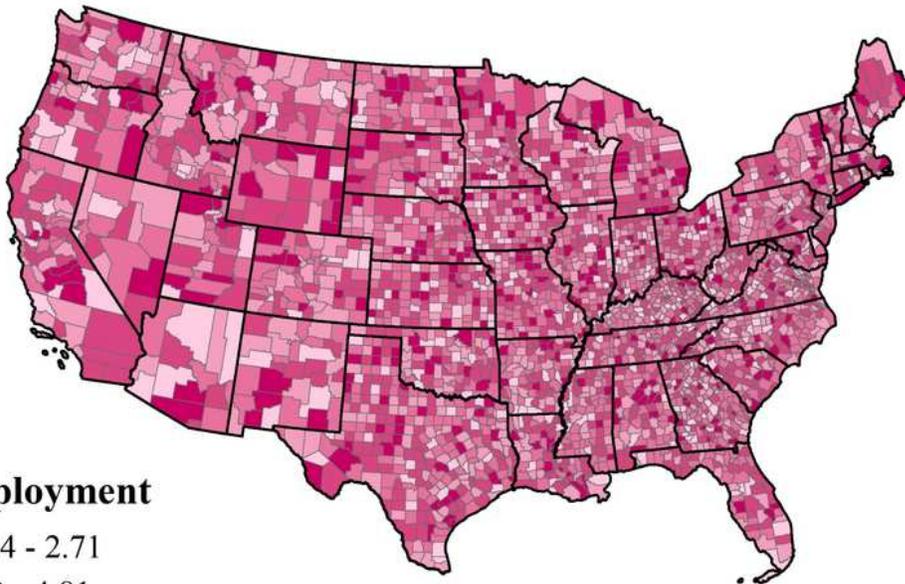
A



Householder with a mortgage



B



Unemployment

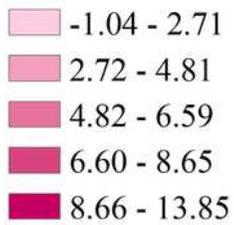
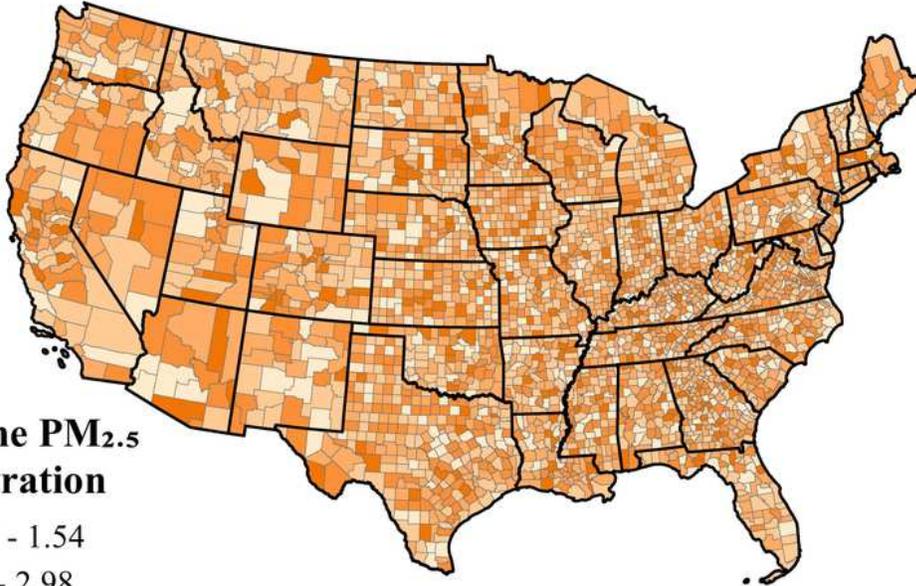


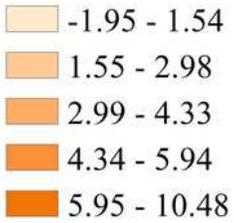
Figure 5

The spatial distribution of the local variable importance of (A) householder with a mortgage, (B) unemployment on COVID-19 death rate in GW-RF model

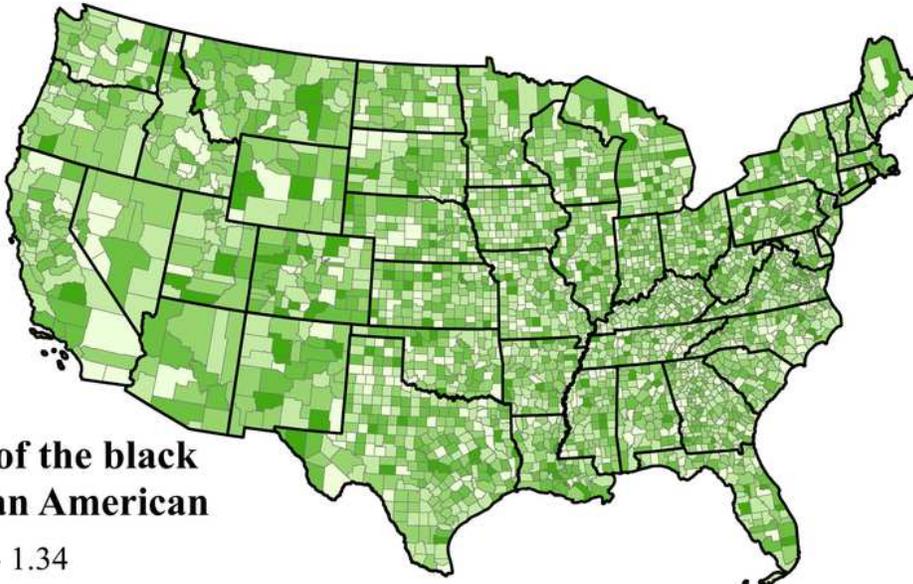
A



**Airborne PM_{2.5}
concentration**



B



**Percent of the black
or African American**

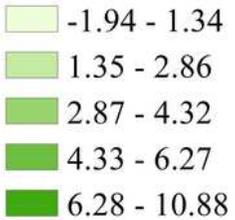


Figure 6

The spatial distribution of the local variable importance of (A) airborne PM 2.5 concentration, (B) percent of the black or African American on COVID-19 death rate in GW-RF model