

1 **Association between COVID-19 mortality and population level health and socioeconomic indicators**

2 Sasikiran Kandula*, Jeffrey Shaman

3 Department of Environmental Health Sciences, Mailman School of Public Health, Columbia University, New York, NY.

4

5

6 *Correspondence:

7 sk3542@cumc.columbia.edu

8 722 West 168th Street, 1104A

9 New York, NY. 10032.

10

11 **Abstract**

12 With the availability of multiple COVID-19 vaccines and the predicted shortages in supply for the near future,
13 it is necessary to allocate vaccines in a manner that minimizes severe outcomes. To date, vaccination strategies
14 in the US have focused on individual characteristics such as age and occupation. In this study, we assess the
15 utility of population-level health and socioeconomic indicators as additional criteria for geographical allocation
16 of vaccines. Using spatial autoregressive models, we demonstrate that 43% of the variability in COVID-19
17 mortality in US counties can be explained by health/socioeconomic factors, adjusting for case rates. Of the
18 indicators considered, prevalence of chronic kidney disease and proportion of population living in nursing
19 homes were found to have the strongest association. In the context of vaccine rollout globally, our findings
20 indicate that national and subnational estimates of burden of disease could be useful for minimizing COVID-19
21 mortality.

22 Introduction

23 By the end of 2020, the COVID-19 pandemic has resulted in 81.5 million documented cases and 1.8 million deaths globally
24 (1). The United States has contributed nearly a quarter of these cases and has lost 1 in every 1000 residents to COVID-19
25 (2). The outbreak has affected all states in the US but with considerable differences in the trajectory and severity of
26 individual outbreaks. Besides this inter- and intra-state geographical variability, the likelihood of adverse outcomes among
27 those infected is reported to be associated with individual's age, gender, race/ethnicity and underlying health conditions
28 (3-6). An estimated 22% of the global population and 28% of the US have one or more of the underlying conditions that
29 pose increased risk of severe outcomes from COVID-19(7).

30 Early studies on clinical characteristics of severe outcomes from COVID-19 were reported from China(5, 8), after the first
31 large outbreak in Wuhan, and concurring estimates were subsequently published from UK, France, US and elsewhere (3, 4,
32 9-12). Guan et al (8) reported that among 1100 of the earliest laboratory confirmed cases of COVID-19 in China, the
33 presence of co-morbidities such as diabetes, hypertension and chronic obstructive pulmonary disease were more prevalent
34 in those with severe outcomes (admission to ICU, requiring mechanical ventilation or death), along with a slightly elevated
35 risk among men and by now well-established risk with increasing age. Using a larger data sample of 45 thousand cases,
36 Deng et al(5) reported that mortality was associated (relative risk (RR) or hazard ratio (HR)) with cardiovascular disease
37 (RR = 6.75, 95%CI = 5.40-8.43), hypertension (HR = 4.48, 95%CI = 3.69-5.45), diabetes (RR = 4.43, 95%CI = 3.49-5.61) and
38 respiratory disease (RR = 3.43, 95%CI = 2.42-4.87, $p < 0.001$).

39 A later, more extensive study(9) from the UK linking 17 million cases to 11 thousand deaths also found association between
40 COVID-19 deaths and kidney disease (HR=2.5, 95%CI = 2.3-2.7), diabetes(HR = 1.95, 95%CI = 1.8-2.1), extreme obesity (HR
41 = 1.9, 95%CI = 1.7-2.1) and several other co-morbidities. From a pooled analysis of 75 studies from multiple countries,
42 Popkin et al(12) summarized that individuals with obesity are at increased risk of death (OR = 1.48; 95% CI, 1.22-1.80)
43 hospitalization (OR = 2.13; 95% CI, 1.74-2.60) and ICU admission (OR = 1.74; 95% CI, 1.46-2.08). Based on these findings
44 and the known prevalence of co-morbidities that existed in the population before the emergence of the pandemic, the
45 populations at risk of severe COVID-19 outcomes at county-level in the US(13) and in several countries have been
46 estimated(7). Other studies have examined the associations of socioeconomic characteristics including poverty, income and
47 race/ethnicity (14-16).

48 Over the past year, public health attempts to reduce transmission largely centered on non-pharmaceutical interventions
49 such as social distancing, face coverings and hand hygiene. In the US, these interventions have had limited success and part
50 of this failure stems from their dependence on collective compliance. The recent availability of high efficacy vaccines gives
51 individuals an additional tool to protect themselves (vaccine supply permitting), and importantly, does not require
52 cooperation from collective public.

53 The availability of vaccines also implies an opportunity to refocus our efforts at reducing infections from efforts at reducing
54 severe outcomes by prioritizing vaccination for those at a higher risk of severe outcomes. To date such strategies have been
55 largely guided by individual characteristics such as age and occupation. We hypothesize that population level
56 characteristics can also guide the optimal allocation and distribution of vaccines geographically. This points to a potential
57 two-layered approach of first identifying high-risk communities within which high-risk individuals can be prioritized.

58 Here, we assess the feasibility of the first part of such an approach and evaluate the extent to which the geographical
59 variability of mortality in US can be explained by population characteristics that predate the epidemic. Our outcome of

60 interest is COVID-19 associated mortality rates at county resolutions, which we attempt to model as a function of population
61 health and socioeconomic indicators. An initial set of indicators associated with COVID-19 mortality as reported in peer-
62 reviewed studies, and data sources for estimates of these indicators were identified. A smaller subset of the variables were
63 selected based on the correlation between the variables and their independent effects on the response.

64 Conventional regression models assume that observations are independent of one another, which in the case of spatial data
65 translates to assuming observations in nearby locations are no more closely related than those farther away. Given the
66 transmission dynamics of COVID-19, counties nearby are likely to have similar case and death rates and spatial
67 dependence rather than spatial independence is a more appropriate assumption. This spatial dependence also extends to
68 health and socioeconomic indicators and potentially latent and unobservable characteristics that effect mortality.

69 Spatial autoregressive (SAR) models offer a parsimonious way to augment basic regression models with spatial dependence
70 between locations (17), and are an extensively studied family of analytical approaches with applications ranging from
71 econometrics, environmental studies and health sciences (18-20). In the current study, we first establish the presence of
72 spatial autocorrelation in the response and explanatory variables, thus motivating the need for spatial models. We apply
73 three forms of SAR models, show that they explain a greater proportion of the variability in mortality than linear models
74 and report effect estimates from each.

75 **Data and Methods**

76 County level indicators of population's health and social status were retrieved from public sources including the US census
77 and large population surveys. In cases where the survey data are not available at county resolutions, data from prior studies
78 on small-area estimates were used. We tried to limit the number of source dependencies and when alternative estimates
79 were available from multiple sources, we preferred estimates from the US Centers for Disease Control and Prevention
80 (CDC). See Table 1 for a list of sources and descriptions for each variable; Figure 1 presents summary statistics.

81 *New York Times*

82 Counts for cumulative cases and deaths through December 31, 2020 were retrieved from New York Times public repository
83 (21). These data included both confirmed and probable cases and deaths at the US county-level and is based on Times'
84 monitoring and analyses of news conferences, data releases and communications with public officials. The determination
85 of cases and deaths as either confirmed or probable is made per definitions laid out in the position statement of the Council
86 of State and Territorial Epidemiologists (22). But as the application can vary across local agencies, here we treat both
87 confirmed and probably categories identically and use total cases and deaths. Case and death rates as a proportion of
88 residents are based on county population estimates from the American Community Survey (ACS) 2014-2018(23).

89 County-specific data for the 5 counties in New York City were retrieved from USAFACTS (24) as the Times' data source was
90 found to combine counts for these five counties into a single entity. Figure 2 shows maps of reported county case and death
91 rates.

92 *Population Level Analysis and Community Estimates (PLACES)*

93 From the PLACES study(25), a collaboration between the CDC and Robert Wood Johnson Foundation, estimates for
94 population-level health and behavioral indicators were retrieved. These small-area estimates of population health
95 outcomes across the US at county resolutions were generated using data collected through the Behavioral Risk Factor

96 Surveillance System (BRFSS)(26), the US decennial 2010 census and the ACS, following a multi-level regression and post-
97 stratification approach (27, 28).

98 Of the 27 indicators available in PLACES, we extracted 5 measures of population level prevalence of health conditions that
99 are reported to have individual level associations with COVID-19 outcomes, namely obesity, diabetes, chronic obstructive
100 pulmonary diseases and chronic heart and kidney diseases. In addition, three related health indicators, the prevalence of
101 high blood pressure and high cholesterol and proportion of residents uninsured were also included.

102 *Social Vulnerability Index*

103 CDC's Social Vulnerability Index (SVI) is a measure of a county's relative vulnerability to hazardous events (29, 30) and is
104 intended to help public officials and planners better prepare for such events. Overall, county ranks are based on fifteen
105 socioeconomic indicators collected in the ACS. Three of the factors in the SVI, namely county population density, median
106 per capita income and proportion of the population that is older than 65 years of age, are hypothesized to be associated
107 with COVID-19 mortality (14, 15, 31, 32). As association between the other variables in SVI and COVID-19 is uncertain, we
108 limited inclusion to the raw estimates of these 3 variables and ignore the other variables in SVI and the overall index.

109 *County Health Rankings*

110 Two additional variables derived from the ACS 2014-18 and available through the County Health Rankings (CHR)(33), are
111 hypothesized to be measures of socioeconomic disparities in a county and included in this study: ratio of the 80th percentile
112 income to 20th percentile as a measure of income inequality, and the proportion of non-White to White residents as a
113 measure of racial diversity. We observed that estimates for these variables in a small percentage (~1.5%) of counties were
114 missing and used the following three-step process to impute missing values: a) the mean of neighboring (defined in later
115 sections) counties that have estimates; b) if there are no neighbors with estimates, the median of all counties in the state
116 for which estimates are available; c) if estimates are missing for all counties in a state, the median across all counties in the
117 US for which estimates are available.

118 *US 2010 Census*

119 It has also been reported that COVID-19 clusters occur in facilities in which people live in group quarters, where the
120 increased vulnerability can result from either the living conditions in such facilities (difficulty to social distance in
121 correctional facilities or on college campuses, for example) or the characteristics of the residents (elderly in nursing homes
122 with underlying health conditions)(34, 35). As mortality from COVID-19 is known to be less likely in younger populations,
123 we focused instead on elderly living in group quarters. An estimate of proportion of the population living in nursing homes
124 or facilities with skilled-nursing in each county was included in this analysis (Table 1).

125 ***Methods overview***

126 We first built linear univariate models for each predictor with county-level COVID-19 mortality as outcome, adjusting for
127 county case rates. These models inform both the individual effects and the proportion of variance in mortality explained by
128 each of these predictors. We followed this with a linear multivariate model, again adjusting for case rates. In both univariate
129 and multivariate models, observational independence is inappropriate because of spatial autocorrelation in both the
130 response and predictors. We verify this by standard tests on the residual of the multivariate model. We finally build spatial
131 simultaneous autoregressive models and report effect estimates.

132 *Spatial weight matrix and spatial autocorrelation*

133 As introduced in an earlier section, a key assumption in standard ordinary least square regression (OLS) models is the
134 independence of observations that does not hold because COVID-19 cases and deaths in a county are related to cases and
135 deaths in other counties (spatial dependence) and often counties adjacent to it (spatial autocorrelation). Models that do not
136 account for spatial dependence and autocorrelation are shown to have inflated type I errors (20, 36).

137 To establish adjacency of counties in the US, we define a simple spatial $n \times n$ matrix, \mathbf{W} , using shape files that list county
138 boundaries as an ordered set of geocoded reference points (37). County adjacency is defined by queen congruity (at least
139 one shared boundary point) and the spatial weight matrix is row standardized i.e. for each county i , the weight of link to
140 county j , w_{ij} , is the inverse of the number of neighbors of i , if j is adjacent to i , and 0 otherwise; $\sum_j w_{ij} = 1$. A county is
141 assumed to not be a neighbor of itself i.e. $w_{ij} = 0$ when $i = j$.

142 Moran's I(38, 39), a commonly used measure of global spatial autocorrelation, is calculated as:

143
$$I = \frac{n * \sum_{i=1}^n \sum_{j=1}^n w_{ij} (x_i - \bar{x}) (x_j - \bar{x})}{(\sum_i \sum_j w_{ij})(\sum_{i=1}^n (x_i - \bar{x})^2)}$$

144 where n is the number of counties, x_i is the variable of interest for county i , \bar{x} is the mean across all counties and w_{ij} is as
145 defined by the spatial weight matrix, \mathbf{W} . Here, as \mathbf{W} is row standardized $\sum_i \sum_j w_{ij} = n$ and the above equation can be
146 simplified. The significance of the statistic was tested under the randomization assumption i.e. x_i are draws from a random
147 distribution and there is no spatial association. A related measure to identify specific regions within the study region that
148 exhibit spatial autocorrelation, the Local Moran's I, was also estimated. Figure 3 shows Moran's I and counties with
149 significant(40) Local Moran's I and for each predictor and outcome.

150 We were also interested in determining whether spatial autocorrelation, if present, resided in the response or in the
151 residual, as this also informs the choice of the spatial model. To identify this we used robust Lagrange Multiplier tests that
152 can detect possible autocorrelated residuals in the presence of an omitted lagged response and vice versa (41, 42). The
153 statistics reported here are from implementations of these tests in the *spdep* (39, 43) R(44) library.

154 *Variable pruning*

155 As the variables selected for inclusion are related, we calculated Spearman's correlation between pairs of variables (Figure
156 4) and found some of the variables to be very highly correlated. Hence, it would not be appropriate to include these pairs
157 together in models. We used the results of the univariate analysis, to aid variable selection by only retaining those variables
158 that have a correlation of less than .75 with variables of a higher R^2 . This led to the elimination of five variables – prevalence
159 indicators for diabetes, heart disease, high blood pressure (all highly collinear with kidney disease), high cholesterol
160 (collinear with COPD) and median per capita income. The linear multivariate model and the spatial models were built using
161 this smaller set of predictors ($n = 9$).

162 *Spatial simultaneous autoregressive models (SAR)*

163 The general form of an autoregressive model in spatial statistics is given by (17, 20, 45):

164
$$y = \mathbf{X}\beta + \rho \mathbf{W}y + \lambda \mathbf{W}u + \varepsilon$$

165 where y is a $n \times 1$ vector of the response variable, \mathbf{X} is a $n \times k$ matrix of k predictors for n counties, \mathbf{W} is the $n \times n$ spatial
166 weight matrix, ρ is the spatial autoregressive lag coefficient and λ the spatial error coefficient and β , u the coefficient and
167 error vectors respectively. When $\lambda = 0$, the autoregressive process is assumed to occur in the response only (captured by
168 $\rho\mathbf{W}$) and the model is referred to as a spatial lag model. When $\rho = 0$, the autoregressive process is assumed to occur only in
169 the errors (captured by $\lambda\mathbf{W}$), and the model referred to as spatial error model. Model implementations are per *spatialreg*
170 (43, 45, 46) library in R.

171 Results

172 Results from the univariate analysis indicate that the selected variables individually explain 24-29% of the variability in
173 mortality, adjusting for case rates. Mortality is estimated to increase by 43 per thousand residents (95% CI: 37-49) for every
174 1% increase in prevalence of chronic kidney disease, and by 10.4 (95% CI: 8-13) for chronic heart disease, 7.4 (95% CI: 6-
175 8) for diabetes, 4.4 (95% CI: 3-5.8) for COPD, 3.7 (95% CI: 2.6-5.8) for high cholesterol, 2.8 (95% CI: 2.2-3.3) for high blood
176 pressure and 2.6 (95% CI: 2-3.2) for obesity prevalence respectively (Figure 5). These health indicators also explain 28%,
177 25.5%, 27.5%, 24.6%, 24.6%, 25.9% and 25.3% of the variability respectively.

178 Among socioeconomic indicators, the largest association was seen with the nursing home variable (Adjusted R^2 : 29%) with
179 an estimated increase of 39 deaths per thousand (95% CI: 34-44) for every 1% increase in percent living in nursing homes.
180 Mortality rates are estimated to increase by 2.8 (95% CI: 2.3-3.4) and 2.4 (95% CI: 2-2.9) for each 1% increase in percentage
181 of the population who are elderly (65+ years) and uninsured 18-64 year olds, respectively. In contrast, mortality rate is
182 estimated to decrease by 1.5 (95% CI: 1.05-1.87) for every thousand dollar increase in per capita income. On average, the
183 R^2 estimates for socioeconomic indicators are lower than for health indicators.

184 Following variable pruning to correct for collinearity, the multivariate model explained 38% of the variability in mortality
185 with a few changes in effect estimates. Obesity's association is not statistically significant in the presence of kidney disease
186 and COPD's association is counterintuitively negative (Table 2).

187 Moran's I test for spatial autocorrelation in residuals of the above model was found to be statistically significant (18.4, $p <$
188 $1E-6$). Both robust LM tests were found to be significant indicating possible autocorrelation in both the error (28.7, $p <$ $1E-$
189 6) and response (33.5, $p <$ $1E-6$). Hence, three model forms, the general SAR model, spatial lag and spatial error models
190 were attempted.

191 The proportion of variability explained by the SAR models is about 14% higher than the linear model (Figure 6). The spatial
192 error model had an Nagelkerke R^2 (47) of 43.5% with an estimated autocorrelation error coefficient (λ) of .418 (95% CI: .37
193 - .46). The spatial lag model and the general model were observed to have an R^2 nearly identical to that of the error model,
194 The autocorrelation coefficient in response (ρ) was found to be .347 (95% CI: .31-.39) for the spatial model, but when both
195 coefficients were estimated simultaneously in a general model, the lag coefficient was found to be not significant: $\lambda = .336$
196 (95% CI: .244-.429); $\rho = .083$ (95% CI: -.007 - .174; $p = .07$). Figure 7 shows the spatial lag model's estimates and residuals.

197 The Global Moran's test on the residuals of all three models found no significant spatial autocorrelation ($p > .05$). Effect
198 estimates for inequality variable were found to be not significant ($p > .05$) in any of the spatial models (Figure 6). The
199 negative association of COPD seen in the linear model is also observed with the spatial models. Obesity, as in the linear
200 model, was found to be not statistically significant in two of the spatial models.

201 To test for sensitivity of models' R^2 to the variable pruning method, we additionally subset variables using alternative
202 spearman correlation thresholds of .5, .65 and .85 and built linear and spatial models with each. Figure 8 shows that R^2 was
203 not sensitive to the value of threshold and the spatial models have a consistently higher R^2 than the linear model.

204 Discussion

205 We have built models to estimate COVID-19 mortality rates for given case rates and population health and socioeconomic
206 characteristics. Our results indicate that together these indicators can explain 43% of the variability in US county mortality
207 rates, when spatial autocorrelation is accounted for. We found that among health indicators considered the prevalence of
208 chronic kidney disease and among socioeconomic indicators the proportion living in nursing homes, have the largest
209 associations with mortality.

210 The choice and timeliness of control strategies in response to an outbreak do affect its progress and caseload. Our findings
211 here show that differential risks of severe outcomes from COVID-19 across populations can be in part estimated from the
212 structures and contexts in which the outbreak occurs, for example, a population's quality of health, its access to healthcare
213 and the disparities therein. With the availability of vaccines, these population level indicators can serve as criteria for
214 prioritizing geographical allocation of vaccines.

215 These findings may also be relevant to low- and middle-income countries (LMIC). It has been reported that almost all of the
216 Pfizer-BioNTech and Moderna vaccine doses to be manufactured through the end of 2021 have been purchased and are
217 reserved for distribution in the US, Canada, UK and the EU (48, 49). Of the 42 countries that have rolled out vaccines by
218 early January 2021, only 6 are middle-income countries and none are low-income countries (50). The COVAX initiative with
219 participation from governments of several LMIC countries, the WHO and partner non-governmental organizations, aims to
220 achieve equitable and affordable access to vaccines globally through a common vaccine purchase and allocation framework
221 (51). When allocation decisions need to span multiple countries, national and subnational socioeconomic indicators and
222 burden of disease estimates can potentially be leveraged to reduce overall risk of severe outcomes from COVID-19 as our
223 findings demonstrate.

224 This study has a few limitations. Case and death counts were retrieved a week after the end of the study period. Given the
225 lags in data reporting, particularly with deaths, events occurring at the end of the study period may not have been recorded
226 and the rates used are underestimates. Similarly, the outcomes may not yet be known for cases recorded near the end of
227 the study period.

228 The adjacency based spatial weight matrix that was used in this study does not sufficiently capture the spread of COVID-
229 19. Cases that occur in a county are not only correlated with those in counties geographically adjacent to it but also with
230 counties with which it has strong population mixing; for example, counties with metropolitan centers into which
231 commuters travel from the suburbs, or counties with major airports. Spatial weight matrices that capture mobility patterns
232 may be more appropriate and lead to better spatial models. Similarly, methods that can explicitly account for spatial
233 autocorrelation in predictors remain to be explored.

234 Finally, the model structure presented may not be parsimonious in the number of predictors. Although we dropped a third
235 of the predictors initially considered (to correct observed collinearity), model forms with a smaller subset of independent
236 variables may yield near identical R^2 and need to be explored. This is also belied by the lack of significance of some of the
237 predictors included in the spatial models. One approach could start with a minimal set of predictors, incrementally add

238 predictors while evaluating goodness of the resulting model in each iteration and terminating when the improvement is
239 below a threshold. Similarly, the variable pruning discussed above is ad hoc; the variables included in the model may be
240 interchangeable with those discarded with only marginal change in model performance.

241 References

- 242 1. WHO Coronavirus Disease (COVID-19) Dashboard 2020 [Available from: <https://covid19.who.int/>.
- 243 2. CDC Case Task Force. United States COVID-19 Cases and Deaths by State over Time 2021 [Available from:
244 <https://data.cdc.gov/Case-Surveillance/United-States-COVID-19-Cases-and-Deaths-by-State-o/9mfq-cb36>.
- 245 3. Chow N, Fleming-Dutra K, Gierke R, Hall A, Hughes M, Pilishvili T, et al. Preliminary estimates of the prevalence of
246 selected underlying health conditions among patients with coronavirus disease 2019—United States, February 12–
247 March 28, 2020. 2020.
- 248 4. CDC. Coronavirus disease 2019 (COVID-19). Evidence used to update the list of underlying medical conditions that
249 increase a person’s risk of severe illness from COVID-19. 2020.
- 250 5. Deng G, Yin M, Chen X, Zeng F. Clinical determinants for fatality of 44,672 patients with COVID-19. *Critical Care*.
251 2020;24(1):1-3.
- 252 6. Stokes EK, Zambrano LD, Anderson KN, Marder EP, Raz KM, Felix SEB, et al. Coronavirus disease 2019 case
253 surveillance—United States, January 22–May 30, 2020. *Morbidity and Mortality Weekly Report*. 2020;69(24):759.
- 254 7. Clark A, Jit M, Warren-Gash C, Guthrie B, Wang HH, Mercer SW, et al. Global, regional, and national estimates of the
255 population at increased risk of severe COVID-19 due to underlying health conditions in 2020: a modelling study. *The*
256 *Lancet Global Health*. 2020;8(8):e1003-e17.
- 257 8. Guan W-j, Ni Z-y, Hu Y, Liang W-h, Ou C-q, He J-x, et al. Clinical characteristics of coronavirus disease 2019 in China.
258 *New England journal of medicine*. 2020;382(18):1708-20.
- 259 9. Williamson EJ, Walker AJ, Bhaskaran K, Bacon S, Bates C, Morton CE, et al. Factors associated with COVID-19-related
260 death using OpenSAFELY. *Nature*. 2020;584(7821):430-6.
- 261 10. Docherty AB, Harrison EM, Green CA, Hardwick H, Pius R, Norman L, et al. Features of 16,749 hospitalised UK patients
262 with COVID-19 using the ISARIC WHO Clinical Characterisation Protocol. *medrxiv*. 2020.
- 263 11. Simonnet A, Chetboun M, Poissy J, Raverdy V, Noulette J, Duhamel A, et al. High prevalence of obesity in severe acute
264 respiratory syndrome coronavirus-2 (SARS-CoV-2) requiring invasive mechanical ventilation. *Obesity*. 2020.
- 265 12. Popkin BM, Du S, Green WD, Beck MA, Algaith T, Herbst CH, et al. Individuals with obesity and COVID-19: A global
266 perspective on the epidemiology and biological relationships. *Obesity Reviews*. 2020;21(11):e13128.
- 267 13. Razzaghi H, Wang Y, Lu H, Marshall KE, Dowling NF, Paz-Bailey G, et al. Estimated county-level prevalence of selected
268 underlying medical conditions associated with increased risk for severe COVID-19 illness—United States, 2018.
269 *Morbidity and Mortality Weekly Report*. 2020;69(29):945.
- 270 14. Adhikari S, Pantaleo NP, Feldman JM, Ogedegbe O, Thorpe L, Troxel AB. Assessment of community-level disparities in
271 coronavirus disease 2019 (COVID-19) infections and deaths in large US metropolitan areas. *JAMA network open*.
272 2020;3(7):e2016938-e.
- 273 15. Baena-Díez JM, Barroso M, Cordeiro-Coelho SI, Díaz JL, Grau M. Impact of COVID-19 outbreak by income: hitting hardest
274 the most deprived. *Journal of Public Health*. 2020;42(4):698-703.
- 275 16. Townsend MJ, Kyle TK, Stanford FC. Outcomes of COVID-19: disparities in obesity and by ethnicity/race. *Nature*
276 *Publishing Group*; 2020.
- 277 17. LeSage JP. An introduction to spatial econometrics. 2008(123):19-44.

- 278 18. Fischer MM, Getis A. Handbook of applied spatial analysis: software tools, methods and applications: Springer Science
279 & Business Media; 2009.
- 280 19. LeSage JP, Pace RK. Spatial econometric models. Handbook of applied spatial analysis: Springer; 2010. p. 355-76.
- 281 20. F. Dormann C, M. McPherson J, B. Araújo M, Bivand R, Bolliger J, Carl G, et al. Methods to account for spatial
282 autocorrelation in the analysis of species distributional data: a review. *Ecography*. 2007;30(5):609-28.
- 283 21. Coronavirus (Covid-19) Data in the United States [Internet]. 2020 [cited 1/11/2021]. Available from:
284 <https://github.com/nytimes/covid-19-data>.
- 285 22. Council of State and Territorial Epidemiologists. Coronavirus Disease 2019 (COVID-19) 2020 Interim Case Definition
286 [Available from: [https://wwwn.cdc.gov/nndss/conditions/coronavirus-disease-2019-covid-19/case-](https://wwwn.cdc.gov/nndss/conditions/coronavirus-disease-2019-covid-19/case-definition/2020/08/05/)
287 [definition/2020/08/05/](https://wwwn.cdc.gov/nndss/conditions/coronavirus-disease-2019-covid-19/case-definition/2020/08/05/)].
- 288 23. U.S. Census Bureau. 2014-2018 American Community Survey 5-year Public Use Microdata Samples 2019 [
- 289 24. USAFACTS. Coronavirus data resource [Available from: [https://usafacts.org/visualizations/coronavirus-covid-19-](https://usafacts.org/visualizations/coronavirus-covid-19-spread-map/)
290 [spread-map/](https://usafacts.org/visualizations/coronavirus-covid-19-spread-map/)].
- 291 25. Centers for Disease Control and Prevention. PLACES Project [Available from: <https://www.cdc.gov/places>].
- 292 26. Remington PL, Smith MY, Williamson DF, Anda RF, Gentry EM, Hogelin GC. Design, characteristics, and usefulness of
293 state-based behavioral risk factor surveillance: 1981-87. *Public health reports*. 1988;103(4):366.
- 294 27. Zhang X, Holt JB, Lu H, Wheaton AG, Ford ES, Greenlund KJ, et al. Multilevel regression and poststratification for small-
295 area estimation of population health outcomes: a case study of chronic obstructive pulmonary disease prevalence using
296 the behavioral risk factor surveillance system. *American journal of epidemiology*. 2014;179(8):1025-33.
- 297 28. Zhang X, Holt JB, Yun S, Lu H, Greenlund KJ, Croft JB. Validation of multilevel regression and poststratification
298 methodology for small area estimation of health indicators from the behavioral risk factor surveillance system.
299 *American journal of epidemiology*. 2015;182(2):127-37.
- 300 29. Flanagan BE, Gregory EW, Hallisey EJ, Heitgerd JL, Lewis B. A social vulnerability index for disaster management.
301 *Journal of homeland security emergency management*. 2011;8(1).
- 302 30. Centers for Disease Control and Prevention. Social Vulnerability Index [Available from:
303 <https://www.atsdr.cdc.gov/placeandhealth/svi/index.html>].
- 304 31. Lamb MR, Kandula S, Shaman J. Differential COVID-19 case positivity in New York City neighborhoods: Socioeconomic
305 factors and mobility. *Influenza & Other Respiratory Viruses*. 2020.
- 306 32. Sannigrahi S, Pilla F, Basu B, Basu AS, Molter A. Examining the association between socio-demographic composition
307 and COVID-19 fatalities in the European region using spatial regression approach. *Sustainable cities society*.
308 2020;62:102418.
- 309 33. Remington PL, Catlin BB, Gennuso KP. The county health rankings: rationale and methods. *Population health metrics*.
310 2015;13(1):11.
- 311 34. Abrams HR, Loomer L, Gandhi A, Grabowski DC. Characteristics of US Nursing Homes with COVID-19 Cases. *Journal of*
312 *the American Geriatrics Society*. 2020.
- 313 35. Barnett ML, Grabowski DC, editors. Nursing homes are ground zero for COVID-19 pandemic. *JAMA Health Forum*;
314 2020: American Medical Association.
- 315 36. Anselin L. Under the hood issues in the specification and interpretation of spatial regression models. *Agricultural*
316 *economics*. 2002;27(3):247-67.
- 317 37. U.S. Census Bureau. Cartographic Boundary Files - Shapefile 2018 [Available from:
318 <https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html>].

- 319 38. Cliff AD, Ord JK. Spatial processes: models & applications: Taylor & Francis; 1981.
- 320 39. Bivand RS, Wong DW. Comparing implementations of global and local indicators of spatial association. *Test*.
321 2018;27(3):716-48.
- 322 40. Anselin L. Local indicators of spatial association—LISA. *Geographical analysis*. 1995;27(2):93-115.
- 323 41. Anselin L. *Spatial econometrics: methods and models*: Springer Science & Business Media; 2013.
- 324 42. Anselin L, Bera AK, Florax R, Yoon MJ. Simple diagnostic tests for spatial dependence. *Regional science urban*
325 *economics*. 1996;26(1):77-104.
- 326 43. Bivand RS, Pebesma EJ, Gómez-Rubio V, Pebesma EJ. *Applied spatial data analysis with R*: Springer; 2008.
- 327 44. Team RC. *R: A language environment for statistical computing*. R Foundation for Statistical Computing: version 3.5. 0.
328 2018.
- 329 45. Bivand R, Piras G, editors. *Comparing implementations of estimation methods for spatial econometrics* 2015: American
330 Statistical Association.
- 331 46. Bivand R, Hauke J, Kossowski T. Computing the Jacobian in Gaussian Spatial Autoregressive Models: An Illustrated
332 Comparison of Available Methods. *Geographical Analysis*. 2013;45(2):150-79.
- 333 47. Nagelkerke NJ. A note on a general definition of the coefficient of determination. *Biometrika*. 1991;78(3):691-2.
- 334 48. LaFraniere S, Thomas K, Weiland N. Trump administration officials passed when Pfizer offered months ago to sell the
335 U.S. more vaccine doses. *The New York Times*. 2020. 12/07/2020.
- 336 49. Mullard A. How COVID vaccines are being divvied up around the world. 2020 11/30/2020.
- 337 50. Director-General's opening remarks at the media briefing on COVID-19 – 8 January 2021 [press release]. 1/8/2021.
- 338 51. COVAX. GAVI Alliance [Available from: <https://www.gavi.org/covax-facility>].

339

Variable	Source	Description; primary source
Deaths	New York Times(21), USAFACTS(24)	Cumulative COVID-19 confirmed and probable deaths through December 31 2020; per thousand residents
Cases	New York Times, USAFACTS	Cumulative COVID-19 confirmed and probable cases through December 31 2020; per 100000 residents
Obesity	PLACES(25)	Proportion of residents 18+ years of age with calculated BMI \geq 30 kg/m ² , based on self-reported weight and height; BRFSS(26)
Diabetes	PLACES	Proportion of residents 18+ years of age who report being told by a doctor/nurse/other health professional that they have type 1 or type 2 diabetes; BRFSS
Chronic kidney disease (CKD)	PLACES	Proportion of residents 18+ years of age who report being told by a doctor/nurse/other health professional that they have kidney disease; BRFSS
Chronic heart disease (CHD)	PLACES	Proportion of residents 18+ years of age who report being told by a doctor/nurse/other health professional that they have angina or coronary heart disease; BRFSS
Chronic obstructive pulmonary disease (COPD)	PLACES	Proportion of residents 18+ years of age who report being told by a doctor/nurse/other health professional that they have COPD, emphysema, or chronic bronchitis; BRFSS
High cholesterol	PLACES	Proportion of residents 18+ years of age who report being told by a doctor/nurse/other health professional that they have high cholesterol ; BRFSS
High blood pressure	PLACES	Proportion of residents 18+ years of age who report being told by a doctor/nurse/other health professional that they have high blood pressure ; BRFSS
Uninsured	PLACES	Proportion of residents 18-64 years of age who report having no health insurance coverage
Population density	SVI(30)	Number of residents per square mile; Census Cartographic Boundary File - U.S. Tracts 2018(37)
Income	SVI	Median per capita income (in 100 thousand US\$); American Community Survey, 2014-2018 (5-year)(23)
Elderly	SVI	Proportion of residents 65+ years of age; American Community Survey, 2014-2018 (5-year)
Group quarters - nursing	US 2010 Census	Proportion of residents living in nursing/skilled-nursing facilities; P042
Inequality	CHR(33)	Ratio of household income at 80 th percentile with income at 20 th percentile; American Community Survey, 2014-2018 (5-year)
Resident diversity	CHR	Proportion of non-White resident to White residents; American Community Survey, 2014-2018 (5-year)

340

341 **Table 1.** Descriptions and sources for variables included in the study.

342

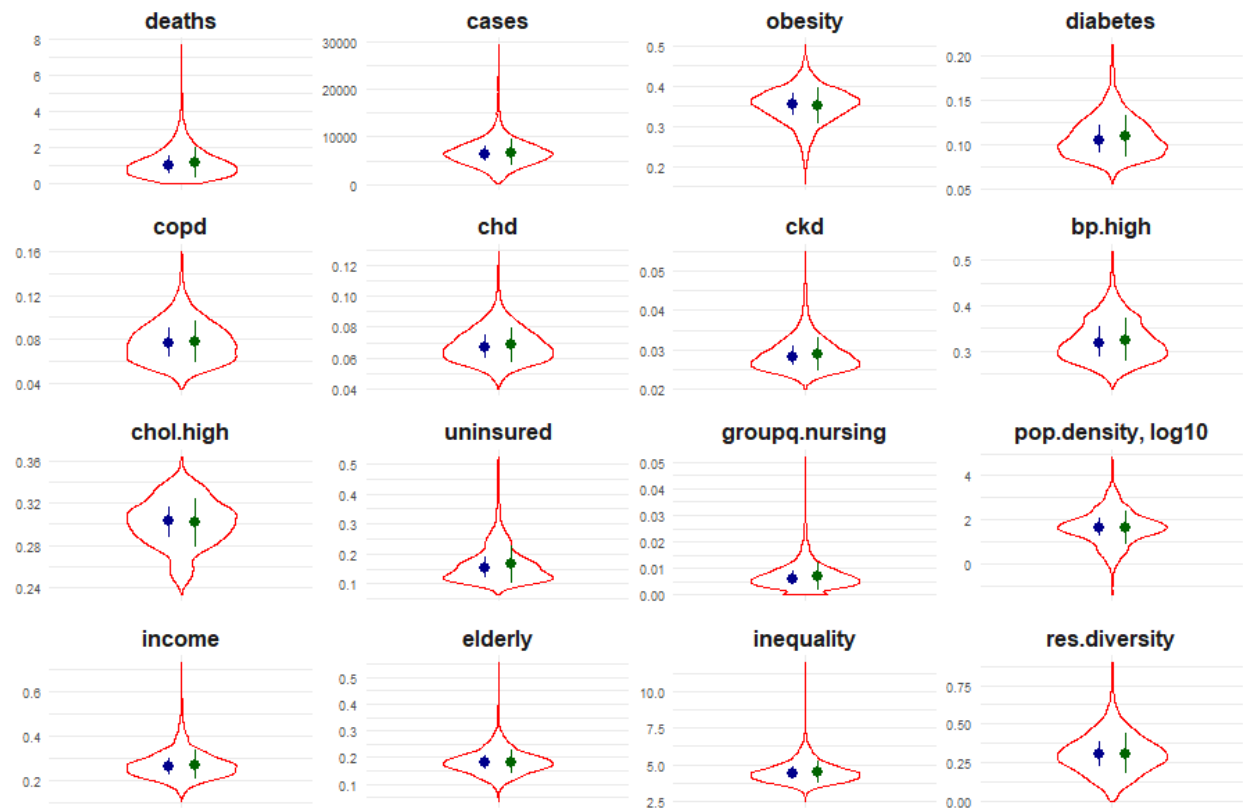
	Estimate	95% CI	p
(Intercept)	-2.007	(-2.27, -1.74)	<0.001
cases	1.26E-4	(1.16E-4, 1.36E-4)	<0.001
obesity	-0.654	(-1.43, .13)	0.101
copd	-4.681	(-6.64, -2.72)	<0.001
ckd	53.48	(41.11, 65.84)	<0.001
groupq.nursing	41.66	(36.27, 47.08)	<0.001
uninsured	1.479	(.92, 2.03)	<0.001
elderly	2.449	(1.87, 3.03)	<0.001
inequality	0.046	(0.005, .087)	0.029
pop.density	3.4E-5	(1.9E-5, 4.8E-5)	<0.001
res.diversity	0.557	(0.36, .75)	<0.001

343

344 **Table 2.** Results of multivariate analysis with linear model, adjusting for case rate. Adjusted R² = .3812; F-
 345 statistic = 194. (p-value: < .001)

346

347



348

349 **Figure 1.** Violin plots of distribution for each variable among counties in the US, along with median (interquartile range) in
350 blue and mean (standard deviation) in green.

351

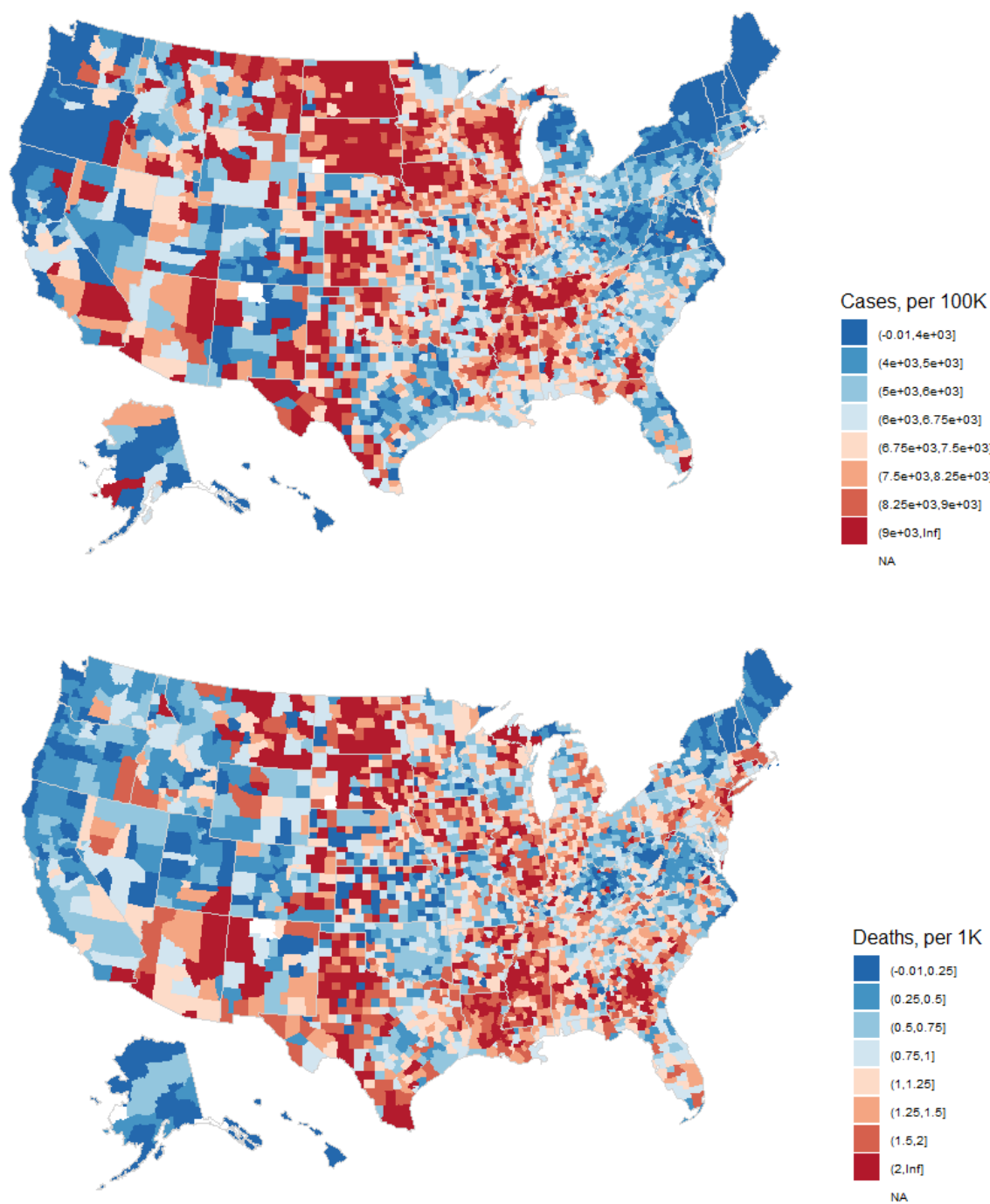
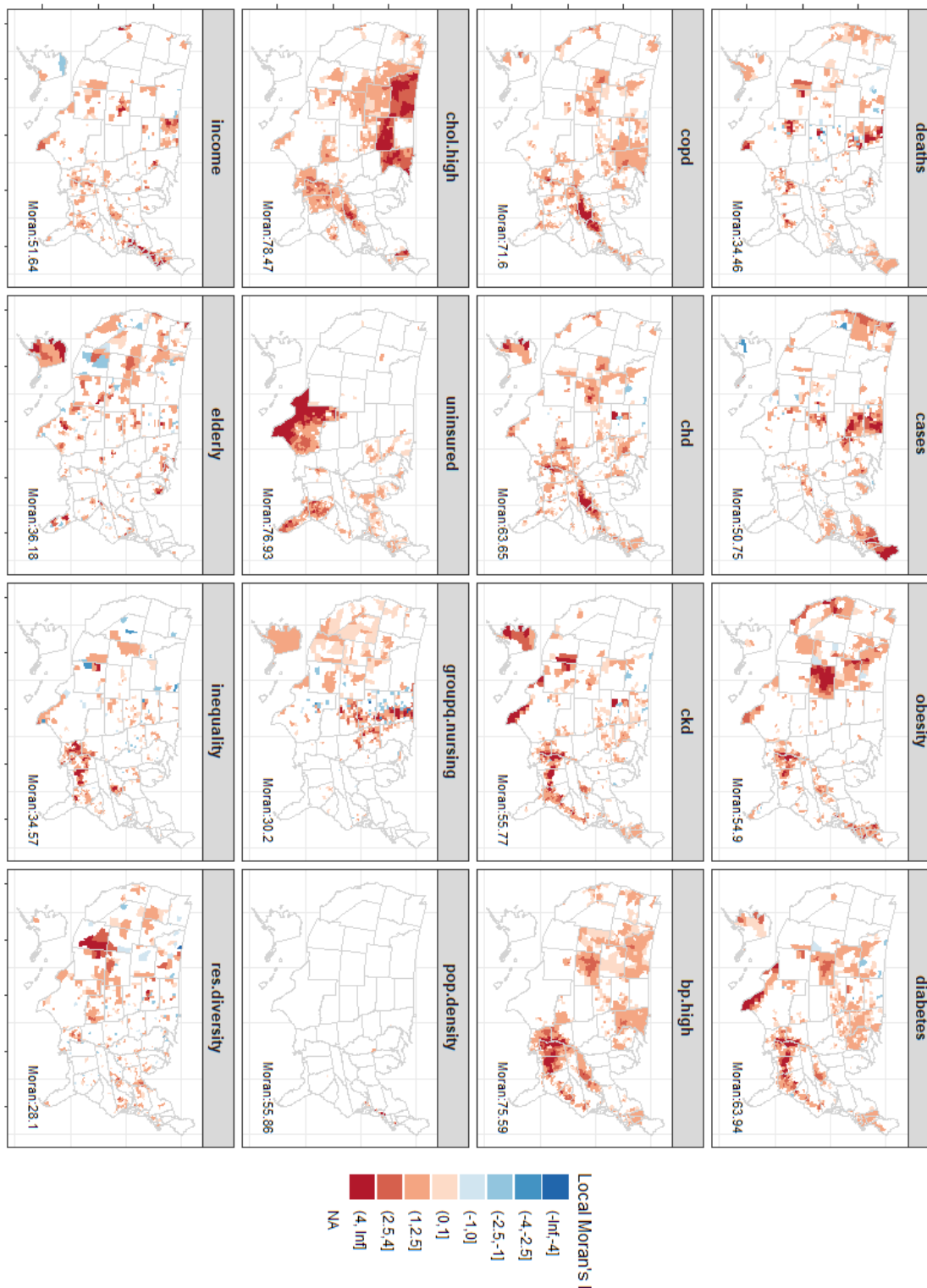


Figure 2. COVID-19 cases (per 100000 residents) and deaths (per 1000 residents) in US counties through December 31, 2020.



356

357 **Figure 3.** Local Moran's I statistic for spatial autocorrelation for all measures and outcome. Only counties where the statistic
 358 is significant ($p < .05$) are shown. Significance is tested under $Pr[1 - E(I)/Var(I)]$ as given by Anselin(40). Global Moran's I
 359 statistic is denoted by the label in each subpanel and found to be statistically significant for all variables.

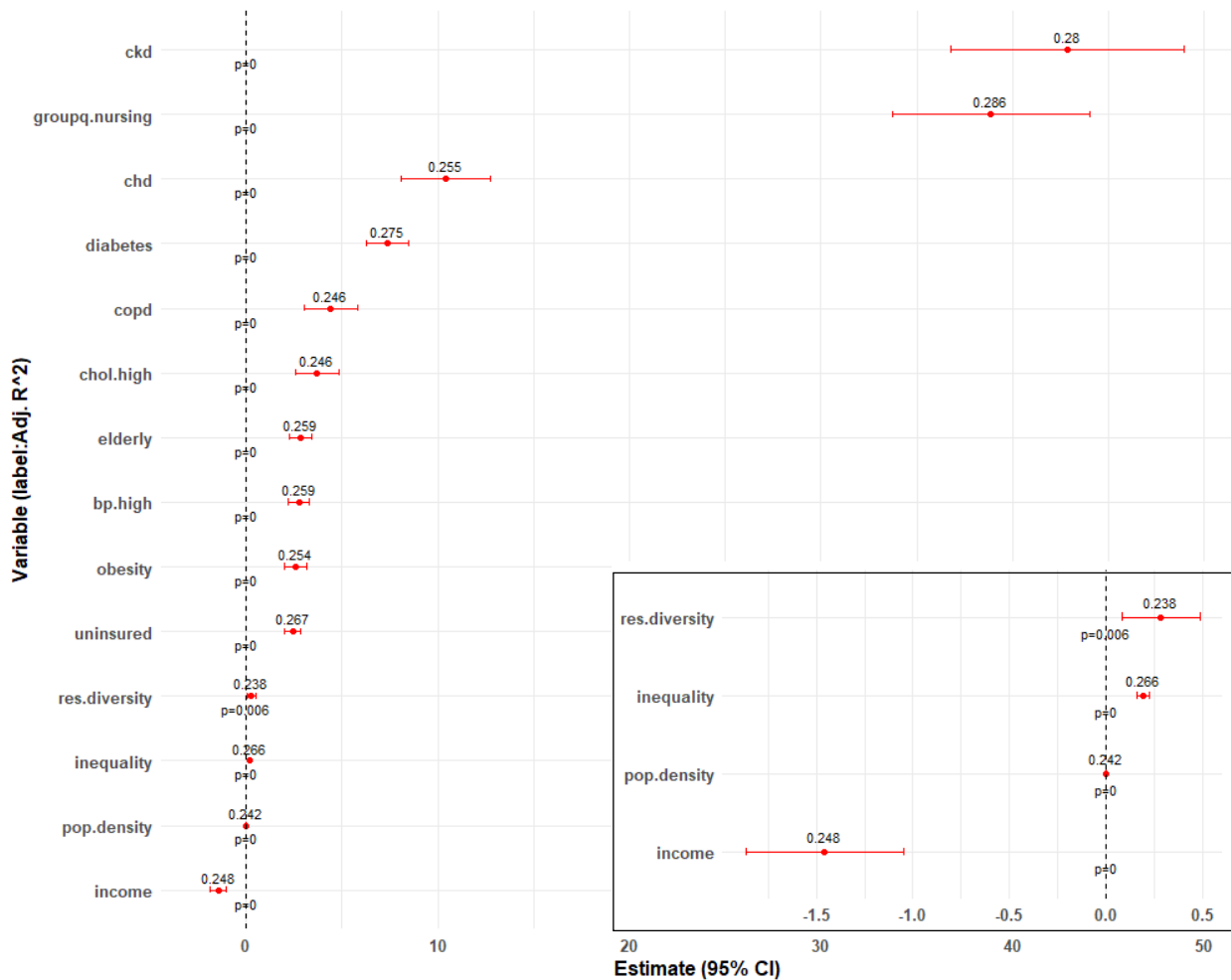
360



361

362 **Figure 4.** Pairwise surface plots (below diagonal), Spearman correlation (above diagonal) and density (diagonal) of
 363 outcome and measures used in the study. * indicates level of statistical significance of the correlation: $p < 0.001$ (***);
 364 $0.001 \leq p < 0.01$ (**); $0.01 \leq p < 0.05$ (*).

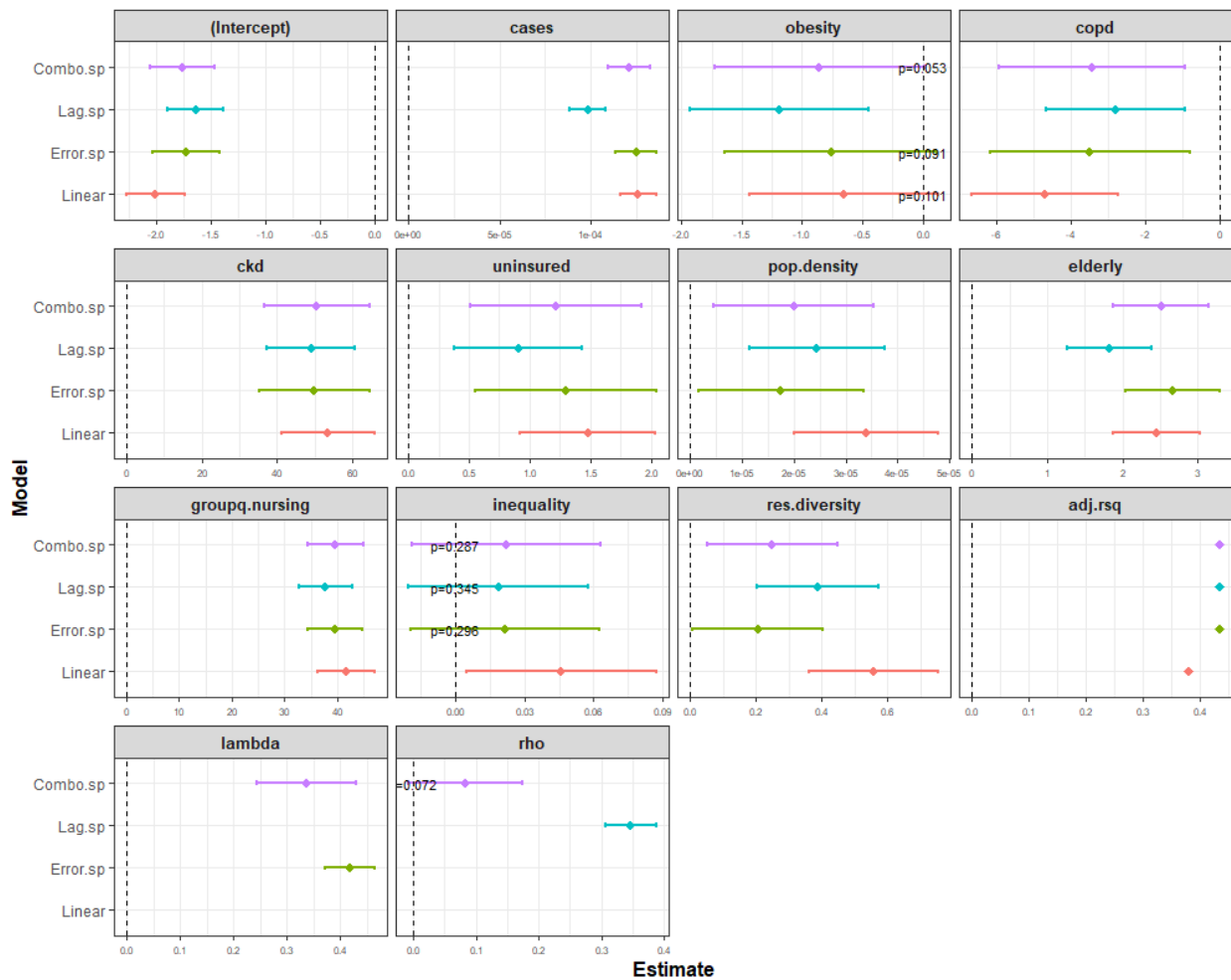
365



366

367 **Figure 5.** Estimates (95% CI) of health and socioeconomic indicators in a linear univariate model with death rate as
 368 outcome and adjusting for COVID-19 case rates. Labels indicate adjusted R². Inset magnifies select variables of smaller
 369 estimates.

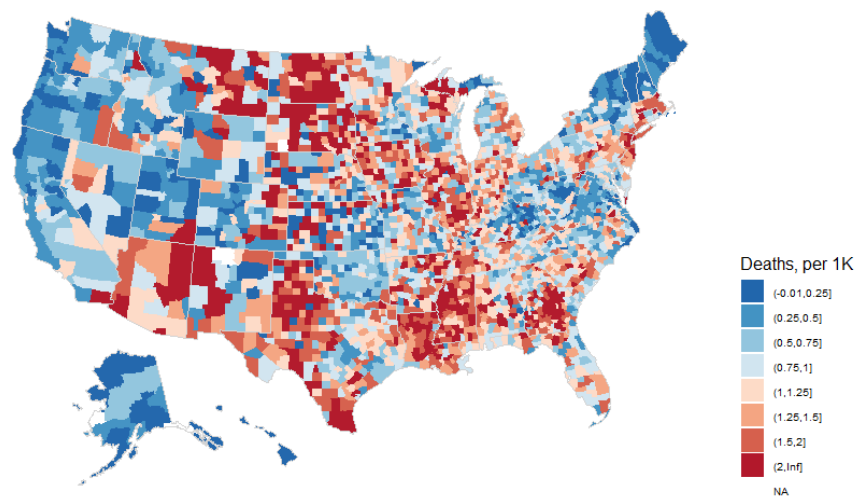
370



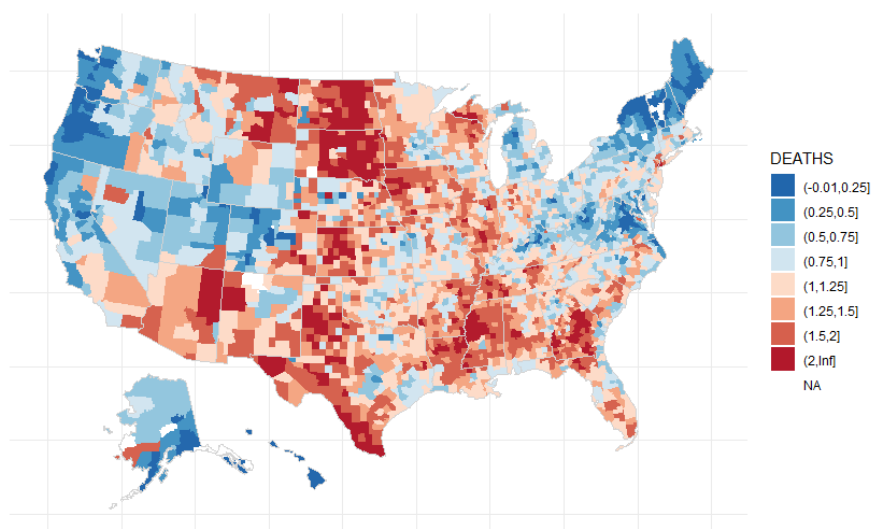
371

372 **Figure 6.** Variables estimates with linear and three spatial regression models. p -values indicated when $p > .05$.

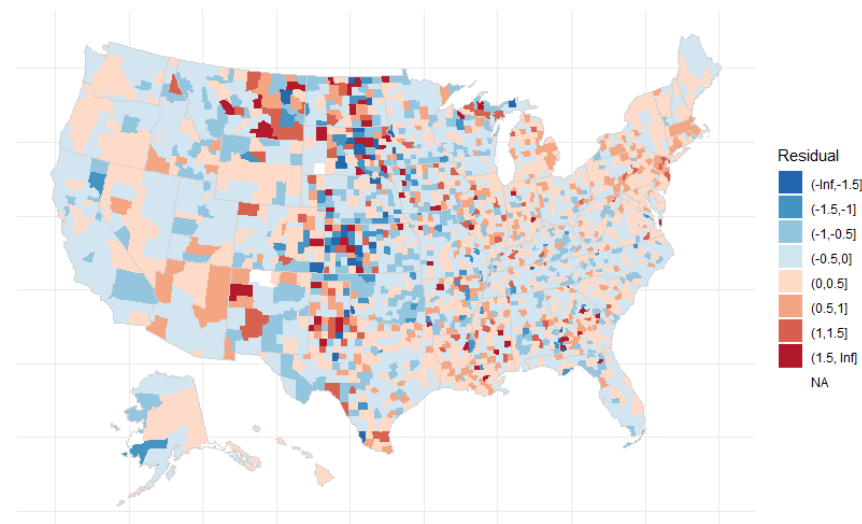
373



374

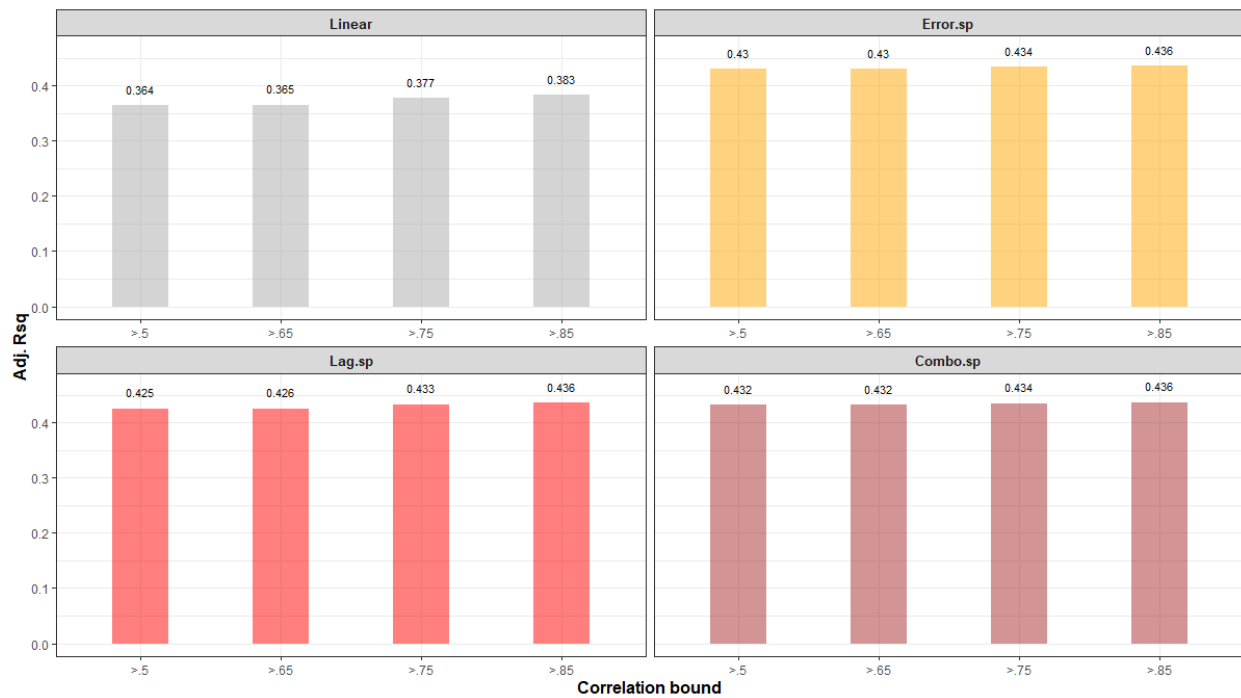


375



376

377 **Figure 7.** Observed death rate (as in Figure 2), model fit and residual of the spatial lag model.



378

379 **Figure 8.** Sensitivity of Adjusted R^2 to Spearman correlation threshold used in variable pruning.