



# Using machine learning to estimate the effect of racial segregation on COVID-19 mortality in the United States

Gerard Torrats-Espinosa<sup>a,b,1</sup>

<sup>a</sup>Department of Sociology, Columbia University, New York, NY 10027; and <sup>b</sup>Data Science Institute, Columbia University, New York, NY 10027

Edited by Douglas S. Massey, Princeton University, Princeton, NJ, and approved December 12, 2020 (received for review July 27, 2020)

This study examines the role that racial residential segregation has played in shaping the spread of COVID-19 in the United States as of September 30, 2020. The analysis focuses on the effects of racial residential segregation on mortality and infection rates for the overall population and on racial and ethnic mortality gaps. To account for potential confounding, I assemble a dataset that includes 50 county-level factors that are potentially related to residential segregation and COVID-19 infection and mortality rates. These factors are grouped into eight categories: demographics, density and potential for public interaction, social capital, health risk factors, capacity of the health care system, air pollution, employment in essential businesses, and political views. I use double-lasso regression, a machine learning method for model selection and inference, to select the most important controls in a statistically principled manner. Counties that are 1 SD above the racial segregation mean have experienced mortality and infection rates that are 8% and 5% higher than the mean. These differences represent an average of four additional deaths and 105 additional infections for each 100,000 residents in the county. The analysis of mortality gaps shows that, in counties that are 1 SD above the Black–White segregation mean, the Black mortality rate is 8% higher than the White mortality rate. Sensitivity analyses show that an unmeasured confounder that would overturn these findings is outside the range of plausible covariates.

COVID-19 | racial segregation | machine learning

COVID-19 caused by severe acute respiratory syndrome coronavirus 2 (SARS CoV-2) infection has led the world to a public health crisis of a scale not seen in a century. Since the first US case was documented in the state of Washington on January 20, 2020, COVID-19 has caused more than 400,000 deaths in the United States and exposed stark racial disparities in disease risk and fatality rates (1). In Chicago and Milwaukee, 70% and 73% of COVID-19 deaths have been among African Americans during the first months of the pandemic. In the states of Louisiana and Michigan, 70% and 40% of deaths have also been among Black residents (2). These rates are more than twice the percentage of Black individuals that make up the population in these areas.

In the United Kingdom, Black individuals are 4.2 times more likely to die from COVID-19 than Whites, and, when adjusting fatality rates for measures of self-reported health and disability, Blacks are still 1.9 times more likely to die from COVID-19 than Whites (3). While death rates by race and ethnicity that adjust for preexisting health conditions have not been reported in the United States, it is plausible to think that an important part of the racial gap in mortality will remain unexplained after individual risk factors and comorbidities are taken into account.

This study examines how racial residential segregation has impacted COVID-19 mortality in 2,174 counties in the United States that include approximately 96% of the country population. The link between racial segregation and health outcomes has

been extensively documented in the social and medical sciences (4, 5). In their systematic review of 39 studies of segregation and health outcomes, Kramer and Hogue (4) report that, for the most part, studies that find an association between Black–White segregation and health outcomes show that segregation is more detrimental for Blacks than it is for Whites (6–9). In five of the studies that they review, Whites are also harmed by the levels of Black–White segregation in their communities (10–14). Other studies have found that the conditions of isolation and population concentration that characterize racially segregated cities have been associated with the spread of infectious diseases such as tuberculosis and HIV/AIDS among racial minorities (15, 16). Racial segregation has also been linked to higher exposure to air pollutants and environmental hazards that affect the respiratory system and increase the risk of mortality among racial minorities (17, 18).

In racially segregated counties, COVID-19 mortality among racial and ethnic minorities may be higher if these groups are more exposed to contextual factors that make them more vulnerable to a highly contagious virus such as SARS-CoV-2. In addition to being more likely to experience preexisting health conditions that increase their vulnerability to COVID-19 (19–22), Blacks and Hispanics are overrepresented in jobs that have been classified as essential during the pandemic (23), are more likely to live in multigenerational households (24), are less likely to have health insurance (25), and live in neighborhoods where essential establishments such as pharmacies and grocery stores are more scarce (26).

Highly clustered friendship networks and social interactions within members of the same racial and ethnic group may increase

## Significance

This study examines the role that racial residential segregation has played during the first 9 mo of the COVID-19 pandemic in the United States. To account for other factors that may explain COVID-19 mortality and infection, I assemble a dataset that includes 50 county-level factors that measure demographics, density, and potential for public interaction, social capital, health risk factors, capacity of the health care system, air pollution, employment in essential businesses, and political views. I use double-lasso regression to guide the selection of the most important controls. Results show that more-segregated counties had higher mortality and infection rates overall and larger mortality rates among Blacks relative to Whites.

G.T.-E. designed research, performed research, analyzed data, and wrote the paper.

The author declares no competing interest.

This article is a PNAS Direct Submission.

Published under the PNAS license.

<sup>1</sup>Email: gerard.torrats@columbia.edu.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2015577118/-DCSupplemental>.

Published February 2, 2021.

the risk of infection among minorities as well (27, 28). If racial minorities are at higher risk of developing COVID-19 because of the higher prevalence of underlying health conditions or because of the environment that surrounds them, more frequent within-race interactions could lead to higher mortality rates and faster transmission rates in places where minorities are more spatially concentrated.

If these mechanisms are at play, the mortality rate among disadvantaged groups like Blacks and Hispanics should be higher than that of Whites in places where Blacks and Hispanics are more segregated from Whites. However, it is also plausible to think that higher levels of segregation could translate into higher mortality rates for the overall population, if the share of minorities in the county is large, and thus drive the overall count of deaths and infections, or if minorities and Whites overlap in public spaces (e.g., public transit and restaurants) so that the virus spills over from minority clusters to the rest of the population through these encounters.

To account for factors that may confound the relationship between residential segregation and COVID-19 mortality, I assemble a dataset that includes 50 county-level factors that are potentially related to residential segregation and COVID-19 infection and mortality. These factors are grouped into eight categories: demographics, density and potential for public interaction, social capital, health risk factors, capacity of the health care system, air pollution, employment in essential businesses, and political views. I use double-lasso regression (29, 30), a machine learning method for model selection and inference, to select the most important predictors in these categories in a statistically principled manner. The double-lasso estimates show that counties that are 1 SD above the segregation mean have experienced a mortality rate that is 8% higher, an infection rate that is 5% higher, and a gap in Black–White mortality that is 8% larger. Sensitivity analyses show that an unmeasured confounder that would overturn these findings is outside the range of plausible explanations.

### Racial Segregation in the COVID-19 Context

Racial residential segregation has been linked to a wide range of public health outcomes. Higher Black–White segregation is associated with elevated mortality among Black individuals in multiple age groups (7, 31). Deaths by homicide among Blacks are higher in racially segregated metropolitan areas (32–34). Beyond mortality and life expectancy outcomes, higher racial segregation is related to higher incidence of tuberculosis and cardiovascular disease among racial minorities (15, 35), lower availability of food establishments serving healthy foods (36), and higher exposure to toxic air pollutants (17, 18).

A large literature on racial and ethnic homophily in social interactions has documented that individuals are more likely to form relationships with others of the same racial or ethnic group (27, 28, 37–40). In a highly segregated setting, these patterns of intragroup interactions may be more frequent, since physical distance between groups makes intergroup connections more difficult. If racial and ethnic minorities are at higher risk of developing COVID-19, the propinquity mechanism may increase rates of infection and mortality in counties where at-risk minorities are more segregated and isolated. These more frequent interactions within minorities that are more vulnerable to COVID-19 may also increase their mortality rates relative to Whites.

Racial and ethnic differences in COVID-19 mortality have emerged in different settings. In US cities and states that report data by race and ethnicity, Blacks and Hispanics are overrepresented in the death and infection counts (2, 41). A fraction of the observed racial and ethnic gaps in infection and mortality rates is likely due to underlying health conditions that make Blacks and Hispanics more vulnerable to COVID-19. A study with 5,700

patients hospitalized with COVID-19 in New York City found that hypertension, obesity, and diabetes were the most common comorbidities (42). Although this study did not report the race and ethnicity of the patients, these health conditions are more prevalent among Blacks and Hispanics (19–22).

Beyond preexisting medical conditions, racial and ethnic minorities may be at higher risk of developing COVID-19 if they experience contextual factors that make them more vulnerable to the virus or make them less likely to get medical care. Blacks and Hispanics are overrepresented in jobs that have been classified as essential during the pandemic (23), are more likely to live in multigenerational households (24), are less likely to have health insurance (25), and live in neighborhoods where essential establishments such as pharmacies and grocery stores are more scarce (26). Furthermore, racial and ethnic minorities may be reluctant to seek medical care or get tested if they fear interacting with the public health system because they are undocumented or an arrest warrant has been issued on them (43, 44).

The key idea is that differences in the likelihood to develop COVID-19 across racial and ethnic groups combined with racial and ethnic homophily in social interactions may produce mortality rates that look very different depending on how the population is spatially distributed. The data collection and empirical strategy described below are designed to test these predictions.

For this study, I measure racial residential segregation using the Relative Diversity Index, a segregation metric capturing the ratio of within-tract diversity to total diversity in the county (*SI Appendix, section 2* explains how the Relative Diversity Index is calculated from Census tract data). The Relative Diversity Index can be interpreted as one minus the ratio of the probability that two individuals from the same tract are members of different racial/ethnic groups to the probability that any two individuals are members of different groups (45). The index can take values from zero to one, with zero indicating that all tracts have the exact same diversity as the county as a whole (i.e., the shares of each racial group are the same across all tracts in the county), and one representing a county where tracts have no diversity (e.g., one set of tracts includes all Black residents and no one else, another set of tracts includes all Hispanic residents and no one else, and so on). I focus on three variants of the Relative Diversity Index: a multigroup version that characterizes overall segregation in the county when all racial/ethnic groups are considered, a Black–White index that captures the extent to which Blacks are segregated from Whites, and a Hispanic–White index that captures the extent to which Hispanics are segregated from Whites. *SI Appendix, Fig. S1* reports bivariate associations between the 50 county attributes and the multigroup Relative Diversity Index. *SI Appendix, Table S1* reports the mean, standard deviation, minimum, median, and maximum values for the three versions of the index. *SI Appendix, section 3* reports results using the Theil's Information Theory Index as the measure of segregation, which lead to the same conclusions as the results shown here (the multigroup Relative Diversity Index and the multigroup Theil Index have a correlation of 0.96 and an  $R^2$  of 0.94, as shown in *SI Appendix, Fig. S2*).

## Results

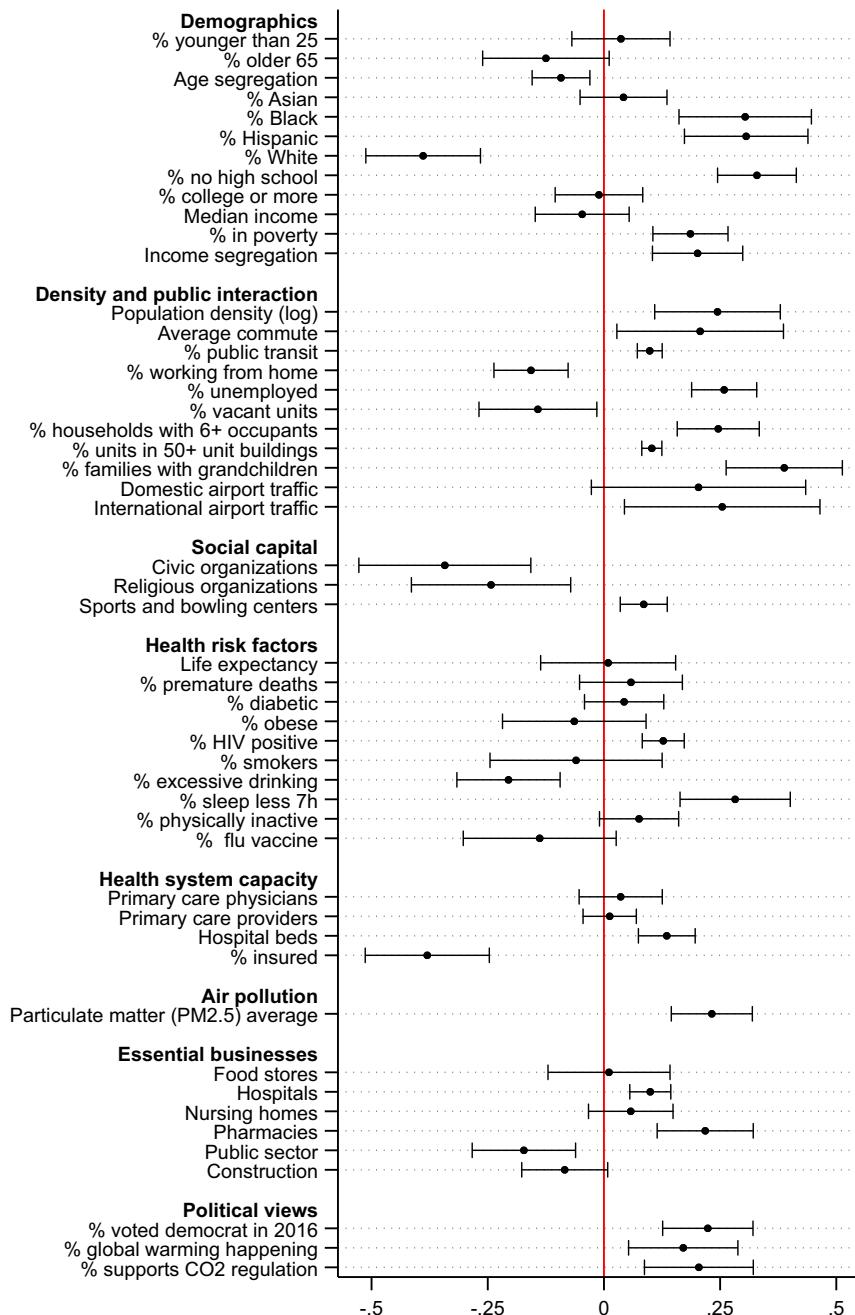
**Correlates of COVID-19.** To account for other county characteristics that may explain higher mortality and infection rates in racially segregated counties, I bring together data from 12 different sources to measure 50 attributes of the counties and their population. I group these 50 county characteristics into the following eight categories: demographics, density and potential for public interaction, social capital, health risk factors, capacity of the health care system, air pollution, employment in essential businesses, and political views. The choice of these eight categories is informed by recent evidence on what we know from the impact of COVID-19 across the US population (46), by prior

studies on the social determinants of health (47, 48), and by sociological evidence on urban inequality, segregation, and health outcomes (49, 50).

Fig. 1 reports standardized bivariate associations between each of the 50 covariates and COVID-19 mortality, net of state fixed effects. *SI Appendix, Table S1* reports the mean, standard deviation, minimum, median, and maximum values for these covariates. *SI Appendix, Table S2* lists all data sources and explains how variables are constructed and when they are measured (all covariates are measured at some point between 2016 and 2020, based on their most recent availability). The analytic sample includes 2,174 counties for which data on all these characteris-

tics are available. This set of counties included 96% of the US population in 2018.

For the most part, the bivariate associations in Fig. 1 reflect the higher incidence of COVID-19 in urban counties. COVID-19 mortality has been higher in counties with larger shares of racial and ethnic minorities, higher poverty rates, higher income segregation, more population density, higher share of the population living in overcrowded housing units, more traffic of passengers in nearby airports, more air pollution, and more progressive political views. These county attributes also correlate strongly with racial segregation (*SI Appendix, Fig. S1*). The empirical strategy described in *Materials and Methods* provides a



**Fig. 1.** Standardized bivariate associations between county attributes and COVID-19 mortality. The death rate is the log of the number of deaths per 100,000 residents in the county as of September 30, 2020. The associations are estimated via OLS regression with state fixed effects, population weights, and standard errors clustered by state. Bars around estimated bivariate associations reflect 95% CIs. The outcome and covariates have been standardized to have mean 0 and SD 1.

statistically principled way to choose the most important controls. It will also enable assessing the magnitude of an unmeasured confounder that would change the estimated relationship between segregation and COVID-19 mortality rates.

**Racial Segregation and Aggregated COVID-19 Mortality and Infection.** Ignoring possible confounding factors, the relationship between multigroup racial segregation and COVID-19 mortality and infection is strong. Fig. 2 plots the bivariate standardized association between the multigroup Relative Diversity Index of segregation and the log rates of COVID-19 deaths and infections per 100,000 residents, net of state fixed effects. A 1-SD difference in the multigroup Relative Diversity Index of segregation predicts a 28% difference in mortality rates and a 17% difference in infection rates.

Fig. 3 shows ordinary least squares (OLS) and double-lasso estimates of the association between multigroup segregation and COVID-19 mortality and infection as of September 30, 2020. For OLS regressions, I show estimates when no controls are included and estimates from models that include state fixed effects. The double-lasso regressions include state fixed effects and the set of 18 controls selected by the lasso approach described in *Materials and Methods*. *SI Appendix, Table S3* shows the full regression output for the double-lasso models.

On the basis of the double-lasso estimates in Fig. 3, a 1-SD difference in segregation predicts an overall mortality rate that is 8% higher (four additional deaths for each 100,000 residents in the county) and an overall infection rate that is 5% higher (105 additional infections for each 100,000 residents in the county).

#### Racial Segregation and Racial and Ethnic Gaps in COVID-19 Mortality.

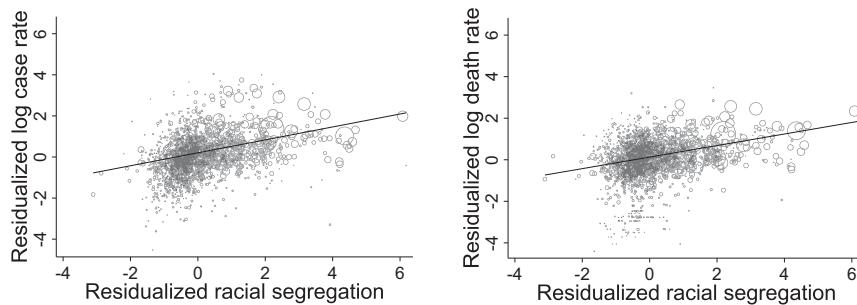
The next set of analyses focuses on the impact of the Black–White Relative Diversity Index and the Hispanic–White Relative Diversity Index on Black–White and Hispanic–White mortality gaps, respectively, as of September 30, 2020. Using two subsamples of counties that reported COVID-19 deaths for Blacks ( $n = 243$ ) and Hispanics ( $n = 218$ ) along with White deaths, I estimate OLS and double-lasso models analogous to those in Fig. 3 in which the outcomes are the difference in the death rate between Blacks and Whites and between Hispanics and Whites. Fig. 4 shows OLS and double-lasso estimates of the association between differences in Black–White and Hispanic–White segregation and differences in Black–White and Hispanic–White mortality gaps. The lasso estimates show that, in counties that are 1 SD above the Black–White segregation mean, the mortality rate among Blacks is 8% higher than the White mortality rate. There is no association between Hispanic–White segregation and the Hispanic–White mortality gap. The full regression output from the lasso models is shown in *SI Appendix, Table S4*.

**Robustness Checks and Sensitivity to Unmeasured Confounding.** *SI Appendix, Tables S3 and S4* show regression results when the racial and ethnic composition of the county is fully accounted for (note that the lasso only selects percent Hispanic and percent White as controls) and when racial and ethnic differences in poverty rates, household median income, unemployment, and life expectancy are included in the regressions. As columns 2 and 5 show, results are unchanged when percent Asian and percent Black are added to the set of controls. Columns 3 and 6 show that controlling for racial and ethnic differences in socioeconomic status and life expectancy does not change the conclusions of the main findings either. *SI Appendix, Figs. S4–S6* are the counterparts to Figs. 2–4 when the Theil’s Information Theory Index is used to measure of segregation. These models yield very similar estimates for the association between segregation and overall mortality, infection rates, and racial/ethnic mortality gaps.

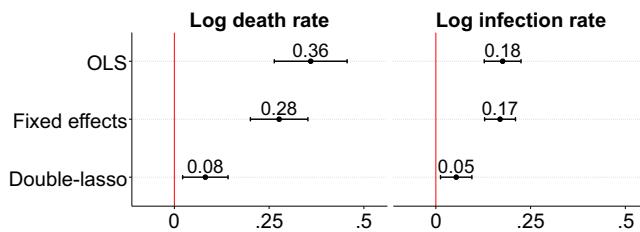
To assess the extent to which the estimate from the Black–White model in Fig. 4 is biased and driven by an unobserved confounder, I implement two sensitivity tests proposed by Frank (51) and Oster (52). *SI Appendix, Fig. S7* shows that an omitted variable that would invalidate the inference in the Black–White model (i.e., make the estimate statistically nonsignificant at the 5% level) would have to be more strongly correlated with the Black–White mortality gap and the Black–White Relative Diversity Index than any of the 50 covariates listed in Fig. 1. Similarly, *SI Appendix, Fig. S8* shows that, for the true effect of Black–White segregation on the Black–White mortality gap to be zero (i.e., the point estimate is zero), there should exist an unobserved covariate that, when added to the regression, increases the  $R^2$  from 0.56 (as shown in *SI Appendix, Table S4*) to 0.80 and is 6 times more predictive of segregation than the 18 controls selected by the lasso. These two tests point to highly implausible scenarios where the findings from Fig. 4 would be overturned.

#### Discussion

The history of the United States is filled with instances where natural disasters and public health crises have disproportionately upended the lives of racial and ethnic minorities and their communities (53). COVID-19 has been a continuation of this pattern. In places that report mortality rates by race and ethnicity, Blacks and Hispanics are overrepresented among the infected, hospitalized, and death due to COVID-19. This study builds on the literature on the social determinants of health (47, 48), to examine how racial residential segregation has impacted mortality and infection rates across US counties. To account for potential confounders of the relationship between segregation and COVID-19 outcomes, I have assembled a dataset with 50 covariates and used double-lasso regression, a machine learning



**Fig. 2.** Relationship between multigroup Relative Diversity Index and COVID-19 infection rates (Left) and mortality (Right) net of state fixed effects. The x axis represents the residuals from a regression of the multigroup Relative Diversity Index on the set of state dummies. The y axis represents the residuals from a regression of the corresponding COVID-19 outcome on the set of state dummies. The size of the dots is proportional to the county population. COVID-19 outcomes are measured as of September 30, 2020.



**Fig. 3.** OLS and double-lasso regression estimates of the relationship between multigroup Relative Diversity Index and COVID-19 death and infection rates. OLS models include no controls. Fixed effects models include state fixed effects. Double-lasso models include the 18 controls selected by the lasso procedure (shown in *SI Appendix, Table S3*) and state fixed effects. All regressions include population weights. Standard errors are clustered by state. Bars around estimated coefficients reflect 95% CIs. The multigroup Relative Diversity Index and covariates have been standardized to have mean 0 and SD 1. Outcomes are in log rates. The sample includes 2,174 counties. COVID-19 outcomes are measured as of September 30, 2020.

method for model selection and inference, to choose the most important controls without making strong a priori assumptions about the functional form.

The findings point to racial segregation as an important driver of mortality and infection rates across the country. A 1-SD difference in racial segregation (measured with the multigroup Relative Diversity Index) predicts an overall mortality rate that is 8% higher and an overall infection rate that is 5% higher. These estimates represent an average of four additional deaths and 105 additional infections for each 100,000 residents in the county. The analyses of mortality by race and ethnicity show that, in counties that are 1 SD above the Black–White segregation mean, the mortality rate among Blacks is 8% higher than the White mortality rate. Hispanics don't exhibit higher COVID-19 mortality rates than Whites in counties where they are more segregated, a finding that is in line with extensive epidemiological evidence on the so-called "Hispanic paradox" (54). Although the research design does not allow for a direct causal interpretation of these findings, sensitivity analyses reveal that an unobserved confounder that would overturn these findings is highly implausible.

The aggregated nature of the data does not allow for a test of mechanisms linking segregation and COVID-19 mortality. In the Introduction, I have speculated that racial homophily in social ties may facilitate the spread of a highly contagious virus such as SARS CoV-2 if racial minorities are segregated and at higher risk of being infected. Results from models that examine mortality gaps between minorities and Whites are consistent with this hypothesis, but future studies using individual-level data that track interactions between individuals at a more granular level should extend this analysis.

Due to data limitations, the estimation of the effect of segregation on racial/ethnic gaps in mortality is not yet possible for all counties. When more data reporting deaths by race become available, extending the analyses presented here should be a priority. Looking at where COVID-19 has had its greatest impact and considering the extensive research that documents the negative effects of segregation on the health outcomes of racial minorities, it is plausible to assume that segregation will be a major contributor to racial gaps in mortality across the country.

The findings from this study should encourage epidemiologists to include features of the built environment in their mathematical models forecasting the spread of diseases. Compartmental models such as the susceptible–infected–removed model (55) are powerful first-order approaches to predict the evolution of infectious diseases, but there is ample room to add complexity

and heterogeneity to these models to better reflect the spatial clustering of population groups.

## Materials and Methods

**Double-Lasso Regression.** Theory and prior evidence have informed the data collection strategy that leads to the set of 50 possible controls shown in Fig. 1. But, from a statistical standpoint, it is unclear whether all of them should be included in a traditional OLS regression. Covariates strongly correlated with COVID-19 outcomes and segregation are necessary to avoid omitted variable bias. Including covariates that correlate primarily with the outcome will remove residual variation in the outcome and ensure a more precise estimate of the segregation coefficient. Including covariates that strongly correlate with segregation but not so much with the outcome will unnecessarily inflate the variance of the estimated coefficient on segregation. And, given the nature and novelty of the outcome, the mortality and spread of a new virus, we may not have clear priors as to which variables must be included. Leaving model selection up to the researcher could raise the suspicion that the selected controls are the ones that fit the statistically significant story that we want to tell.

Recent developments in the econometrics literature leverage the power of machine learning to guide the principled selection of variables when these are many and the correct functional form is unknown. Belloni et al. (29) propose a "double-lasso" approach that identifies the relevant covariates to be included. The method uses lasso regression (56) to select covariates in two stages, first choosing those that predict the outcome and then those that predict the independent variable of interest.

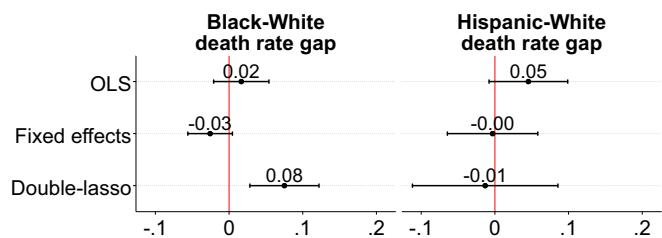
Before going over the two steps in the double-lasso regression, it is important to understand the advantage that the lasso methodology has over OLS regression. Given a model with  $p$  independent variables of this form,

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \epsilon, \quad [1]$$

the OLS solution would find the vector of  $\beta$  coefficients that minimizes the following objective function:

$$\hat{\beta}^{OLS} = \arg \min_{\beta} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}. \quad [2]$$

If  $X_1$  is the segregation covariate,  $\hat{\beta}_1$  will be the estimate of the association between segregation and the corresponding COVID-19 outcome. The problem we face is that the number of  $p$  covariates could be very large. Some of these controls may be redundant and others irrelevant, and, for some others, we may not have any priors as to whether they should be in the model.



**Fig. 4.** OLS and double-lasso regression estimates of the relationship between the Black–White and Hispanic–White Relative Diversity Indices and racial gaps in COVID-19 death rates. OLS models include no controls. Fixed effects models include state fixed effects. Double-lasso models include the 18 controls selected by the lasso procedure (shown in *SI Appendix, Table S4*) and state fixed effects. All regressions include population weights. Standard errors are clustered by state. Bars around estimated coefficients reflect 95% CIs. Segregation is measured with the Black–White and Hispanic–White Relative Diversity Indices. The Black–White (Hispanic–White) death rate gap is the difference between the log death rate for Blacks (Hispanics) and the log death rate for Whites. The Relative Diversity Indices and covariates have been standardized to have mean 0 and SD 1. The sample includes 243 counties for the Black–White model and 218 for the Hispanic–White model. *SI Appendix, Figs. S9 and S10* show results when using equal samples ( $n = 180$ ) across both models. COVID-19 outcomes are measured as of September 30, 2020.

The lasso provides us with a framework to perform a statistically principled selection of the controls that must feature in the model. Instead of minimizing the objective function in Eq. 2, the lasso shrinks the regression coefficients toward zero by adding a penalty term. The penalty term, known as the  $L1$  norm, forces the sum of the absolute values of the regression coefficients to be as small as possible, which reduces overfitting by discouraging complex models. Because of this constraint, the lasso will select the independent variables that contribute the most to minimizing the sum of squared errors, dropping those that contribute nothing and optimally shrinking the rest. Formally,

$$\beta^{\text{lasso}} = \arg \min_{\beta} \left\{ \frac{1}{2} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}. \quad [3]$$

In Eq. 3,  $\lambda$  is the tuning parameter and controls the "strength" of the constraint that forces the sum of the absolute values of the regression coefficients to be as small as possible. Different values of  $\lambda$  will yield different sets of coefficients in  $\beta^{\text{lasso}}$ . If we set  $\lambda$  to zero,  $\beta^{\text{lasso}} = \beta^{\text{ols}}$ . If  $\lambda = \infty$ , all coefficients are set to zero. The choice of the optimal value of the lasso penalty can be made by cross-validation, using a plug-in iterative formula, or using an adaptive approach. I use the plug-in iterative formula because it yields a more parsimonious set of covariates.

The lasso is a powerful regression tool if we are interested in predicting the outcome, but it limits us in a fundamental way if we are interested in making inferences about any of the coefficients in  $\beta^{\text{lasso}}$ . By allowing the lasso to shrink all coefficients toward zero (and setting some to zero), it is possible that the coefficient on segregation,  $\beta_1$ , will be shrunk and thus biased. Similarly, some of the coefficients on the controls that must be in the model to avoid the omitted variable bias problem could also be shrunk, which would also bias the estimate of  $\beta_1$  (some of these controls

1. M. W. Hooper, A. M. Nápoles, E. J. Pérez-Stable, COVID-19 and racial/ethnic disparities. *J. Am. Med. Assoc.* **323**, 2466 (2020).
2. C. W. Yancy, COVID-19 and African Americans. *J. Am. Med. Assoc.* **323**, 1891–1892 (2020).
3. UK Office for National Statistics, "Coronavirus (COVID-19) related deaths by ethnic group, England and Wales" (Office for National Statistics, 2020). <https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/articles/coronavirusrelateddeathsbyethnicgroupenglandandwales/2march2020to10april2020>. Accessed 10 July 2020.
4. M. R. Kramer, C. R. Hogue, Is segregation bad for your health? *Epidemiol. Rev.* **31**, 178–194 (2009).
5. D. R. Williams, P. B. Jackson, Social sources of racial disparities in health. *Health Aff.* **24**, 325–334 (2005).
6. K. D. Hart, S. J. Kunitz, R. R. Sell, D. B. Mukamel, Metropolitan governance, residential segregation, and mortality among African Americans. *Am. J. Public Health* **88**, 434–438 (1998).
7. M. O. Hearst, J. M. Oakes, P. J. Johnson, The effect of racial residential segregation on black infant mortality. *Am. J. Epidemiol.* **168**, 1247–1254 (2008).
8. R. S. Cooper et al., Relationship between premature mortality and socioeconomic factors in black and white populations of us metropolitan areas. *Publ. Health Rep.* **116**, 464–473 (2016).
9. S. A. Jackson, R. T. Anderson, N. J. Johnson, P. D. Sorlie, The relation of residential segregation to all-cause mortality: A study in black and white. *Am. J. Public Health* **90**, 615–617 (2000).
10. C. A. Collins, Racism and health: Segregation and causes of death amenable to medical intervention in major us cities. *Ann. N. Y. Acad. Sci.* **896**, 396–398 (1999).
11. J. Fang, S. Madhavan, W. Bosworth, M. H. Alderman, Residential segregation and mortality in New York city. *Soc. Sci. Med.* **47**, 469–476 (1998).
12. F. B. LeClere, R. G. Rogers, K. D. Peters, Ethnicity and mortality in the United States: Individual and community correlates. *Soc. Forces* **76**, 169–198 (1997).
13. C. M. Masi, L. C. Hawley, Z. H. Piotrowski, K. E. Pickett, Neighborhood economic disadvantage, violent crime, group density, and pregnancy outcomes in a diverse, urban population. *Soc. Sci. Med.* **65**, 2440–2457 (2007).
14. R. A. Rodriguez et al., Geography matters: Relationships among urban residential segregation, dialysis facilities, and patient outcomes. *Ann. Intern. Med.* **146**, 493–501 (2007).
15. D. Acevedo-Garcia, Zip code-level risk factors for tuberculosis: Neighborhood environment and residential segregation in New Jersey, 1985–1992. *Am. J. Public Health* **91**, 734–741 (2001).
16. K. P. Fennie, K. Lutfi, L. M. Maddox, S. Lieb, M. J. Trepka, Influence of residential segregation on survival after AIDS diagnosis among non-Hispanic blacks. *Ann. Epidemiol.* **25**, 113–119 (2015).
17. L. Downey, US metropolitan-area variation in environmental inequality outcomes. *Urban Stud.* **44**, 953–977 (2007).
18. R. Morello-Frosch, B. M. Jesdale, Separate and unequal: Residential segregation and estimated cancer risks associated with ambient air toxics in US metropolitan areas. *Environ. Health Perspect.* **114**, 386–393 (2006).
19. Centers for Disease Control and Prevention, Prevalence of diabetes among Hispanics—Selected areas, 1998–2002. *MMWR Morb. Mortal. Wkly. Rep.* **53**, 941–944 (2004).
20. Centers for Disease Control and Prevention, Differences in prevalence of obesity among black, white, and Hispanic adults—United States, 2006–2008. *MMWR Morb. Mortal. Wkly. Rep.* **58**, 740–744 (2009).
21. R. L. Sacco, W. Hauser, J. Mohr, Hospitalized stroke in blacks and Hispanics in northern Manhattan. *Stroke* **22**, 1491–1496 (1991).
22. R. O. White, B. M. Beech, S. Miller, Health care disparities and diabetes care: Practical considerations for primary care providers. *Clin. Diabetes* **27**, 105–112 (2009).
23. US Bureau of Labor Statistics, Employed persons by detailed industry, sex, race, and Hispanic or Latino ethnicity. Labor Force Statistics from the Current Population Survey. <https://www.bls.gov/cps/cpsaat18.htm>. Accessed 10 July 2020.
24. S. Mollborn, P. Romby, J. A. Dennis, Who matters for children's early development? Race/ethnicity and extended household structures in the United States. *Child Indicators Res.* **4**, 389–411 (2011).
25. A. P. Bartel et al., Racial and ethnic disparities in access to and use of paid family and medical leave: Evidence from four nationally representative datasets. *Mon. Labor Rev.* **10.21916/mlr.2019.2** (2019).
26. M. L. Small, M. McDermott, The presence of organizational resources in poor urban neighborhoods: An analysis of average and contextual effects. *Soc. Forces* **84**, 1697–1724 (2006).
27. S. Curranini, M. O. Jackson, P. Pin, An economic model of friendship: Homophily, minorities, and segregation. *Econometrica* **77**, 1003–1045 (2009).
28. M. McPherson, L. Smith-Lovin, J. M. Cook, Birds of a feather: Homophily in social networks. *Annu. Rev. Sociol.* **27**, 415–444 (2001).
29. A. Belloni, V. Chernozhukov, C. Hansen, Inference on treatment effects after selection among high-dimensional controls. *Rev. Econ. Stud.* **81**, 608–650 (2014).
30. A. Belloni, V. Chernozhukov, Y. Wei, Post-selection inference for generalized linear models with many controls. *J. Bus. Econ. Stat.* **34**, 606–619 (2016).
31. T. C. Yang, S. A. Matthews, Death by segregation: Does the dimension of racial segregation matter? *PloS One* **10**, e0138489 (2015).
32. L. J. Krivo, R. D. Peterson, D. C. Kuhl, Segregation, racial structure, and neighborhood violent crime. *Am. J. Sociol.* **114**, 1765–1802 (2009).
33. R. D. Peterson, L. J. Krivo, Racial segregation and black urban homicide. *Soc. Forces* **71**, 1001–1026 (1993).
34. E. S. Shihadeh, N. Flynn, Segregation and crime: The effect of black social isolation on the rates of black urban violence. *Soc. Forces* **74**, 1325–1352 (1996).
35. R. S. Cooper, Social inequality, ethnicity and cardiovascular disease. *Int. J. Epidemiol.* **30**, S48 (2001).
36. K. Morland, S. Wing, A. D. Roux, C. Poole, Neighborhood characteristics associated with the location of food stores and food service places. *Am. J. Prev. Med.* **22**, 23–29 (2002).
37. J. R. Lincoln, J. Miller, Work and friendship ties in organizations: A comparative analysis of relation networks. *Adm. Sci. Q.* **24**, 181–199 (1979).
38. J. Moody, Race, school integration, and friendship segregation in America. *Am. J. Sociol.* **107**, 679–716 (2001).

could be excluded altogether if they have a small predictive power on the outcome).

The double-lasso regression is designed to solve this problem in two intuitive steps. The idea is to fit two separate lasso regressions, one of the outcome  $Y$  on all covariates and another of the independent variable of interest,  $X_1$ , on all covariates. These two regressions will yield two sets of nonzero coefficients, one for each of the respective lassos. The final step is to fit a linear regression of the outcome  $Y$  on  $X_1$  and the set of covariates with nonzero coefficients selected by any of the two lassos. In more concrete terms and taking Eq. 1 as our model of interest, where  $X_1$  is the measure of segregation and  $\beta_1$  represents the relationship between segregation and COVID-19 mortality, the double-lasso regression proceeds as follows:

- 1) Fit a lasso of  $Y$  on  $X$ , where  $X$  includes all covariates except segregation,  $X_2, \dots, X_p$ . Identify the set of covariates with nonzero coefficients in this regression and label it  $X_y$ .
- 2) Fit a lasso of  $X_1$  on  $X$ . Identify the set of covariates with nonzero coefficients in this regression and label it  $X_s$ .
- 3) Select the union of covariates in  $X_y$  and  $X_s$  and label it  $X_u$ .
- 4) Fit a linear regression of  $Y$  on  $X_1$  and  $X_u$  to obtain an estimate of  $\beta_1$ .

We can intervene in the process of selecting covariates by forcing some of them to be included in  $X_u$ . We may want to do that if we know that some covariates are absolutely necessary to obtain an unbiased estimate of  $\beta_1$ . In this case, I force the state indicators to be included as controls in  $X_u$ . The covariates selected by the lasso procedure are shown in *SI Appendix, Table S3*.

**Data Availability.** Data and code have been deposited in Harvard Dataverse (<https://doi.org/10.7910/DVN/JHFOSE>).

**ACKNOWLEDGMENTS.** I thank Shamus Khan, Bruce Western, Josh Whitford, Andreas Wimmer, and two anonymous reviewers for their helpful comments.

39. E. Stearns, C. Buchmann, K. Bonneau, Interracial friendships in the transition to college: Do birds of a feather flock together once they leave the nest? *Sociol. Educ.* **82**, 173–195 (2009).
40. A. Wimmer, K. Lewis, Beyond and below racial homophily: ERG models of a friendship network documented on Facebook. *Am. J. Sociol.* **116**, 583–642 (2010).
41. Centers for Disease Control and Prevention, Hospitalization rates and characteristics of patients hospitalized with laboratory-confirmed coronavirus disease 2019—COVID-NET, 14 states, March 1–30, 2020. *MMWR Morb. Mortal. Wkly. Rep.* **69**, 458–464 (2020).
42. S. Richardson *et al.*, Presenting characteristics, comorbidities, and outcomes among 5700 patients hospitalized with COVID-19 in the New York City area. *J. Am. Med. Assoc.* **323**, 2052–2059 (2020).
43. S. Brayne, Surveillance and system avoidance: Criminal justice contact and institutional attachment. *Am. Socio. Rev.* **79**, 367–391 (2014).
44. A. Goffman, On the run: Wanted men in a Philadelphia ghetto. *Am. Socio. Rev.* **74**, 339–357 (2009).
45. S. F. Reardon, G. Firebaugh, Measures of multigroup segregation. *Socio. Methodol.* **32**, 33–67 (2002).
46. CDC COVID-19 Response Team, Preliminary estimates of the prevalence of selected underlying health conditions among patients with Coronavirus disease 2019—United States, February 12–March 28, 2020. *MMWR Morb. Mortal. Wkly. Rep.* **69**, 382–386 (2020).
47. L. F. Berkman, Social epidemiology: Social determinants of health in the United States: Are we losing ground? *Annu. Rev. Publ. Health* **30**, 27–41 (2009).
48. M. Marmot, Social determinants of health inequalities. *Lancet* **365**, 1099–1104 (2005).
49. D. Acevedo-Garcia, K. A. Lochner, T. L. Osypuk, S. V. Subramanian, Future directions in residential segregation and health research: A multilevel approach. *Am. J. Public Health* **93**, 215–221 (2003).
50. I. G. Ellen, T. Mijanovich, K. N. Dillman, Neighborhood effects on health: Exploring the links and assessing the evidence. *J. Urban Aff.* **23**, 391–408 (2001).
51. K. A. Frank, Impact of a confounding variable on a regression coefficient. *Socio. Methods Res.* **29**, 147–194 (2000).
52. E. Oster, Unobservable selection and coefficient stability: Theory and evidence. *J. Bus. Econ. Stat.* **37**, 187–204 (2019).
53. B. Bolin, L. C. Kurtz, "Race, class, ethnicity, and disaster vulnerability" in *Handbook of Disaster Research*, H. Rodriguez, W. Donner, J. E. Trainor, Eds. (Springer, 2018), pp. 181–203.
54. J. T. Lariscy, R. A. Hummer, M. D. Hayward, Hispanic older adult mortality in the United States: New estimates and an assessment of factors shaping the Hispanic paradox. *Demography* **52**, 1–14 (2015).
55. W. O. Kermack, A. G. McKendrick, A contribution to the mathematical theory of epidemics. *Proc. Roy. Soc. Lond.* **115**, 700–721 (1927).
56. R. Tibshirani, Regression shrinkage and selection via the lasso. *J. Roy. Stat. Soc. B* **58**, 267–288 (1996).