

Ensemble Machine Learning of Factors Influencing COVID-19 Across US Counties

David McCoy

Division of Environmental Health Sciences, UC Berkeley <https://orcid.org/0000-0002-5515-6307>

Whitney Mgbara

Dept. Env. Science, Policy and Management, UC Berkeley

Nir Horvitz

School of Mathematical Sciences, University of Kwazulu-Natal

Wayne M. Getz

Dept. Env. Science, Policy and Management; School of Mathematical Sciences, University of Kwazulu-Natal

Alan Hubbard (✉ hubbard@berkeley.edu)

Division Biostatistics, UC Berkeley

Research Article

Keywords: COVID-19, Machine Learning, Super Learner, Variable Importance

Posted Date: October 14th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-90547/v1>

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

[Read Full License](#)

Ensemble Machine Learning of Factors Influencing COVID-19 Across US Counties

David McCoy^{a,1}, Whitney Mgbara^{b,1}, Nir Horvitz^c, Wayne M. Getz^{b,c}, Alan Hubbard^a

Corresponding author: Alan Hubbard

¹ *Co-first authors listed alphabetically*

*All code for collecting data, collected data, statistical scripts, up-to-date outcome data, visualizations, and statistical results are available on GitHub:
https://github.com/blind-contours/Getz_Hubbard_Covid_Ensemble_ML_Public.git
email addresses: david_mccoy@berkeley.edu, wmgbara@berkeley.edu, nirh@berkeley.edu, wgetz@berkeley.edu, hubbard@berkeley.edu*

^a *Division Environmental Health Sciences, UC Berkeley, CA 94720, USA*

^b *Dept. Env. Science, Policy and Management, UC Berkeley, CA 94720, USA*

^c *School of Mathematical Sciences, University of Kwazulu-Natal, Durban 4000, South Africa*

^d *Division Biostatistics, UC Berkeley, CA 94720-3114, USA*

Abstract

Background: Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) the causal agent for COVID-19, is a communicable disease spread through close contact. It is known to disproportionately infect certain communities due to both biological susceptibility and inequitable exposure. In this study, we investigate the most important health, social, and environmental factors impacting both the early and later phases of COVID-19 transmission and mortality in US counties.

Methods: We aggregate county-level physical and mental health, environmental pollution, access to health care, demographic characteristics, vulnerable population scores, and other epidemiological data to create a large feature set to analyze COVID-19 outcomes. Because of the high-dimensionality and multicollinearity of the data, we use ensemble machine learning and marginal prediction methods to identify the most salient factors associated with several COVID-19 outbreak measures.

Findings: Our variable importance results show that measures of ethnicity,

public transportation and preventable diseases are the strongest drivers for both incidence and mortality. Specifically, the CDC measures for minority populations, CDC measures for limited English, and proportion of Black/African-American individuals in a county were the most important features for COVID-19 cases at day 25 and to date. For mortality at day 100 and total to date, we find that public transportation use and proportion of Black/African-American individuals in a county are the strongest predictors. The methods predict that, keeping all other factors fixed, a 10% increase in public transportation use increases mortality at day 100 by 2012 (95% CI [1972, 2356]) and likewise a 10% increase in the proportion of Black/African-American individuals in a county increases total deaths to date by 2067 (95% CI [1189, 2654]). In terms of cases to date, ethnicity turns out to almost twice as important as the next most important factors, which are location, disease prevalence, and transit factors.

Interpretations: Our findings indicate that a more focused approach should be taken when managing COVID-19, by considering features of the economy most responsible for transmission and sectors of society most vulnerable to infection and mortality. In particular, our results strongly reinforce others pointing to the disproportionate impact of COVID-19 on minority populations. They also suggests that mitigation measures, including rolling out vaccinations as they become available, will be most efficacious for the US population as a whole when, beyond healthcare workers and first responders, are focused first on the highest-risk communities.

Funding: UC Berkeley, Biomedical Big Data Training Fellowship; NSF Grant 2032264 to WMG and AH.

Keywords: COVID-19, Machine Learning, Super Learner, Variable Importance

1 **1. Introduction**

2 *COVID-19 Background.* Coronavirus disease 2019 (COVID-19), caused by the
3 novel severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), has spread
4 rapidly around the world resulting in a pandemic. [1, 2, 3, 4, 5, 6]. Within a
5 month of discovering the first cluster of cases in Wuhan, China [1, 6, 3], 18
6 additional countries had reported a case of COVID-19 [7]. The World Health
7 Organization declared the resulting outbreaks a Public Health Emergency of
8 International Concern by January 30, 2020 and a pandemic by March 11, 2020
9 [7, 8].

10 In the United States, to the best of our knowledge, the COVID-19 pandemic
11 began to spread on January 21, 2020 when an infected individual traveled from
12 Wuhan, China to Snohomish County, Washington [9]. Community circulation,
13 within US states and territories, followed shortly after [10, 11]. By May 20,
14 2020, the U.S. had the most confirmed active cases and the most deaths of any
15 country, with a reported total of 1,559,750 cases and 92,333 deaths across 55
16 US jurisdictions including 50 states, District of Columbia, Guam, Puerto Rico,
17 the Northern Mariana Islands, and the U.S. Virgin Islands [8].

18 There are a few consistent observations regarding the epidemiology of the
19 COVID-19 pandemic in the US. Most prominent is the relatively high infection
20 and death rates of minority populations, particularly Black/African Americans
21 [12] [13], a disparity researchers have noted occurred in previous pandemics,
22 such as HIV [14]. This disparity has been observed in both adults and children
23 [15]. There is much previous work on the causes of health disparities among
24 Black/African Americans, and others have speculated on which of these causes
25 are related to the differential impact of COVID-19 [16] [17]. Thus, to tease out
26 the impact on vulnerable groups, one needs data on other baseline health factors,
27 such as obesity [18] [19], co-morbidities, age [20] [21] environmental exposures,
28 transportation use and employment factors, including types of occupations [13]
29 [22].

30 *US Health Response and Forecasting.* To control the spread of COVID-19, the
31 US has implemented complex and regionally uneven community-level, non-
32 pharmaceutical interventions, including travel restrictions, social distancing,
33 and stay-at-home orders. These interventions aim to reduce the opportuni-
34 ties for person-to-person and fomite-to-person contact routes for transmission
35 of the virus. Although these interventions have shown to mitigate the commu-
36 nity spread in certain communities, the trend did not hold for all communities.
37 Many counties experienced an uptick in cases after an initial decline. There are
38 several reasons why certain communities continue to see a growing number of
39 cases, including: 1) lifting shelter-in-place or other social distancing restrictions
40 earlier than advised; 2) lax controls on gatherings that resulted in super-spreader
41 events [23], and; 3) unknown affects of plausible seasonality that impacts viral
42 transmission [24]. As such, complex epidemiological contexts have emerged in
43 US communities. The complexity is a result of dynamic environmental factors
44 constituting social and physical environments for US populations that impact
45 an individual’s risk for contracting COVID-19. Thus, to adequately control the
46 spread of COVID-19, it is important to identify the most salient social and phys-
47 ical environmental factors within US communities, driving transmission and and
48 effecting susceptibility.

49 Despite awareness that disease transmission is related to social and physi-
50 cal environmental factors, few studies have rigorously analyzed the underlying
51 drivers of the dynamics of disease emergence for COVID-19 in the US. Though
52 models accounting for the specific vulnerabilities of local populations have been
53 proposed, only a few models exist that assess the importance county-level vari-
54 ation of such variables in fueling COVID-19 outbreaks [25, 26]. Altieri et al.
55 [25], for example, use county-level data from similar sources to this paper, to
56 create an ensemble forecasting model, dubbed “Combined Linear and Exponen-
57 tial Predictors” to predict death counts from COVID-19. Their goal is to curate
58 a data repository that can be used to forecast exponential and sub-exponential
59 cases weeks in advance in order to help nonprofit organization disseminate much
60 needed personal protective devices and respirators to areas projected to have

61 higher mortality rates.

62 In our study, we use similar data sources to model cross-sectional COVID-19
63 outcomes. Our goal, rather than prediction, is to explore the relative importance
64 of different types of social, physical and environmental factors on COVID-19
65 transmission and mortality. Hospital case and mortality data, and seropositive
66 surveillance studies have shown there are subgroups of the population that are
67 more susceptible to higher cases of morbidity and mortality. These include
68 people older than 65 and communities affected by racial disparities. We
69 attempt in this paper to expand on previous studies of county level variation
70 in COVID-19 (e.g., see [27]) by evaluating additional socio-environmental data
71 to understand if these disparities have direct effects on COVID-19 outcomes or
72 are indirect through additional risk factors, such as diabetes, food security, air
73 pollution or access to health care. Likewise, no study has been carried out to
74 determine which of these factors are most associated with COVID-19 outcomes
75 while still controlling nonparametrically for all other factors.

76 *Added Value of this Study.* We aggregate 5 types of COVID-19 outcomes, (i)
77 day of the first case in a county relative to the first case reported in the U.S.
78 (Snohomish County, Washington), (ii) number of cases 25 days after the initial
79 case in a county, (iii) number of all-cause deaths 100 days after the first case in
80 a county, (iv) total number of cases in a county to-date after initial case and (v)
81 total number of all-cause deaths in a county to-date after first case in a county.
82 From many sources including the CDC, U.S. Census Bureau, USA facts, google
83 mobility data, and others, we collect a large number of pertinent variables for
84 COVID-19. Using ensemble machine learning (ML) we create models that make
85 no assumptions on the distributions of the data, these models are thereby non-
86 parametric and allow all degrees of interactions and distributions in order to
87 make the best fit.

88 In each of these models we identify the most important variables in the model
89 by removing the variable and measuring the difference in model risk (model er-
90 ror). We then make marginal predictions for the number of cases and mortalities

91 from COVID-19 when increasing or decreasing these top variables while control-
92 ling for all other factors. Confidence intervals (CIs), significance and robustness
93 of findings are measured via bootstrapping the model. Additionally, we in-
94 vestigate the predicted number of cases and mortalities from our model when
95 controlling for only variables outside the target variable category (i.e ethnicity,
96 public transportation subcategories) and univariate predictions (not controlling
97 for other variables). Given this approach, our contributions are, rather than
98 predictive forecasting the number of cases, an approach to measure the relative
99 importance of risk factors for COVID-19 from an environmental perspective.

100 There are many known risk factors measured from case hospitalization data,
101 including diabetes and heart disease. As such, we hypothesize these factors will
102 show high variable importance, specifically for case mortality. Additionally, we
103 hypothesize that environmental dynamics, which increase exposure time to the
104 virus (i.e., number of occupations in a county, public transit use), will also be
105 strongly associated with COVID-19 cases. When such factors are identified,
106 this information could also be used to update or improve the public health
107 response to specifically target factors related to high case counts in order to
108 further mitigate new cases or prevent a resurgence of cases. Finally, we estimate
109 variable importance of the high dimensional assembled county features without
110 constraining (and bias-inducing), using a combination of machine learning and
111 intuitive substitution estimators. Lastly, we make both the data and methods
112 used in this paper accessible to others, thereby providing open source access
113 and enhancing the utility of our results. All code, data, and county variables
114 used are available and outcome data are updated daily on GitHub.

115 **2. Methods**

116 *Data Sources.* For all outcomes and predictors we compiled publicly available
117 data for US jurisdictions reported at the state-level (e.g. Google Mobility Data)
118 and county-level, excluding Alaska, Hawaii, and Puerto Rico. Our final dataset
119 includes county-level case counts, death counts and a wide variety of county-level

120 demographic, epidemiological, health, and environmental data used as predic-
121 tors in our analysis. Our analysis was based on the cumulative confirmed cases
122 and deaths of COVID-19 in US counties, starting on January 22 2020 (referred
123 to as Day 1) until July 14 2020 (USA Facts). The sources for our county-level
124 features were: USA Facts, Bureau of Economic Analysis, American Community
125 Survey, Tiger/Line GeoDatabase, CDC Interactive Atlas of Heart Disease and
126 Stroke, County Health Rankings, Centers of Medicare and Medicaid Services,
127 National Centers for Environmental Information, CDC Vulnerability Index, Bu-
128 reau of Labor Statistics, MIT Election Lab, Google Community Mobility Re-
129 ports; a total of 12 different sources which were joined on county FIPS codes.

130 *Outcomes.* We use five COVID-19 outcome scenarios: 1) we transformed the
131 case data into day of the first case after the first confirmed case on January 21
132 2020; 2) we determined the number of cumulative cases on the 25 days after
133 the day of the first case for each county (i.e., day 25 of the outbreak in each
134 county); 3) we used the number of all-cause deaths on day 100 after the day of
135 the first case of the outbreak for each respective county up to July 14 2020; 4)
136 we determine the total number of cases to date after the day of the first case for
137 each county and likewise scenario five is the total number of all-cause mortalities
138 to-date after the day of the first case for each county.

139 For each outcome variable excluding number of cases at day 25, we divide the
140 counts by population size for each county to create per-capita COVID-19 case
141 and all-cause mortality outcome. This was done to remove population size bias.
142 Likewise, all predictor variables measured in counts were also standardized by
143 population size. All-cause deaths were used rather than reported fatalities due
144 to COVID-19 for several reasons related to unreliable case data, differences in
145 testing, and co-morbidities between COVID-19 and other fatal acute diseases.
146 By using all-cause deaths measured since day of first case reported in a county,
147 we hope to get a better estimate on the impact of COVID-19 on mortality.
148 As discussed, our predictors cover a wide scope, **Table 1** gives a review of
149 the data sources and variables collected from each source. **Table S1** in the

150 supplementary section gives more details of this process.

151 *Predictors.* Data on our predictor variables include demographics, health re-
152 source availability, health risk factors, social vulnerability, and other COVID-
153 19-related information. The predictor variables used are collected from different
154 sources including Altieri et al. [28] and Killeen et al. [29]. Because the aims
155 of this paper are not purely predictive, but focus on understanding the relative
156 impact each variable has on COVID-19 outcomes, our data curation process is
157 different when compared to these two papers. We aggregate data from differ-
158 ent stratified variables to create an overall public transportation use feature.
159 Likewise for social vulnerability scores, we attempt to include core variable that
160 represent a specific type of risk feature. For example, given our interest are
161 variable importance measures and marginal predictions, if we included both the
162 aggregated CDC vulnerability index with many features collected from other
163 sources that are proxy measures for this index (percent non-English speaking,
164 poverty levels etc.) then findings for aggregate vulnerability index would be
165 conservative given these other variables are also included in the model. That
166 being said, given the large overlap and interactions of all variables collected, it
167 is likely that all variable importance estimates are conservative. In our cura-
168 tion process, only unstratified variables are used (not stratified by age or sex),
169 we also create sub-categories for variables (ethnicity, geography, disease preven-
170 tion, etc.) These sub-categories we use in later analysis to explore possible over-
171 correction of the model for a respective target variable.

172 Briefly, some predictor variables include proportions of individuals by poverty
173 level, gender, age distribution, race distribution, household income, healthcare
174 access, occupation type, and so on, these were collected from USA Facts and the
175 Census Bureau. Airport data, including, distance from county polygon centroid
176 to airports were calculated from the Federal Aviation Administration. The 2020
177 county health rankings and Center for Medicaid and Medicare services were used
178 to gather information on a range of health data including smoking, diabetes,
179 obesity, air pollution and many other physical and mental health metrics. Pre-

180 cipitation by month was gathered from the National Oceanic and Atmospheric
181 Association. Vulnerability index scores from each theme (described in **Table**
182 **S1**) were aggregated from the Center of Disease Control. In total, over 150
183 predictor variables were gathered before curation. This curated data along with
184 all relevant code, documentation and results are provided on our GitHub page.

185 *Data Cleaning and Curation.* The data curation process is described in more
186 detail in the supplementary information alongside the data dictionary. All re-
187 sulting data was numerical (no factor variables). In addition, we screen out any
188 variables with more than 70% missing values. Similarly, we removed variables
189 with close to zero variance. Variables that were nearly perfectly correlating were
190 removed (Pearson correlation = 0.95). Missing data in this cleaned dataset were
191 imputed with the mean. For Google mobility data, we scraped data from the
192 published mobility trend reports from Feb 16 2020 to March 29 2020. These
193 data represent the general increase or decrease in movement to the respective
194 destination (grocery stores, parks etc.) compared to baseline (pre-pandemic
195 period). To create an aggregate score representing the mobility trend for each
196 movement category for each county, we use the slope from linear regression
197 to measure this trend over time. The slope for each movement category was
198 included in our Super Learner models.

199 *Exploratory Analysis.* To graphically represent how our feature data are related
200 to one another, and likewise how counties are related to one another through
201 these variables, we use unsupervised hierarchical agglomerative clustering of
202 both county features and counties. We present the results of this clustering as a
203 heatmap. Our first goal of this method is to understand if there are counties that
204 have a trend for early first case reported, high COVID-19 case rates at day 25,
205 and COVID-19 case rates to-date, as well as high all-cause mortalities at day 100
206 and to-date. If this trend was seen, we next wanted to investigate what variables
207 were ‘highly expressed’ in these counties. As such, all feature data was z-score
208 standardized. We then took the quantiles of each outcome to create factor
209 dummy variables that can be plotted alongside clustering of counties. Clustering

210 was done for both counties and county features and reordered accordingly using
211 Euclidean distances. As this is an unsupervised approach, outcome data are
212 not included in this machine learning method but are simply plotted alongside
213 the clustering results to visually identify correspondences. Groups of county
214 features that were found to be associated with groups of COVID-19 outcomes
215 are presented.

216 *Machine Learning Pipeline.* Although our ultimate goal is not to use our final
217 models for forecasting and prediction, one still needs to estimate a regression
218 model in order to determine our measures of variable importance. By estimat-
219 ing this model as accurately as possible, one can better estimate the variable
220 importance measures that rely on the prediction model. As such, instead of
221 choosing one machine learning (ML) algorithm to model county features for
222 each outcome, we used an ensemble approach (Super Learner) to fit a prediction
223 function for each of our outcomes. The Super Learner combines the predictive
224 probabilities of COVID-19 outcomes across many ML algorithms. The Super
225 Learner finds the optimal combination of a collection of algorithms by mini-
226 mizing the cross-validated risk [30, 31]. This method is an improvement over
227 methods using only one ML algorithm because no one algorithm is universally
228 optimal. The Super Learner has been shown in theory to be at least as good as
229 the best performing algorithm in the ensemble and often times performs con-
230 siderably better [32]. Given the high-dimensionality and complex relationships
231 of the county data and the fact that there are no known distributions for these
232 relationships, we chose a wide range of algorithms for the ensemble in order to
233 optimize performance. For COVID-19 cases at day 25 and to-date per-capita
234 rates and all-cause mortalities at day 100 and total to-date per capita, we use
235 a large number of linear Gaussian based algorithms including conditional mean
236 (control algorithm) simple generalized linear model, a series of penalized re-
237 gressions setting alpha at levels to create ridge regression, lasso regression, and
238 elastic net regression. Similarly, we use a number of gradient boosted decision
239 trees that differed in depth and shrinkage. Because these algorithms require

240 hyper-parameter tuning for optimal performance, we create a grid of all possible
241 hyper-parameters and choose algorithms across this grid for inclusion in
242 the ensemble. The same procedure was applied for day-of-first-case in a county
243 relative to-day-of-first case in the U.S. Instead of using Gaussian algorithms,
244 however, custom learners were made for the Super Learner environment that
245 model Poisson outcome data. The same parameters were chosen for this set of
246 learners to create the Poisson ensemble. To address possible over-fitting and
247 to get cross-validated risks for each algorithm in the ensemble sets, five-fold
248 cross-validation was used for internal SL cross-validation both to build optimal
249 models with each classifier and to determine optimal weighting across classifiers
250 in the ensemble.

251 *Unpacking the black box.* The algorithm used to create our predictor given co-
252 variates has desirable optimality properties, being asymptotically guaranteed to
253 have a fit as good as any of the candidate members (the "oracle property"),
254 with no risk of overfitting. If a library of both smooth (e.g., parametric models)
255 and flexible, nonparametric learners, then one can find it hard to outperform
256 [33, 31]. However, the result is a black box that creates predictions as a complex
257 ensemble of different learners, some having their own internal variable selection
258 process and model selection framework. Thus, the resulting black box needs
259 to be intelligently queried to estimate the independent impact of the various
260 predictors used in the model. We do so in two ways. One, is using a straightfor-
261 ward leave-one-variable out method and re-examining the change in prediction
262 accuracy. However, this provides no information about the direction of the im-
263 pact, which is why we follow with a query inspired by causal inference methods.
264 In that case, we use the model to forward model situations where we change
265 the distribution of predictors across the counties in sequential fashion and then
266 calculate the marginal predicted counts (so called substitution estimators, or G-
267 computation - [34, 35]). The combination of these two versions of non-parametric
268 variable importance measures provide both the importance of the variable (or
269 sub-category of variable) to the resulting predictor as well as an intuitive mea-

270 sure of the adjusted association of single variables.

271 *Variable Importance.* We built Super Learners from the same county-level data
272 for each of the COVID-19 outcomes. To measure variable importance in each
273 model, we take the fitted model and make predictions using all county features
274 and measure the model risk (average in-squared differences in model prediction
275 versus truth, or mean-squared error). We then scramble (sets) of variables and
276 re-do the prediction and derive the new MSE. The plots (**Figures 2 and 3**)
277 show the resulting ranked list of variables (most to least change in the MSE by
278 scrambling). We use a risk-ratio (MSE-ratio) for each variable to measure its
279 relative importance in the model for each outcome. The risk-ratio is the risk
280 in the model without the respective variable (numerator) over the risk when
281 the variable is included in the model (denominator). As such, a risk-ratio of
282 1.5 indicates that the model MSE rises by 50% when the variable is scrambled
283 while controlling for all other variable affects. We use a similar approach to
284 measure the variable sub-category importance on each outcome. Each variable
285 was given a sub-category (described in supplementary material) resulting in a
286 total of 15 categories. Blocks of variables in each category were scrambled and
287 the model risk-ratio measured to attain information on category importance.

288 *Marginal Predictions.* Given that we fit a black-box to derive our prediction
289 models, we have to unpack the black-box to understand what it implies about
290 the adjusted relationship of the important variables to the outcomes. We thus
291 use substitution methods to evaluate the predicted change in the mean if county
292 characteristics are changed. We examine how the mean outcomes would be pre-
293 dicted to change if the inputs of the specific variable of interest are modified,
294 such as reducing a variable in some equivalent way across counties. Other modi-
295 fications of the inputs could be used to examine these variable importance plots,
296 but we looked at % changes in the variable across counties. Using these models
297 we then make marginal predictions on the predicted number of COVID-19 cases
298 and all-cause mortalities when increasing (or decreasing) the top variables found
299 by the variable importance procedure.

300 Suppose a particular observation $O_i = (W_i, A_i, Y_i)$ in county i , $i = 1, \dots, n$,
 301 depends on explanatory variable A_i , other adjustment covariates W_i , and out-
 302 come Y_i , all in county i . If we wish to generate an estimate that characterizes
 303 the association of Y across all counties with A adjusting for W , but does not
 304 rely on a linear approximation, then we do this through plotting the estimate

$$\phi(\pi) = E_{A,W}\{E(Y|A = \max(\tilde{A}, (1 - \pi) * A), W)\} \quad (1)$$

305 as a function of π , where \tilde{A} is the minimum observed value of A across all
 306 counties in the data and π is interpreted as the proportional reduction in the
 307 county-specific value of the variable. To avoid extrapolation, we truncate A at
 308 the minimum observed value for the variable among counties. In essence, we
 309 exam the resulting predicted mean outcome across all counties as though the
 310 particular variable were reduced by $\pi\%$ in all counties. This plot then provides
 311 a relevant function of the importance of the variable to the outcome.

312 Under several strong assumptions, including that the other covariates (the
 313 W , being either all other predictors or all but the ones in same sub-category)
 314 contain all the confounding information, sufficient experimentation (no positiv-
 315 ity violations [[36, 35?]]), and independence of outcomes across counties, one
 316 could interpret (1) as identifying the marginal mean had, contrary to fact, all
 317 counties been set at the stochastic value, $\max(\tilde{A}, (1 - \pi) * A)$. Of course, we
 318 have time-ordering among the covariates that we can assert, so we treat these
 319 plots as a nonparametric form of association measure.

320 To derive inference, we use the nonparametric bootstrap (randomly sam-
 321 pling counties with replacement) for each π reduced mean. The procedure is
 322 as follows: (i) fit each outcome using the aforementioned set of learners, (ii)
 323 iteratively remove variables and determine risk-ratios, (iii) set the variable with
 324 the highest risk-ratio as the target variable for marginal predictions, (iv) for
 325 each percentage from 0-1.0 at 0.10 intervals (a) resample the county data with
 326 replacement, (b) refit the Super Learner with this resampled data, (c) reduce
 327 or increase the target variable for the respective percentage, and (d) predict
 328 the expected number of cases or mortalities with this new fit. Here, in step

329 (iv) we bootstrap this procedure 1000 times by resampling, refitting, and mak-
330 ing marginal predictions, in order to create confidence intervals (CI) at each
331 percent change to the target variable.

332 Note that we use three different models to estimate the regression $E(Y|A, W)$
333 depicted in (1): fully adjusted, adjusted only variables not in the sub-category
334 of the variable of interest, and unadjusted. To evaluate the performance of our
335 model, we compare each marginal prediction to that of a univariate model of
336 the target variable (found by variable importance) for each outcome. Here, a
337 general additive model was retrained at each iteration of the bootstrap for each
338 reduction in the target variable and cases/mortalities were predicted through
339 this univariate model. Likewise, to investigate possible over-corrections of our
340 Super Learners, we also remove similar variables that may have strong multi-
341 collinearity with the target variable. This was done by removing variables in
342 the target variable sub-category and training a Super Learner on this set of
343 covariates through the bootstrap. Results are presented as line plots and the
344 actual observed average or sum for each outcome are plotted as horizontal lines
345 for comparing actual outcomes to model predictions.

346 All data aggregation, curation, cleaning, exploratory analysis, ML pipeline,
347 and marginal predictions were performed in R [37]. All coding scripts are avail-
348 able on our GitHub page for open-access and use: ([https://github.com/blind-
349 contours/Getz_Hubbard_Covid_Ensemble_ML_Public.git](https://github.com/blind-contours/Getz_Hubbard_Covid_Ensemble_ML_Public.git)).

350 **3. Results**

351 *COVID-19 Outcomes and County Feature Distributions.* There are 3,142 coun-
352 ties in the U.S. After our data cleaning and curation process, there were 2,620
353 counties included in the analysis and 101 county-level features. As such, our
354 analysis covers 83% of the U.S. population as represented by counties. In these
355 counties, as of July 15, 2020 there were 243,065 cases at day 25, 2,531,134 cases
356 to-date, 53,018 all-cause deaths at day 100, 111,991 all-cause deaths to-date,
357 and the average number of days to the first case in a county relative to the first

358 case in Snohomish County, Washington was 68 days. A description of the vari-
359 ables used and their sources are given in supplementary materials (**Table S1**).
360 A breakdown of these features with the respective mean, standard deviation,
361 and range of values are also given in the supplementary materials (**Table S2**).

362 *Exploratory.* The heatmap for exploring the patterns in these data are given
363 in **Figure 1**. The marked section of the heatmap show an outcome trend for
364 1: First quantile of day of first case (Q1 = earlier days of first case): 2. high-
365 est total deaths to-date (Q4): 3. highest deaths at day 100 (Q3): 4. highest
366 total COVID-19 cases to-date (Q4): and 4. highest COVID-19 cases at day
367 25 (Q4). The distribution of the outcome quantiles across the county dendro-
368 gram groups are provided in the supplementary materials (**Table S8**). The
369 cluster of counties with most severe outcomes is marked on the left in **Figure**
370 **1**. These patterns indicate there are counties that cluster together based on
371 similar characteristics and these counties correspond with an earlier first case
372 in the county and higher COVID-19 case and mortality rates. For a breakdown
373 of the number of counties in each state in this cluster see the supplementary
374 material; briefly, however, the states with the highest number of counties in
375 this cluster with highest outcomes are 1. Virginia (36), 2. Florida (33), and 3.
376 Texas (30). The highest column values (red and orange pixels in the heat map)
377 in this highest county row cluster occur occur in branches 3-16 and branches
378 46-53 of the column-wise dendrogram. A full list for each cluster is given in
379 the plot but the highest county features for this subset were: 1. obesity, 2.
380 sexually transmitted diseases, 3. income inequality, 4. food environment index,
381 5. CDC limited English scores, 6. latitude, 7. poverty income ratio, 8. GDP,
382 9. preventable hospital stays, 10. arthritis, 11. asthma, and 12. ischemic heart
383 disease.

384 *Super Learner.* As discussed, we use cross-validation to generate a coefficient
385 that defines the weight for a respective learner in the ensemble. This procedure
386 is done for each outcome and the same learners are used for each outcome
387 (outside of day of first case where Poisson outcomes were defined). **Tables**

388 **S3-S7** in supplementary give a detailed breakdown of how each algorithm was
389 used in the Super Learner, the risk of the respective algorithm and the overall
390 risk of the Super Learner. Overall, our Super Learners were able to achieve
391 good fits by utilizing multiple algorithms. Results show that for each outcome
392 the Super Learners were largely and consistently built from multiple elastic net
393 models, multiple xgboost models, and random forest. **Table 2** shows resulting
394 risk for each Super Learner for each outcome. Because the learners were fit to
395 the per-capita standardized outcome data, we multiply each risk by the total
396 population in the dataset to get absolute error based on total numbers of cases
397 or mortalities. We also calculate the r-squared for each Super Learner to show
398 variance explained by each model.

399 *Variable Importance.* The top variable categories for each outcome were: 1) Day
400 of first case in a county: demography; 2) COVID-19 cases at day 25: ethnicity,
401 transit and preventable disease; 3) total COVID-19 cases to-date: ethnicity and
402 preventable disease; 4) mortalities at day 100: transit and ethnicity; 5) total
403 mortalities to-date: ethnicity and transit. The top individual variables across
404 all the models were: overall population of a county, CDC vulnerability scores
405 for minority and limited English, public transportation use, and proportion of
406 Black/African American individuals in a county. To visualize results, we present
407 the risk-ratio (RR) results for each COVID-19 outcome collectively in **Figure**
408 **2** and **Figure 3** as a series of dot-plots (RR threshold at 1.01). Based on these
409 figures, it can be seen that for day of the first case in a county the total pop-
410 ulation (RR: 1.38) was the most important variable. For per-capita COVID-
411 19 cases at day 25, the top variable is the CDC minority score (RR: 1.04).
412 For per-capita COVID-19 cases to-date, the CDC's score for limited English
413 speaking (1.40) the CDC's score for minority populations (1.17) and proportion
414 Black/African-American individuals in a county (RR: 1.12) were the top vari-
415 ables. For per-capita all-cause mortalities at day 100, proportion taking public
416 transportation (RR: 1.14) and proportion Black/African-American individuals
417 (RR: 1.07) were the top variables. For per-capita all-cause mortalities to-date,

418 proportion Black/African-American individuals (RR: 1.08) and proportion tak-
419 ing public transportation (RR: 1.03) were the top county features.

420 *Marginal Prediction Results.* For the day of first case, population size was
421 clearly most important predictor. For examining the association of cases and
422 deaths of COVID-19 we choose two outcomes (total deaths and cases by July
423 14, 2020) and three of the most consistently important variables: two related to
424 demographic features of the population (CDC Minority Score and Proportion of
425 Black/African-American individuals) and one related to transportation (metric
426 of public transportation use). We estimate the relationship of proportional re-
427 ductions in each of the predictor variables on the marginal outcome (using the
428 substitution estimator of (1)) based upon a machine learning fit when control-
429 ling for: 1) all other variables (including possibly strongly collinear variables
430 within the same sub-category), 2) only variables outside the target variable
431 sub-category, and 3) nothing (unadjusted). The latter estimator of $E(Y|A)$ is
432 based upon a smooth regression of the outcome versus the continuous covariates,
433 specifically using general additive model with identify link (GAM;[38]).

434 **Figure 4** shows the predicted average day of first case across all counties
435 for each proportion of population size reduced across the counties. Generally,
436 smaller county populations are predicted to have a delay relative to counties
437 with higher population sizes, *all other factors in the model staying constant*
438 (the current average being 68 days after the initial U.S. case in Washington) for
439 all three estimators, with somewhat larger effects in the adjusted models. For
440 all models, a 0.50 proportional reduction in population size across the counties
441 suggest a delay of around 2 to 3 days.

442 **Figure 5** shows plotted results for two outcomes (total counts and deaths),
443 both versus proportion reductions in the three predictor variables. Proportional
444 decreases in CDC Minority Score is associated with a decrease in COVID-19
445 cases and death. The large attenuation of the relationship in the adjusted
446 models suggest strong confounding by the other covariates, where in the fully
447 (perhaps over) curve approaches the null line. Using the curves not adjusting for

448 other variables in the sub-category, the association suggest a 0.5 proportional
449 reduction in the score would predict a reduction of 750,000 cases (out of around
450 2,400,000 total number by July 14) and approximately 10,000 deaths (out of
451 around 115,000). Note that, particularly with deaths, the unadjusted curve is
452 quite different from the actual number (blue line is substantially below the hor-
453 izontal black line at the point of no intervention), this should be equal to that
454 value if the model fits the data well. This is due to a few counties with extreme
455 large counts (of both death and cases), which are poorly predicted by the bi-
456 variate smooths resulting in large positive residuals. Thus, when one estimates
457 counts based upon reductions in the variable of interest (CDC Minority Score in
458 this case), you get a prediction of the count that underestimates the true count.
459 Note, by substantially reducing residual variation, the adjusted curves tend to
460 be much closer to the observed count at 0 reduction.

461 Reducing public transportation suggests significant reduction in deaths, but
462 little impact on case counts. For cases, we again have a poor fit of the bivariate
463 smooth (unadjusted), along with the suggestion of significant confounding by
464 other covariates. For deaths, a reduction of 50% suggest a reduction in deaths
465 of 10,000, but the fact that the intercept for both adjusted curves (and unad-
466 justed) is less than the observed count suggest again the influence of outliers.
467 Bootstrapped linear regression of the marginal predictions showed that for a
468 10% reduction in public transit use, total deaths reduce by 2012 (95% CI [1972,
469 2356]).

470 Finally, for reductions in the proportion of Black/African-American individ-
471 uals, there appears to be quite different estimates between the adjusted curves
472 of COVID-19 counts, suggesting that variables in the sub-category of this vari-
473 able create the possibility of over-adjustment in the full model. For deaths, the
474 curves are nearly identical and imply that a reduction of the disparity between
475 Blacks and Whites in 50% of the population (one way to interpret an actual
476 reduction in this variable) suggest a reduction of about 9,000 total deaths. Like-
477 wise bootstrapped linear regression of these predictions show a 10% increase in
478 the proportion of Black/African-American individuals in a county increases to-

479 tal deaths to date by 2067 (95% CI [1189, 2654]).

480 **4. Discussion**

481 In this paper, we took a semi-parametric (machine learning) approach to
482 evaluating a wide range of county-level features which may impact the spread,
483 number of cases, and deaths of COVID-19 in the U.S. Our contributions are the
484 following: (i) curating an open-source data repository that includes variables
485 from many sources, categorized into sub-groups and filtered such that strongly
486 collinear variables are removed for statistical analysis; (ii) demonstrating the
487 use of these variables in ensemble machine learning to build 5 Super Learners
488 for each COVID-19 outcome measure; (iii) evaluating the features to identify
489 the top variables that influence each outcome; (iv) adjusting the top variables
490 in our model to make marginal predictions for each COVID-19 outcome, while
491 controlling for all other factors to establish the strength and directionality of
492 the relationships; (v) constructing confidence intervals around all effects via
493 bootstrapping to evaluate significant trends from baseline and between model-
494 ing approaches. Overall, using 101 county-level features our models show very
495 good fits to the outcomes (all observed outcomes were within model confidence
496 intervals apart from mortality which were slightly outside our CIs at baseline).

497 These fits establish that our models are able to accurately predict each out-
498 come given the county-level feature variables. Our variable importance measures
499 for each model fit generally show a trend that the total population drives day
500 of first case in a county and the proportion of Black/African-American indi-
501 viduals in a county and CDC minority scores are most important independent
502 contributions of COVID-19 cases and deaths as of mid July.

503 Causal inference pertaining to the individual relationships of these variables
504 on each outcome is speculative at best given that these study variables are
505 ecological and also are a static snap-shot of county variables collected before
506 the pandemic hit the U.S. However, the general trend in these results seem to
507 represent what is anecdotally or univariately reported as the U.S. faces this

508 continuing pandemic. That is, for day of first case in a county, the total popula-
509 tion as the most important variable makes sense given the larger the population
510 the higher the probability of someone being infected traveling to the respective
511 county. Likewise, CDC minority scores and Black/African-American individuals
512 are correlated with reports that suggest that minority populations and People
513 of Color are disproportionately impacted by COVID-19 [13] [15]. In addition,
514 we also show a significant potential impact of baseline public transportation use
515 and mortality [39][40] [27]. This could indicate that there is higher probability
516 of exposure in counties where travel on public transportation leads longer and
517 closer duration contains. This may also lead to higher infectious doses that may
518 possibly increase severity of infection and consequent mortality, a phenomenon
519 thought to be the case for influenza [41].

520 Our finding that counties with larger CDC minority population measures
521 have higher COVID-19 outcomes, even when controlling for one-hundred other
522 county-level variables (Table 1 and S1), show the value of such measures when
523 trying to determine the impact of risk level based on social factors for those
524 disproportionately impacted by COVID-19. Additional measures, however, are
525 necessary to understand the reasons for this. For instance, although our models
526 adjust for income, access to health, and occupation types, our data are limited
527 to reported factors that may not account for systematic or institutional levels of
528 cultural/societal factors placing individuals at risk. Such factors may confound
529 or modify others for which have been adjusted and may place certain individuals
530 at greater risk for SARS-CoV-2 infection. The main difference between our
531 findings and those reported to date is that our analysis controls for many other
532 possible mediating factors (e.g., access to health-care, smoking, diabetes, heart-
533 disease, and food security).

534 **5. Conclusion**

535 The goal of our study is to identify the most salient factors that put pop-
536 ulations at risk for COVID-19, thereby providing some guidance to individu-

537 als making difficult policy decisions at this critical time to quell the evolving
538 pandemic. Specifically, racial composition of counties and intensity of public
539 transportation use therein seem to be the most important risks factors for both
540 the initial rapid growth and subsequent high incidence, and also help explain
541 variations in mortality rates across counties. More work, however, is needed to
542 establish causal rather than purely statistical relationships. Future work with
543 detailed individual data will be important for getting more robust estimates of
544 the individual impact of the factors examined. Whether causal or statistical,
545 these results should be taken into account when developing policies for lifting
546 restrictions. Additionally, as efforts continue to disseminate services and fund-
547 ing, and to roll out vaccination programs, once effective vaccines have been
548 developed, consideration of these factors will facilitate the efficacious allocation
549 of resources to the benefit of the US population as a whole.

Source	N Var.	Var. Examples
USAFacts	6	COVID-19 Outcome Data, Population
Bureau of Economic Analysis (BEA)	1	GDP
5-Year American Community Survey (ACS), 2014 - 2018	14	County percentages by Sex and Ethnicity, Employment, Household Income, use of Public Transportation
TIGER/Line Geodatabases	7	Latitude, Longitude, Land Area
TIGER/Line Geodatabases; Federal Aviation Administration (FAA)		Distance to Airports
Interactive Atlas of Heart Disease and Stroke (2014-2016)	4	Number of Hospitals, Stroke, Access to Parks
County Health Rankings & Roadmaps	21	Life Expectancy, Smoking, Obesity,, Food Access, Mental Health, Physicians, Household Overcrowding etc.
Centers for Medicare & Medicaid Services (CMS)	15	Drugs Abuse, Hypertension, Hyperlipidemia, Osteoporosis, etc.
National Centers for Environmental Information	1	Precipitation
CDC's Social Vulnerability Index (SVI)	11	Percentile over 65 or under 17, Minority Scores, Limited English, Low Income Housing Estimates, Number Institutionalized
Quarterly Census of Employment and Wages	14	Labor force types, farming/mining, private industry, education/healthcare etc.
MIT election lab	1	Calculated Proportion Voted Republican 2016
Google	6	Google mobility to location type, Residence, Grocery etc.

Table 1: Number of variables used from respective sources with some examples given, complete list with distributions given in supplementary material

COVID-19 Outcome	Model Risk (per capita)	Model Risk (counts)	R-squared
Day of First Case	NA	159.58	0.75
COVID-19 Cases at Day 25	4.22 e-05	10539.64	0.59
Total COVID-19 Cases to-date	5.22 e-05	13053.75	0.87
All-Cause Death at Day 100	2.80 e-08	7.00	0.57
All-Cause Death at to-date	1.42 e-07	35.52	0.59

Table 2: Cross Validated Super Learner Risk Across COVID-19 Outcomes

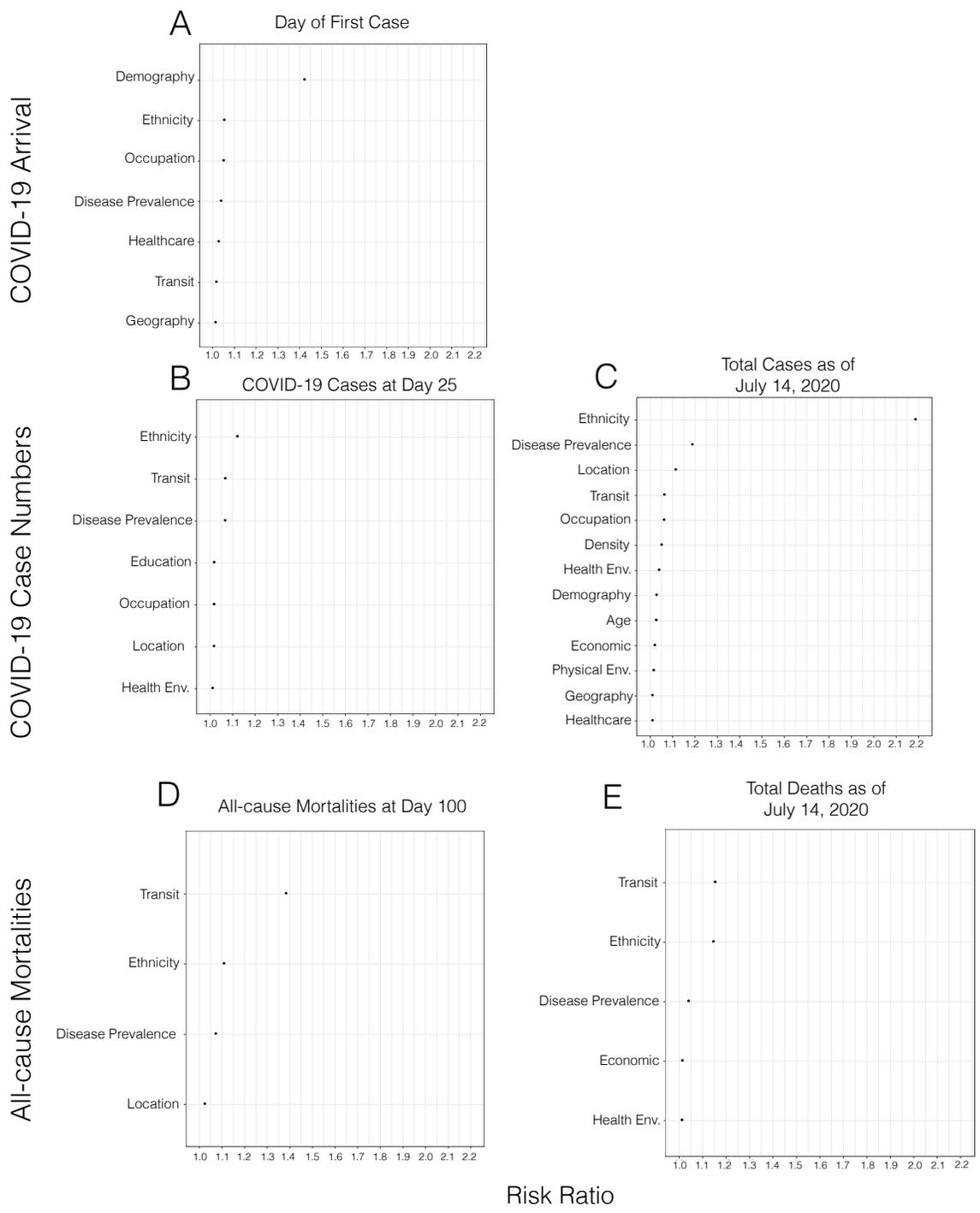


Figure 2: Variable importance as indicated by the relative increase of mean-squared error when the block of variables is removed

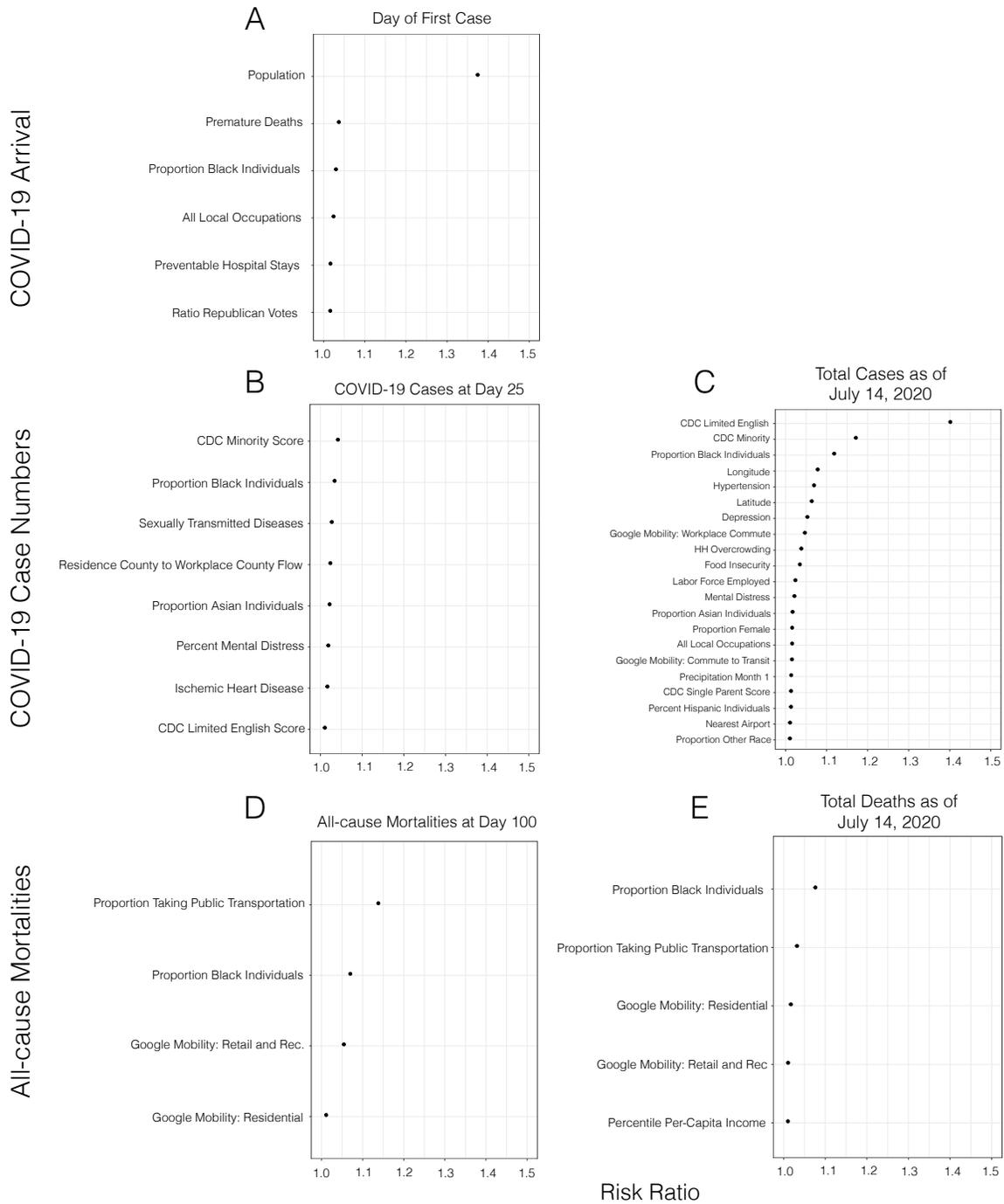


Figure 3: Variable importance as indicated by the relative increase of mean-squared error when a single variable is removed

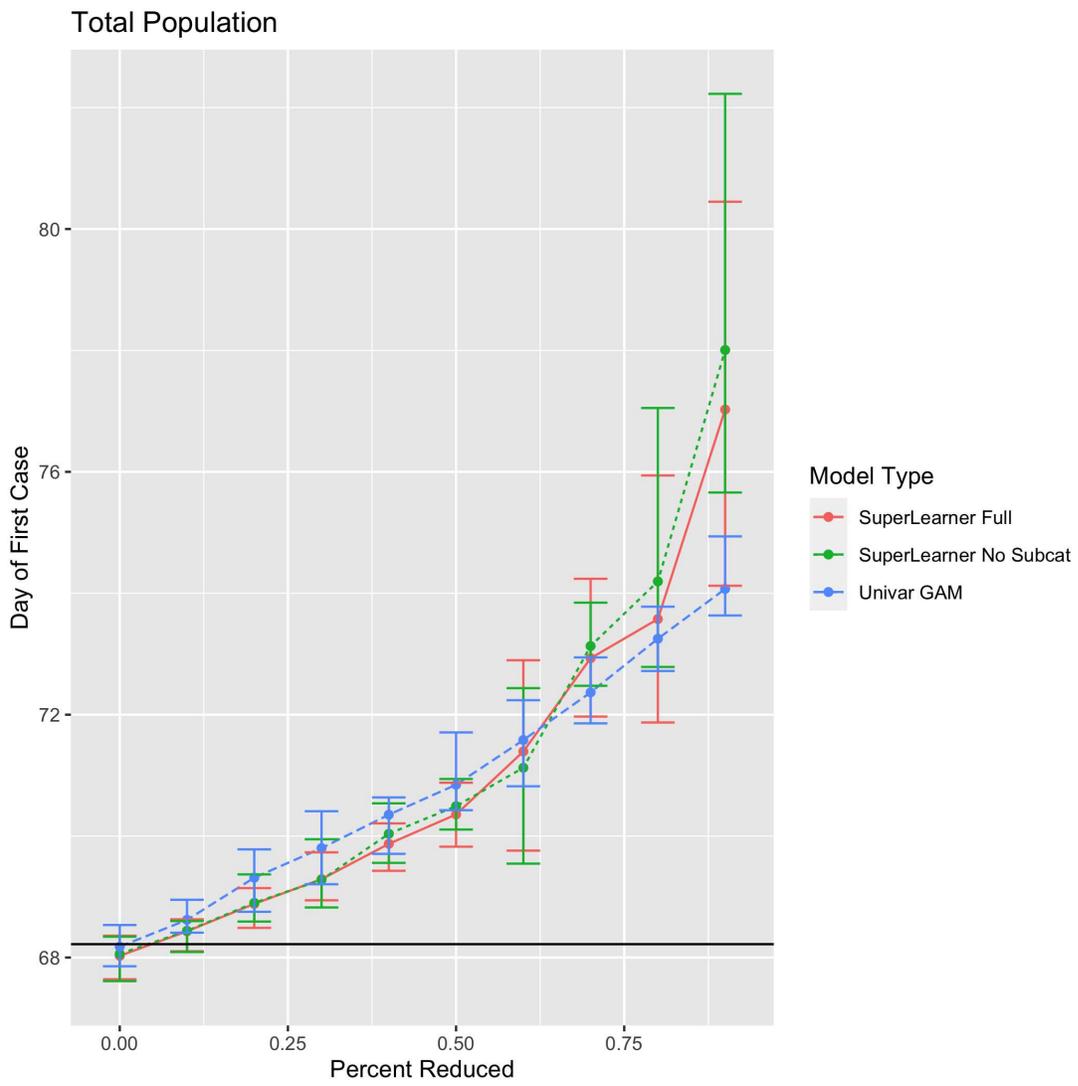


Figure 4: Marginal predictions of day of first case (relative to index time) for different proportional reductions of total population size for models adjusting for all other covariates, only covariates not in sub-category (see supplement table 1) and unadjusted.

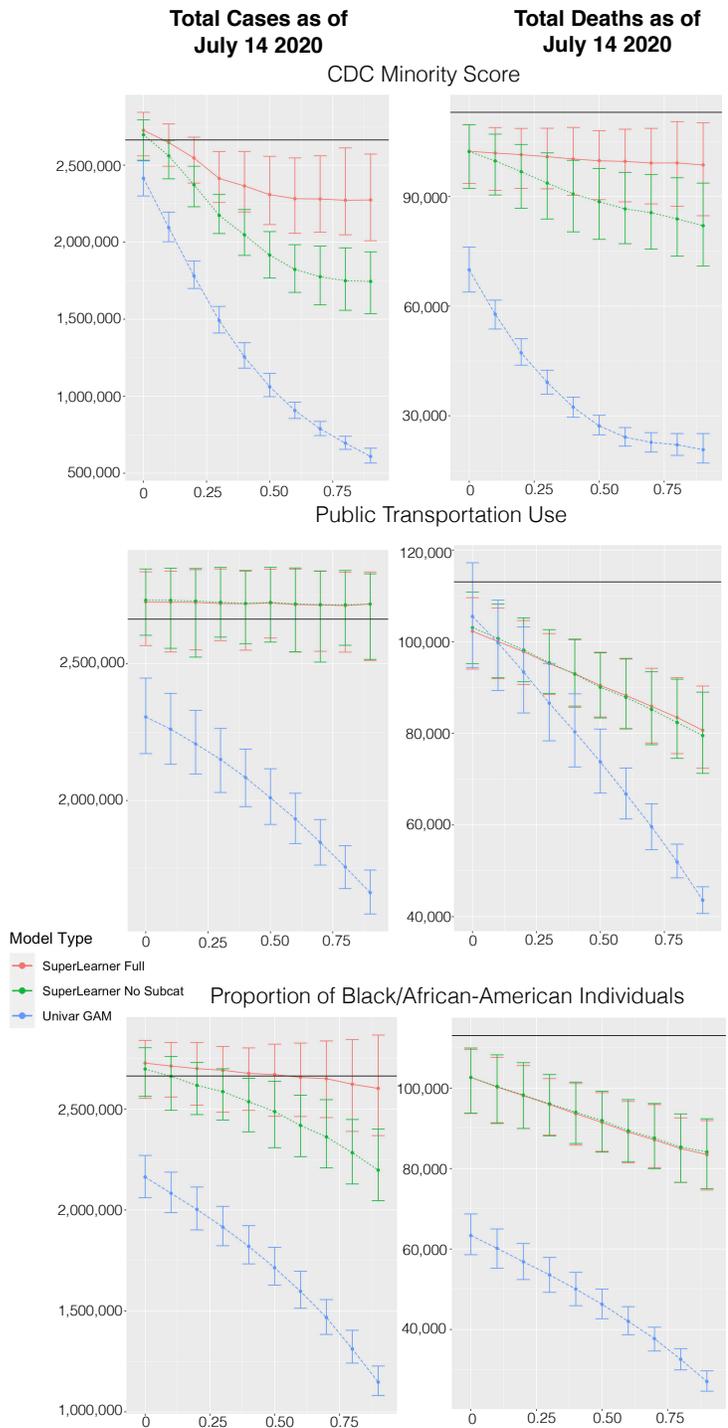


Figure 5: Marginal predictions of total cases and deaths by July 14, 2020) for three of the most consistently important variables in predicting the count outcomes: CDC minority score, proportion of Black individuals and a metric of public transportation use. X-axis is different proportional reductions of each of the three predictors, the Y-axis is the marginal predicted total counts for models adjusting for all other covariates, only covariates not in sub-category (see supplement table 1) and unadjusted.

552 **References**

- 553 [1] N. Zhu, D. Zhang, W. Wang, X. Li, B. Yang, J. Song, X. Zhao, B. Huang,
554 W. Shi, R. Lu, et al., A novel coronavirus from patients with pneumonia
555 in china, 2019, *New England Journal of Medicine*.
- 556 [2] F. Wu, S. Zhao, B. Yu, Y.-M. Chen, W. Wang, Z.-G. Song, Y. Hu, Z.-
557 W. Tao, J.-H. Tian, Y.-Y. Pei, et al., A new coronavirus associated with
558 human respiratory disease in china, *Nature* 579 (7798) (2020) 265–269.
- 559 [3] H. A. Rothan, S. N. Byrareddy, The epidemiology and pathogenesis of
560 coronavirus disease (covid-19) outbreak, *Journal of autoimmunity* (2020)
561 102433.
- 562 [4] Z.-W. Ye, S. Yuan, K.-S. Yuen, S.-Y. Fung, C.-P. Chan, D.-Y. Jin, Zoonotic
563 origins of human coronaviruses, *International journal of biological sciences*
564 16 (10) (2020) 1686.
- 565 [5] P. Zhou, X.-L. Yang, X.-G. Wang, B. Hu, L. Zhang, W. Zhang, H.-R. Si,
566 Y. Zhu, B. Li, C.-L. Huang, et al., A pneumonia outbreak associated with a
567 new coronavirus of probable bat origin, *nature* 579 (7798) (2020) 270–273.
- 568 [6] R. Lu, X. Zhao, J. Li, P. Niu, B. Yang, H. Wu, W. Wang, H. Song,
569 B. Huang, N. Zhu, et al., Genomic characterisation and epidemiology of
570 2019 novel coronavirus: implications for virus origins and receptor binding,
571 *The Lancet* 395 (10224) (2020) 565–574.
- 572 [7] WHO, Who timeline - covid-19, [http://precog.iiitd.edu.in/people/
573 anupama](http://precog.iiitd.edu.in/people/anupama) (April 2020).
- 574 [8] CDC, Coronavirus (covid-19), [https://www.cdc.gov/coronavirus/
575 2019-ncov/index.html](https://www.cdc.gov/coronavirus/2019-ncov/index.html) (2020).
- 576 [9] M. L. Holshue, C. DeBolt, S. Lindquist, K. H. Lofy, J. Wiesman, H. Bruce,
577 C. Spitters, K. Ericson, S. Wilkerson, A. Tural, et al., First case of 2019
578 novel coronavirus in the united states, *New England Journal of Medicine*.

- 579 [10] I. Ghinai, T. D. McPherson, J. C. Hunter, H. L. Kirking, D. Christiansen,
580 K. Joshi, R. Rubin, S. Morales-Estrada, S. R. Black, M. Pacilli, et al., First
581 known person-to-person transmission of severe acute respiratory syndrome
582 coronavirus 2 (sars-cov-2) in the usa, *The Lancet*.
- 583 [11] R. M. Burke, Active monitoring of persons exposed to patients with con-
584 firmed covid-19—united states, january–february 2020, *MMWR. Morbidity*
585 *and mortality weekly report* 69.
- 586 [12] C. W. Yancy, Covid-19 and african americans, *JAMA*. Published online.
- 587 [13] S. S. Coughlin, J. X. Moore, V. George, J. A. Johnson, J. Hobbs, Covid-19
588 among african americans: From preliminary epidemiological surveillance
589 data to public health action, *American Journal of Public Health* (2020)
590 1157–1159.
- 591 [14] S. Kodidela, A. Kumar, K. Gerth, S. Kumar, C. Walker, Lessons learned
592 from health disparities among african ameri-cans in the hiv epidemic: What
593 to expect for covid-19 and potential approaches to mitigate health disparity,
594 *Emerg Infect Dis Diag J: EIDDJ-100021* 2 (03).
- 595 [15] S. Bandi, M. Z. Nevid, M. Mahdavinia, African american children are at
596 higher risk for covid-19 infection, *Pediatric Allergy and Immunology*.
- 597 [16] D. B. G. Tai, A. Shah, C. A. Doubeni, I. G. Sia, M. L. Wieland, The
598 disproportionate impact of covid-19 on racial and ethnic minorities in the
599 united states, *Clinical Infectious Diseases*.
- 600 [17] E. E. Wiemers, S. Abrahams, M. AlFakhri, V. J. Hotz, R. F. Schoeni, J. A.
601 Seltzer, Disparities in vulnerability to severe complications from covid-19
602 in the united states, *Tech. rep.*, National Bureau of Economic Research
603 (2020).
- 604 [18] N. N. Pettit, E. L. MacKenzie, J. Ridgway, K. Pursell, D. Ash, B. Patel,
605 M. T. Pho, Obesity is associated with increased risk for mortality among
606 hospitalized patients with covid-19, *Obesity*.

- 607 [19] J. Lighter, M. Phillips, S. Hochman, S. Sterling, D. Johnson, F. Francois,
608 A. Stachel, Obesity in patients younger than 60 years is a risk factor for
609 covid-19 hospital admission, *Clinical Infectious Diseases*.
- 610 [20] C. Covid, C. COVID, C. COVID, N. Chow, K. Fleming-Dutra, R. Gierke,
611 A. Hall, M. Hughes, T. Pilishvili, M. Ritchey, et al., Preliminary estimates
612 of the prevalence of selected underlying health conditions among patients
613 with coronavirus disease 2019—united states, february 12–march 28, 2020,
614 *Morbidity and Mortality Weekly Report* 69 (13) (2020) 382.
- 615 [21] M. Lipsitch, D. L. Swerdlow, L. Finelli, Defining the epidemiology of covid-
616 19—studies needed, *New England journal of medicine* 382 (13) (2020) 1194–
617 1196.
- 618 [22] M. G. Baker, T. K. Peckham, N. S. Seixas, Estimating the burden of united
619 states workers exposed to infection or disease: a key factor in containing
620 risk of covid-19 infection, *PLoS One* 15 (4) (2020) e0232452.
- 621 [23] J. O. Lloyd-Smith, S. J. Schreiber, P. E. Kopp, W. M. Getz, Superspread-
622 ing and the effect of individual variation on disease emergence, *Nature*
623 438 (7066) (2005) 355.
- 624 [24] M. M. Sajadi, P. Habibzadeh, A. Vintzileos, S. Shokouhi, F. Miralles-
625 Wilhelm, A. Amoroso, Temperature and latitude analysis to predict po-
626 tential spread and seasonality for covid-19, Available at SSRN 3550308.
- 627 [25] N. Altieri, R. Barter, J. Duncan, R. Dwivedi, K. Kumbier, X. Li, R. Net-
628 zorg, B. Park, C. Singh, Y. S. Tan, et al., Curating a covid-19 data repos-
629 itory and forecasting county-level death counts in the united states, arXiv
630 preprint arXiv:2005.07882.
- 631 [26] J. O. Ferstad, A. J. Gu, R. Y. Lee, I. Thapa, A. Y. Shin, J. A. Salomon,
632 P. Glynn, N. H. Shah, A. Milstein, K. Schulman, et al., A model to forecast
633 regional demand for covid-19 related hospital beds, medRxiv.

- 634 [27] K. Desmet, R. Wacziarg, Understanding spatial variation in covid-19 across
635 the united states, Tech. rep., National Bureau of Economic Research (2020).
- 636 [28] N. Altieri, R. L. Barter, J. Duncan, R. Dwivedi, K. Kumbier, X. Li, R. Net-
637 zorg, B. Park, C. Singh, Y. S. Tan, T. Tang, Y. Wang, B. Yu, Curating
638 a COVID-19 data repository and forecasting county-level death counts in
639 the United States, arXiv:2005.07882.
640 URL <http://arxiv.org/abs/2005.07882>
- 641 [29] B. D. Killeen, J. Y. Wu, K. Shah, A. Zapaishchykova, P. Nikutta,
642 A. Tamhane, S. Chakraborty, J. Wei, T. Gao, M. Thies, M. Unberath, A
643 County-level Dataset for Informing the United States' Response to COVID-
644 19, arXiv:2004.00756.
645 URL <http://arxiv.org/abs/2004.00756>
- 646 [30] E. C. Polley, M. J. van der Laan, Super Learner in Prediction, U.C. Berkeley
647 Division of Biostatistics Working Paper (2010) 1–19.
648 URL <http://biostats.bepress.com/ucbbiostat/paper266/>
- 649 [31] M. J. Van Der Laan, E. C. Polley, A. E. Hubbard, Super learner, Sta-
650 tistical Applications in Genetics and Molecular Biology 6 (1). doi:
651 10.2202/1544-6115.1309.
- 652 [32] E. LeDell, M. J. Van Der Laan, M. Peterson, AUC-Maximizing Ensembles
653 through Metalearning, International Journal of Biostatistics 12 (1) (2016)
654 203–218. doi:10.1515/ijb-2015-0035.
- 655 [33] S. Dudoit, M. J. V. D. Laan, Asymptotics of Cross-Validated Risk Estima-
656 tion in Estimator Selection and Performance Assessment Asymptotics of
657 Cross-Validated Risk Estimation in Estimator Selection and Performance
658 Assessment, U.C. Berkeley Division of Biostatistics Working Paper Series
659 Year.
- 660 [34] M. A. Hernán, J. M. Robins, Estimating causal effects from epidemiological

- 661 data, *Journal of Epidemiology and Community Health* 60 (7) (2006) 578–
662 586. doi:10.1136/jech.2004.029496.
- 663 [35] J. M. Robins, Addendum To 'a New Approach To Causal Inference in Mor-
664 tality Studies With a Sustained Exposure Period - Application To Con-
665 trol of the Healthy Worker Survivor Effect', *Computers & mathematics*
666 with applications 14 (9-12) (1987) 923–945. doi:10.1016/0898-1221(87)
667 90238-0.
- 668 [36] Maya L Petersen, Kristin E Porter, Susan Gruber, Yue Wang, M. J.
669 van der Laan, Diagnosing and responding to violations in the positivity
670 assumption, *Stat Methods Med Res* 21 (1) (2012) 1–13. doi:10.4018/
671 978-1-4666-2661-4.ch001.
- 672 [37] R Core Team, *R: A Language and Environment for Statistical Computing*,
673 R Foundation for Statistical Computing, Vienna, Austria (2013).
674 URL <http://www.R-project.org/>
- 675 [38] T. J. Hastie, R. J. Tibshirani, *Generalized additive models*, Chapman and
676 Hall.
- 677 [39] R. S. Loomba, G. Aggarwal, S. Aggarwal, S. Flores, E. G. Villarreal, J. S.
678 Farias, C. J. Lavie, Disparities in case frequency and mortality of coron-
679 avirus disease 2019 (covid-19) among various states in the united states,
680 medRxiv.
- 681 [40] M. W. Tenforde, E. B. Rose, C. J. Lindsell, N. I. Shapiro, D. C. Files,
682 K. W. Gibbs, M. E. Prekker, J. S. Steingrub, H. A. Smithline, M. N. Gong,
683 et al., Characteristics of adult outpatients and inpatients with covid-19—11
684 academic medical centers, united states, march–may 2020, *Morbidity and*
685 *Mortality Weekly Report* 69 (26) (2020) 841.
- 686 [41] A. C. Paulo, M. Correia-Neves, T. Domingos, A. G. Murta, J. Pedrosa,
687 Influenza infectious dose may explain the high mortality of the second and
688 third wave of 1918–1919 influenza pandemic, *PLoS One* 5 (7) (2010) e11655.

Figures

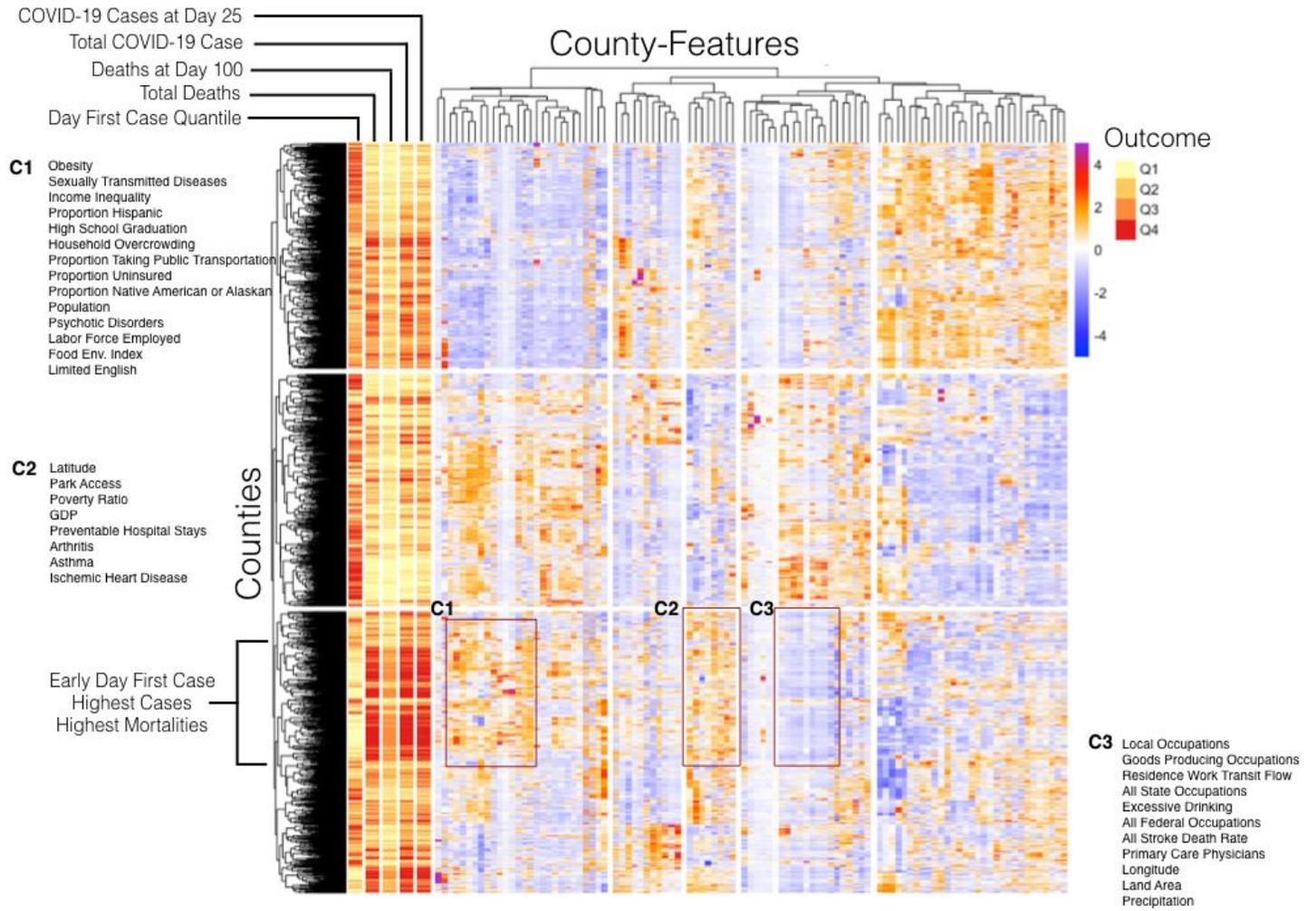


Figure 1

COVID-19 heatmap visualization of the distribution of county-level data.

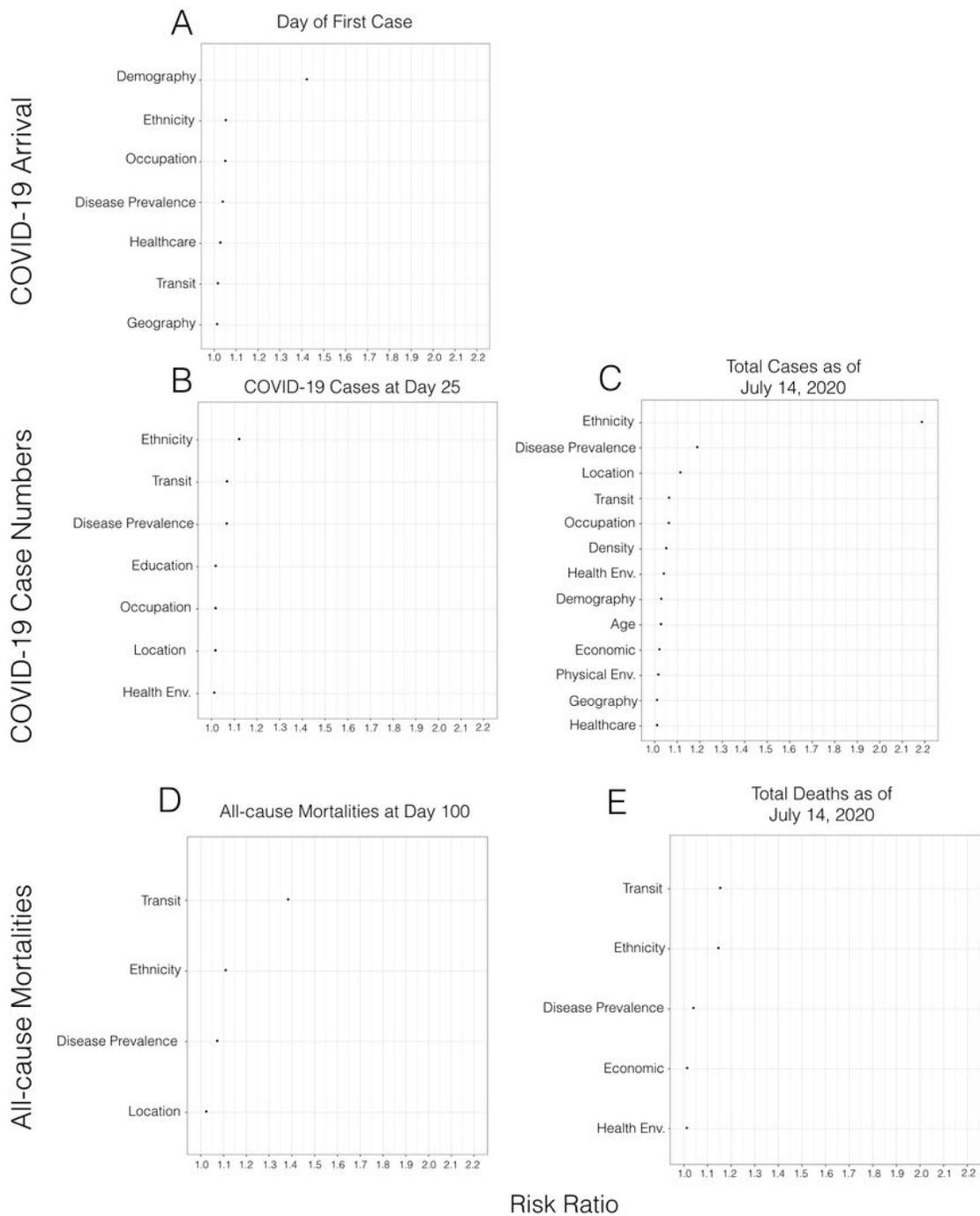


Figure 2

Variable importance as indicated by the relative increase of mean-squared error when the block of variables is removed.

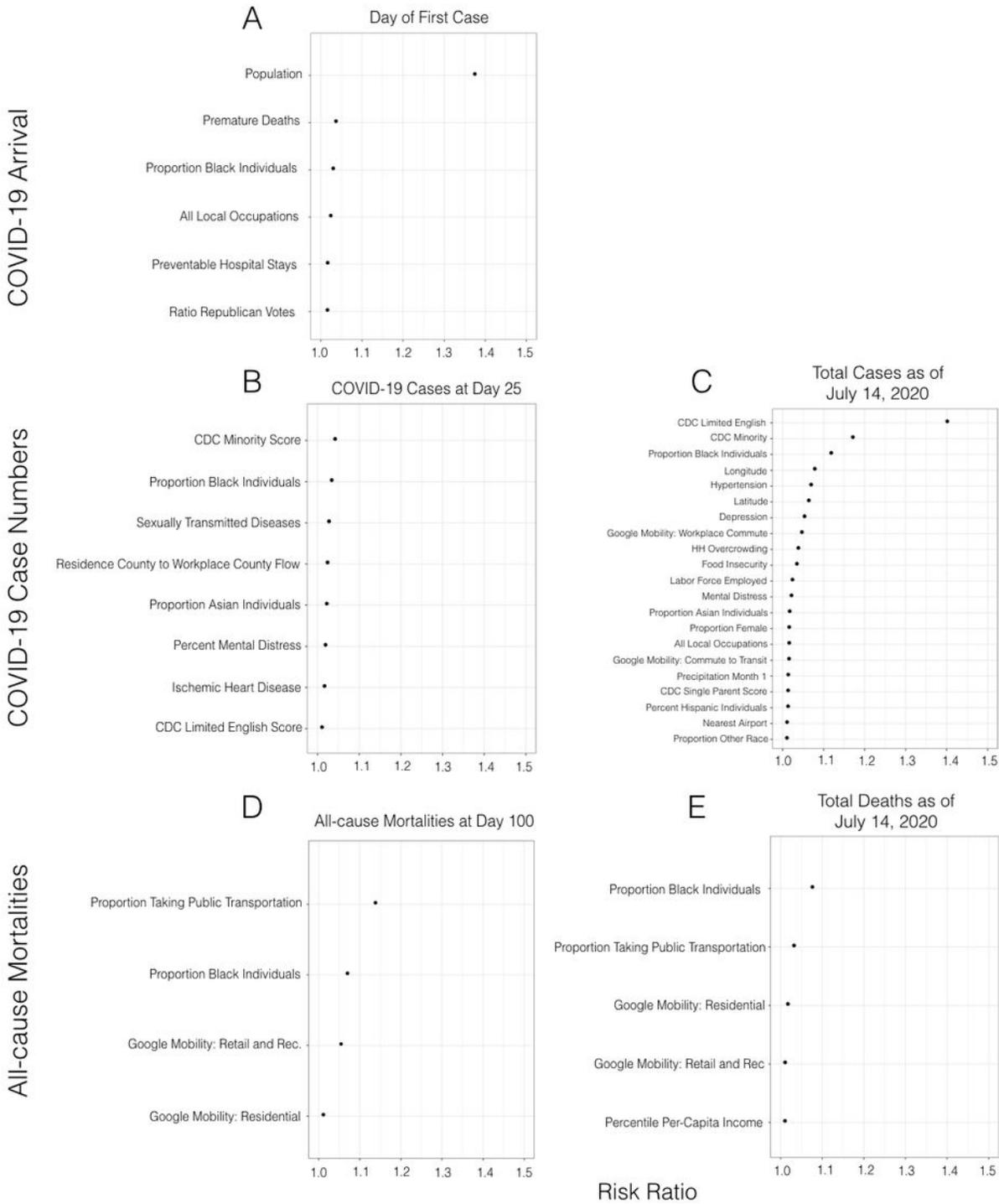


Figure 3

Variable importance as indicated by the relative increase of mean-squared error when a single variable is removed.

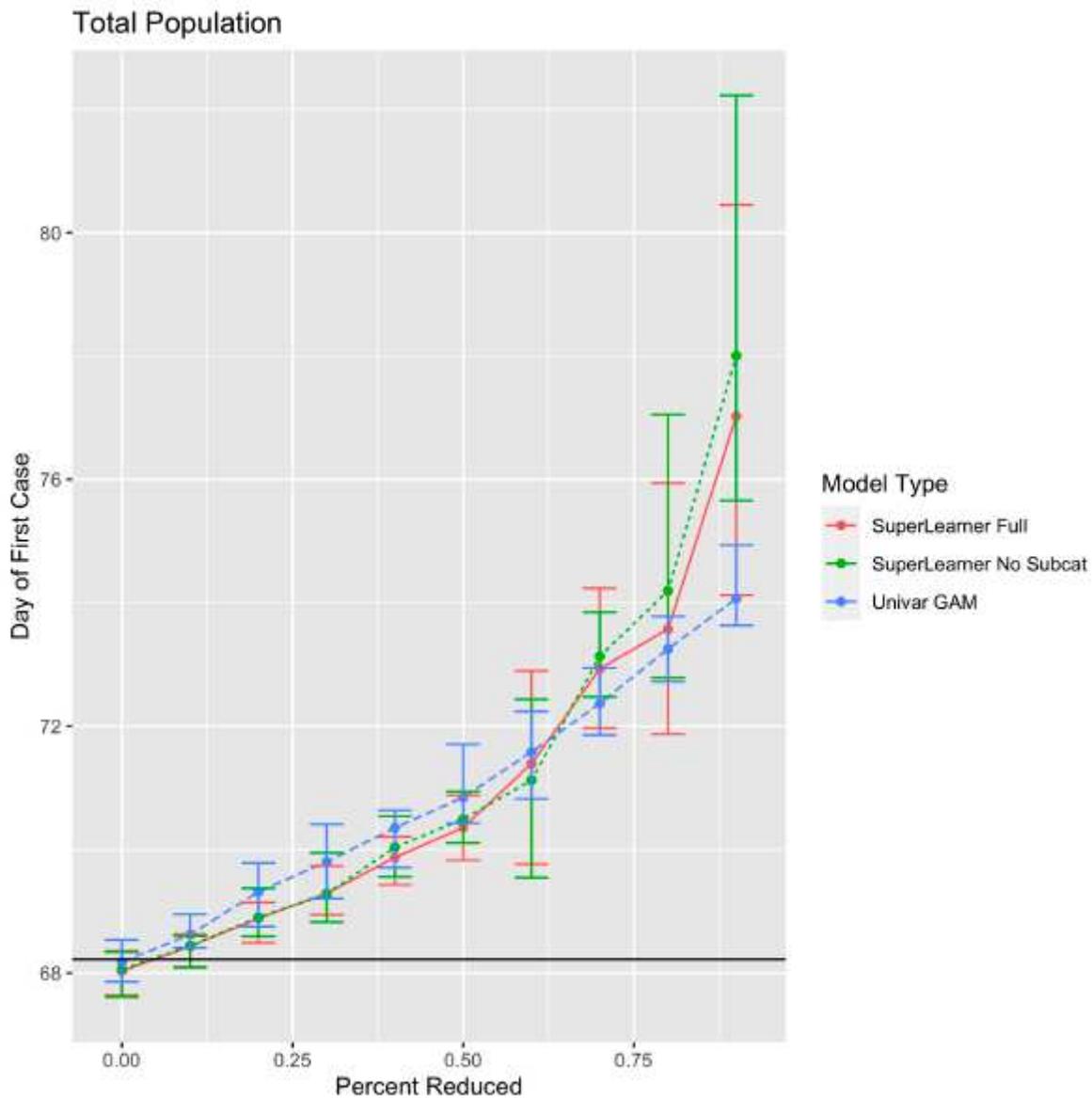


Figure 4

Marginal predictions of day of first case (relative to index time) for different proportional reductions of total population size for models adjusting for all other covariates, only covariates not in sub-category (see supplement table 1) and unadjusted.

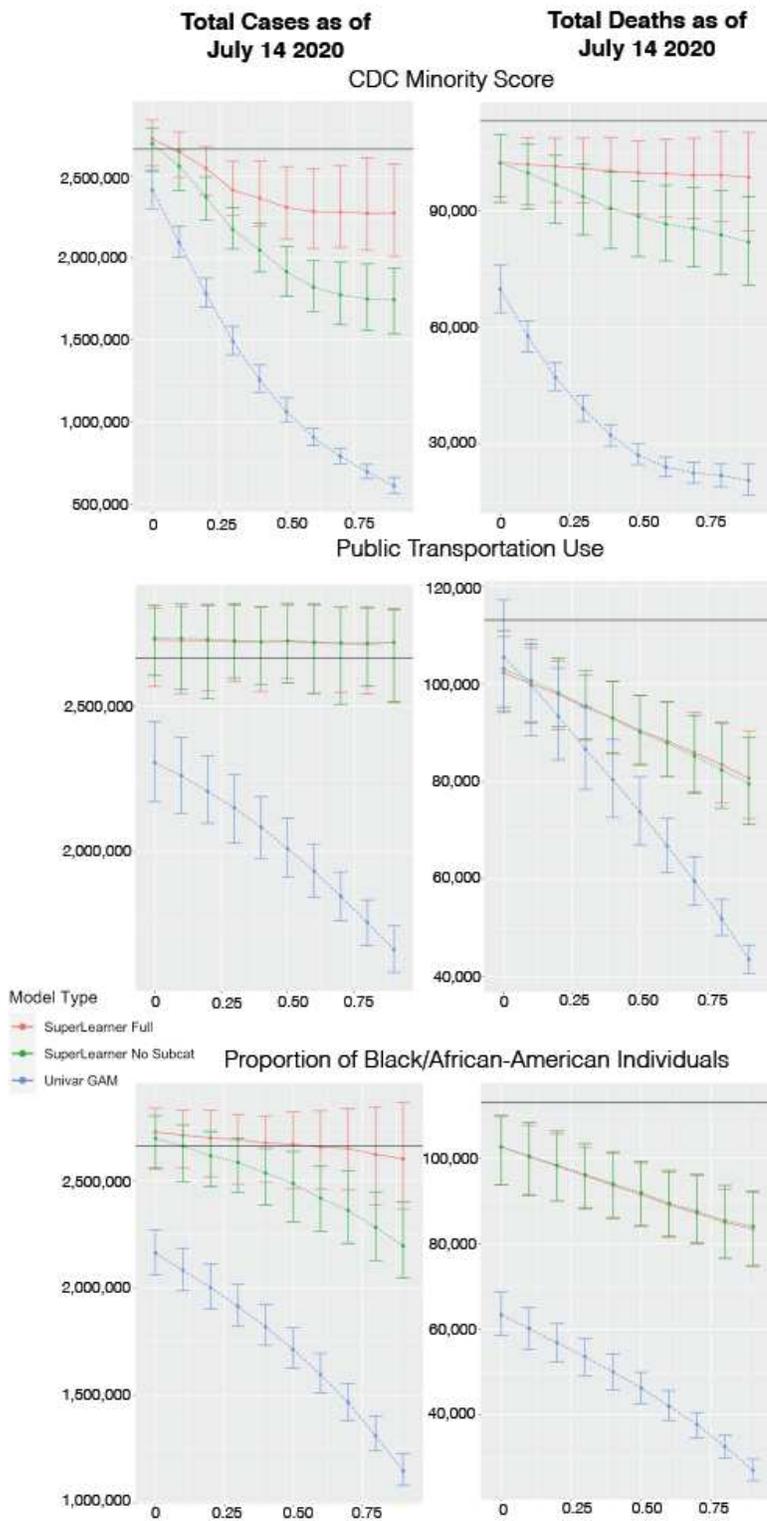


Figure 5

Marginal predictions of total cases and deaths by July 14, 2020) for three of the most consistently important variables in predicting the count outcomes: CDC minority score, proportion of Black individuals and a metric of public transportation use. X-axis is different proportional reductions of each of the three predictors, the Y-axis is the marginal predicted total counts for models adjusting for all other covariates, only covariates not in sub-category (see supplement table 1) and unadjusted.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [GetzHubbardTransportation4.pdf](#)