

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

Predicting severe COVID-19 outcomes for triage and resource allocation

Anthony Onde Morada^{1,2}, Caleb Scheidel³, Jennifer Brown¹, Jeremy Albright³, Victor Kolade¹, and Burt Cagir¹

¹ Guthrie Robert Packer Hospital, Sayre, Pennsylvania, United States of America

² Geisinger Commonwealth School of Medicine, Scranton, PA, United States of America

³ Methods Consultants, Ypsilanti, Michigan, United States of America

*Corresponding Author:

E-mail: burt.cagir@guthrie.org (BC)

17 **Abstract**

18 Background: While numerous studies have identified factors associated with severe
19 COVID-19 outcomes, they have yet to quantify these characteristics. Therefore, our
20 study's purpose is to stratify these risk factors and use them to predict outcomes.

21 Study Design: This is a retrospective review of the CDC COVID-19 Surveillance Data.
22 Logistic regression models calculated risk estimates for independent variables, and
23 random forest models predicted the chance of severe outcomes.

24 Results: Our sample of 3,798,261 patients with COVID-19 consisted mainly of females
25 (51.9%), 10- to 69-year-olds, and White/Non-Hispanics (34.9%). Most were not
26 healthcare workers (90.6%) and did not have preexisting medical conditions (47.1%).

27 Age had an increased risk of severe outcomes that grew every decade of life. White
28 patients had a decreased occurrence of severe outcomes than Non-Whites, except for
29 Pacific Islanders with comparable mortality. The variable selection algorithm detected
30 that three outcomes were more accurate without healthcare worker classification:
31 mechanical ventilation/intubation, pneumonia, and ARDS Acute respiratory distress.

32 However, providers had a decreased risk of severe outcomes overall. Also, patients
33 with preexisting conditions demonstrated an increased risk in all outcomes. Compared
34 to the logistic regressions, the predictive models had a higher performance (AUC>0.8).

35 The death model had the best metrics, followed by hospitalization and ventilation. We
36 amassed these predictive models into the Severe COVID-19 Calculator web application
37 that estimates the probability of severe outcomes.

38 Conclusions: Several patient social and medical demographics recorded by the CDC
39 significantly affect severe COVID-19 outcomes suggesting a multifactorial influence. To

40 account for these variables, a generated Severe Covid-19 Calculator can accurately
41 predict the chance of severe outcomes in citizens that may contract or have COVID-19.
42

43 **Introduction**

44 As the U.S. transitions into the vaccine era of the coronavirus disease-19 (COVID-19)
45 pandemic, health care providers and administrators have to plan allocation of COVID-19
46 treatments and reconsider reopening elective hospital services and surgeries.
47 Therefore, hospitalists and surgeons must balance a patient's risk of nontreatment
48 against potential hospital-acquired COVID-19 infection as hospitals return to pre-
49 pandemic volumes [1–3]. Simultaneously, the government charged hospital
50 administrators and local health care leaders of each state, tribe, and territory to develop
51 their vaccine distribution plan [4–6]. Since the U.S. Food and Drug Administration (FDA)
52 granted emergency use authorization of the Moderna and Pfizer-BioNTech COVID
53 vaccines in late December 2020, the U.S. continues to trail other countries in the
54 proportion of citizens that received at least one vaccine dose [4–9].

55
56 While investigators and journalists have yet to agree, most experts have suggested that
57 technical challenges, lack of federal involvement, and strict adherence to state and CDC
58 priority groups have contributed to the vaccine distribution and administration gap – on
59 January 6th, this gap was as large as 3.7 million to 603,000. [10] Others attribute this
60 vaccination gap to confusing and ambiguous state guidance or poor distribution
61 infrastructure [11–13]. Part of this problem stems from a long-standing ethical question
62 in medicine: How do we distribute a limited treatment equitably and fairly? The diverse
63 array of COVID-19 symptoms complicates this question as severe outcomes, such as
64 hospitalization and death, occur more often in select subgroups, while others are
65 asymptomatic. Researchers also discovered that not unlike our healthcare system pre-

66 pandemic, the social determinants of health amplified these deadly outcomes [14]. A
67 report from the Proceedings of the National Academy of Sciences (PNAS) of the United
68 States of America demonstrated that Black and Latino Americans would be
69 disproportionately left out the CDC Phase 1b recommendation of persons ≥ 75 years of
70 age as minority groups die younger than their white counterparts [15,16].

71
72 While numerous studies have qualified risk factors associated with severe COVID-19
73 outcomes, the current literature has yet to quantify these objective characteristics.
74 Therefore, the purpose of this study is to determine which citizens would most likely
75 develop severe COVID-19 if they contracted the disease. Similar to how the Model for
76 End-Stage Liver Disease (MELD) aids the distribution of livers, our goal is to create an
77 objective and multifactorial algorithm that can stratify patients at risk for severe COVID-
78 19 outcomes [17]. To fulfill these goals, we conducted a retrospective review of the
79 CDC COVID-19 Case Surveillance Restricted Access Detailed Data to enumerate
80 severe COVID-19 outcomes and create a prediction algorithm with machine learning
81 that stratifies a citizen's risk for severe COVID-19 outcomes.

82

83 **Methods**

84 **Data Collection**

85 To determine the rates of severe COVID-19 outcomes, we obtained the CDC COVID-19
86 Case Surveillance Restricted Access Detailed Data after completion and approval
87 through the Registration Information and Data Use Restrictions Agreement (RIDURA)
88 and restricted data access process [18]. The six severe COVID-19 outcomes recorded

89 in the CDC data include hospitalization, intensive care unit (ICU) admission, mechanical
90 ventilation or intubation, pneumonia, acute respiratory distress syndrome (ARDS), and
91 death.

92

93 To ensure our sample best represents the COVID-19 population, we only included
94 patients with laboratory-confirmed COVID-19, thereby excluding suspected cases.

95 Additionally, we chose a study cut-off date of December 15th, 2020, which corresponds
96 to the first day of vaccination in the U.S., to eliminate any potential vaccination effects.

97 Therefore, only positive cases with a date of first specimen collection before December
98 15th, 2020, were included in our study. We also converted the CDC state and region

99 data into the Census Bureau Regions, which reclassifies each patient's residence into

100 one of the five regions: The Northeastern, Midwestern, Southern, Western, and U.S.

101 Territory Regions to account for potential geographic confounding [19].

102

103 The information presented in the database reflects the CDC Human Infection with 2019

104 Novel Coronavirus Person Under Investigation (PUI) and Case Report Forms at an

105 individual patient level. The variable "pre-existing medical conditions" also reflects this

106 reporting form and is defined as having any of the following conditions: chronic lung

107 disease, diabetes mellitus, cardiovascular disease, chronic renal disease, chronic liver

108 disease, immunocompromised condition, neurologic/neurodevelopmental/intellectual

109 disability, other current diseases, current pregnancy status, and current or former

110 smoker [18].

111

112 The CDC Surveillance Review and Response Group (SRRG), as part of the CDC's
113 COVID-19 Emergency Response, maintains the COVID-19 Case Surveillance
114 Restricted Access Detailed Dataset. While the CDC approves database access and
115 study proposals through evaluation of the RIDURA. Institutional review board approval
116 was not required as this study involves data that is a collection of publicly available data
117 with information recorded by the SRRG in such a manner that the subjects cannot be
118 identified directly or through identifiers linked to the subjects (45 Code of Federal
119 Regulations (CFR) 46.101(b)). The CDC does not take responsibility for the scientific
120 validity or accuracy of methodology, results, statistical analyses, or conclusions
121 presented.

122

123 Determining Risk Factors

124 Baseline demographic variables and clinical outcomes were presented with frequencies
125 and percentages. The following outcomes were considered: hospitalization, ICU
126 admission, mechanical ventilation/intubation, pneumonia, acute respiratory distress, and
127 mortality. For each outcome, an initial multivariable logistic regression model fit all
128 baseline demographics and symptoms as covariates. Forest plots showing odds ratios
129 and 95% confidence intervals for each covariate in the model were presented, and
130 model fit metrics including Akaike information criterion (AIC), Bayesian information
131 criterion (BIC), Area under the curve (AUC), and pseudo-R² measures were shown.
132 The Least Absolute Shrinkage and Selection Operator (LASSO) logistic regression was
133 then used as a variable selection technique to assess if any variables were
134 uninformative to predicting each outcome. The LASSO algorithm's penalty

135 hyperparameter was tuned using 5-fold cross-validation to determine its optimal value
136 based on the area under the receiver operating characteristic curve (AUC). Variables
137 that were shrunk to zero in the LASSO models were dropped, and another logistic
138 regression model was fit for each outcome with the variables kept by the LASSO model.
139 Model fit metrics were then presented for these models to determine if dropping
140 uninformative variables improved model fit. Statistical significance was defined as a p-
141 value < 0.05.

142

143 Predictive Model Creation

144 For each outcome, a random forest model was fit with baseline demographics, including
145 age, sex, race/ethnicity, region, healthcare worker, pre-existing medical conditions, and
146 month of positive COVID-19 test as predictive features.

147

148 Before modeling, the data was first split into training and testing sets. Due to the large
149 sample size, we took a random sample of 3% of the data for computational speed. We
150 then took a random sample of 80% of that data to train the model and used the
151 remaining 20% of the data as the test set. Data pre-processing included creating
152 indicator variables for all non-numeric categorical values of the predictor variables and
153 using indicators of the month of the positive test date for accounting treatment
154 improvement over time. Since fewer patients experienced each outcome than those
155 who didn't, we upsampled each model's outcome to balance the classes.

156

157 The random forest algorithm's hyperparameters were tuned using 3-fold cross-
158 validation on the training set, with the number of trees contained in the ensemble fixed
159 at $N = 1000$. The hyperparameters correspond to the number of predictors randomly
160 sampled at each split when creating the tree models (m_{try}) and the minimum number
161 of data points in a node required for the node to be split further (min_n). The optimal
162 hyperparameters were chosen based on which values maximized the area under the
163 receiver operating characteristic curve (AUC). The AUC metric measures fit from 0.5
164 (no better than a coin flip) to 1.0 (perfect prediction).

165
166 Final models were fit on the entire training set and evaluated on the test set using the
167 following metrics to assess performance: accuracy, precision, recall, and AUC.
168 Accuracy is the number of correct predictions divided by the total number of predictions
169 or the percentage of predictions from the correct model. Precision is the number of true
170 positives divided by the sum of the true and false positives, demonstrating how precise
171 a model is out of those it predicted positive. The recall is the number of true positives
172 divided by the true positives and false negatives, accounting for the percentage of
173 relevant results correctly classified by the model.

174
175 The random forest models' relative predictive performance was compared to the logistic
176 models' performance in risk factor estimation. Feature importance plots for the random
177 forest models were produced to show which variables were most important to making
178 the predictions for each outcome. All analyses were performed in R Version 4.0.3 [20].

179

180 Results

181 The total sample contained 3,798,261 patients with laboratory-confirmed cases of
182 COVID-19 between January 1st, 2020, and December 15th, 2020. Table1 presents the
183 sample demographics. Our sample of patients with COVID-19 were primarily women
184 between the ages of 10 - 69, white/non-Hispanic, and who lived in the Northeast CDC
185 Census Regions. While most patients' status regarding healthcare work or pre-existing
186 medical conditions was missing or unknown, most subjects with COVID-19 were not
187 healthcare workers and did not have any pre-existing medical conditions (Table 1).

188

189 **Table 1. Demographics of patients with COVID-19.**

Variable	Label	Patients with COVID-19 (N = 3,798,261)
<hr/>		
Sex		
	Male	1,795,160 (47.6%)
	Female	1,954,915 (51.9%)
	Other or Unknown	18,709 (0.50%)
Age Group		
	0 – 9 Years	141,591 (3.70%)
	10 – 19 Years	379,438 (10.0%)
	20 – 29 Years	711,770 (18.8%)
	30 – 39 Years	611,965 (16.2%)

40 – 49 Years	562,714 (14.9%)
50 – 59 Years	555,063 (14.7%)
60 – 69 Years	403,679 (10.7%)
70 – 79 Years	231,412 (6.1%)
80+ Years	185,817 (4.9%)
Unknown	14,812 (0.39%)

Race/Ethnicity

White, Non-Hispanic	1,324,461 (34.87%)
American Indian/Alaska Native, Non-	
Hispanic	15,431 (0.41%)
Asian, Non-Hispanic	67,567 (1.78%)
Black, Non-Hispanic	288,992 (7.61%)
Hispanic/Latino	489,383 (12.88%)
Multiple/Other, Non-Hispanic	111,481 (3.10%)
Native Hawaiian/Other Pacific Islander,	
Non-Hispanic	11,060 (0.29%)
Unknown	1,489,886 (39.23%)

CDC Census

Region

Northeast	1,325,064 (34.89%)
Midwest	745,419 (19.63%)
South	1,050,551 (27.7%)
Pacific	669,141 (17.62%)

	U.S. Territory	7,282 (0.19%)
	Unknown	804 (0.02%)
Health Care Worker		
	No	1,255,055 (33.04%)
	Yes	130,446 (3.43%)
	Unknown	2,412,760 (63.52%)
Pre-existing Medical Conditions		
	No	711,859 (18.74%)
	Yes	633,792 (16.69%)
	Unknown	2,452,610 (64.57%)

190

191 Table 2 presents the six severe COVID-19 outcomes labeled in the CDC dataset. We
 192 found that providers hospitalized 13.52% of patients with COVID-19, admitted 4.98% of
 193 patients to the ICU, and intubated or mechanically ventilated 3.02% of patients.
 194 Additionally, 6.36% of patients with COVID-19 eventually developed pneumonia, 1.78%
 195 developed ARDS, and 4.87% eventually died.

196

197 **Table 2. Severe COVID-19 outcomes.**

Outcome	Label	Patients with COVID-19 (N = 3,798,261)	Missing Data
Hospitalization			53.5%
	No	1,526,273 (86.5%)	

	Yes	238,638 (13.5%)	
ICU Admission			79.3%
	No	747,280 (95.0%)	
	Yes	39,152 (5.0%)	
Mechanical Ventilation/Intubation			84.1%
	No	585,626 (97.0%)	
	Yes	18,244 (3.0%)	
Pneumonia			79.1%
	No	741,769 (93.6%)	
	Yes	50,373 (6.4%)	
Acute Respiratory Distress			79.2%
	No	777,033 (98.2%)	
	Yes	14,054 (1.8%)	
Death			54.6%
	No	1,638,857 (95.1%)	
	Yes	83,810 (4.9%)	

198

199 To understand the risks of severe COVID-19 outcomes in special populations, we
 200 performed a univariate subgroup analysis on healthcare providers and subjects with
 201 pre-existing conditions. Table 3 demonstrates that patients with pre-existing medical
 202 conditions had a higher risk of developing all six severe COVID-19 outcomes than the
 203 opposing group. Likewise, this table also illustrates that healthcare workers have a
 204 significantly decreased risk of obtaining all six severe COVID-19 outcomes than non-

205 healthcare workers. To determine the risk of severe outcomes by age, we plotted the
206 age groups by the proportion of patients with severe COVID-19 in Fig 1. This figure
207 demonstrates that all six outcomes increase in older subjects, while death and
208 hospitalizations escalate at a steeper rate.

209

210 **Table 3. Severe COVID-19 outcomes in special populations.**

Outcomes	Medical Conditions (N = 633,792)	Healthcare Workers (N = 130,446)	Overall (N = 3,798,261)
Hospitalized	125,449 (23.9%) *	4,191 (4.2%) *	238,638 (13.5%)
ICU Admission	25,656 (8.9%) *	774 (1.6%) *	39,152 (5.0%)
Mechanical Ventilation	12,964 (5.3%) *	301 (0.7%) *	18,244 (3.0%)
Pneumonia	34,935 (10.9%) *	2,020 (4.1%) *	50,373 (6.4%)
Acute Respiratory Distress	10,397 (3.3%) *	485 (1.0%) *	14,054 (1.8%)
Death	46,449 (10.3%) *	233 (0.2%) *	83,810 (4.9%)

211 * Proportions are different compared to overall population (p < 0.001)

212

213 **Fig 1. Risk of severe COVID-19 by age group**

214

215 We accounted for potential confounding factors in a multivariate logistic regression
216 model for each of the six severe outcomes controlling for the same seven independent
217 variables of sex, age group, race/ethnicity, healthcare worker status, U.S. Census
218 region, pre-existing medical condition status, and month of positive COVID-19 test

219 (Figure 2). We compared the logistic regression model using all variables to the LASSO
220 selected models and presented the best model metrics to evaluate predictive capacity.

221
222 **Fig 2. Logistic regression models for severe COVID-19 outcomes.** Forest plot for
223 each severe COVID-19 outcome. a) Hospitalization – all variables, b) ICU admission –
224 all variables, c) Mechanical ventilation/intubation – LASSO selected variables, d)
225 Pneumonia – LASSO selected variables, e) Acute respiratory distress – LASSO
226 selected variables, and f) Mortality – all variables. Male is the reference group for sex.
227 Age 20-29 Years is the reference group for age. Non-Hispanic White is the reference
228 group for race/ethnicity. Northeast is the reference group for the region. The reference
229 group for the month of the positive test is January-March.

230
231 Figure 2a illustrates the results of the logistic regression model created for the
232 hospitalization outcome. The LASSO algorithm did not drop any of the variables, and
233 therefore, we did not create its associated model (S1 Table). When controlling for all
234 other covariates, we found the three largest effect sizes were in patients over 80,
235 between the ages of 70-79, and with pre-existing medical conditions. The forest plot
236 shows that patients over 80 and 70-79 years old had seventeen- and fourteen-times
237 higher odds of being hospitalized than those with ages 20-29, respectively. Patients with
238 a pre-existing medical condition had a three-fold increased risk of being hospitalized.

239
240 Additionally, the LASSO algorithm did not drop any variables while creating the ICU
241 admission logistic regression model (Fig 2b, S2 Table). Patients ≥ 80 and between the

242 ages of 70-79 had a seventeen- and twenty-times increased odds of admission to the
243 ICU, respectively. American Indian/Alaska Native, non-Hispanic patients had a four
244 times greater risk of ICU admission likelihood than white, non-Hispanic patients.
245 Patients with a pre-existing medical condition had a four times greater risk of ICU
246 admission.

247
248 In the outcome of mechanical ventilation or intubation, the LASSO algorithm dropped
249 the variable healthcare worker status (Fig 2c). However, compared to the logistic
250 regression model with all variables, the model fit metrics were nearly equivalent,
251 suggesting that the LASSO model did not significantly improve the model fit for
252 mechanical ventilation (S1 Fig, S3 Table). Fig 2c illustrates the forest plot for the
253 LASSO selected variables for the mechanical ventilation logistic regression models. The
254 characteristics with the most significant effect sizes include patients 70-79 years old, 60-
255 69 years old, and those with pre-existing medical conditions. Patients in their seventh
256 and eighth decades of life had a 17- and 12-fold increase risk of being intubated or put
257 on mechanical ventilation, respectively. Subjects with pre-existing medical conditions
258 had a threefold increase of being intubated or ventilated.

259
260 The LASSO model for pneumonia also dropped healthcare workers as a variable to
261 include; however, compared to the all-variables model, the metrics were relatively
262 equivalent, suggesting that dropping the variable did not improve model fit (S2 Fig S4
263 Table). Fig 2d presents the forest plot for the LASSO selected logistic regression model.
264 Patients over the age of 80 and between the ages of 70-79 had the most significant

265 effect sizes with an increased odds of developing pneumonia by 17- and 12-fold,
266 respectively. Patients with pre-existing medical conditions had the third-highest effect
267 size with an OR of 2.8.

268
269 The LASSO algorithm for ARDS also dropped healthcare workers from its model;
270 however, compared to the logistic regression with all variables, the variables' effect
271 sizes and p-values were similar (S3 Fig, S5 Table). The forest plot in Fig 2e presents
272 the ORs of the logistic regression for the LASSO selected variables. Patients over 80
273 and in their seventh decade of life had a fifteen- and eleven-times increased risk of
274 developing ARDS compared to the 20–29-year age group, respectively. After age
275 groups, the third-largest effect size demonstrates that patients with pre-existing medical
276 conditions have increased ARDS odds by three.

277
278 The LASSO model for mortality did not drop any variables (S6 Table). We present the
279 mortality regression models with all variables in Fig 2f. Patients older than 80 years old
280 and between the ages of 70-79 had the two largest effect sizes. Patients over 80 and in
281 the seventh decade of life have 356- and 122-times increased risk of mortality,
282 respectively, compared to patients aged 20-29. Pre-existing medical conditions had the
283 third-highest effect size with a fourfold increased risk of mortality.

284
285 Next, we created a prediction model using the same independent variables to calculate
286 a patient's risk of severe COVID-19 outcomes if they were to contract the disease. We
287 present the hyperparameter tuning results from 3-fold cross-validation on the training

288 sets in S7 Table. For each outcome, an initial grid search across a broader range of
289 values led to an optimum number of variables available for splitting at each tree node,
290 or mtry value of 1 and a minimum number of data points required for the node to be split
291 further, or min_n value of 18 leading to the best performance. The grid search was then
292 repeated on a narrower range of hyperparameter values to produce the ROC curves for
293 each outcome.

294
295 The ROC curves for the final random forest models fit the entire training set using the
296 optimal hyperparameter values for each outcome shown in Figure 3. All of the models
297 had a high AUC > 0.8. The model with death as the outcome had the highest
298 performance (AUC = 0.953), followed by hospitalization (AUC = 0.892) and mechanical
299 ventilation (AUC = 0.882).

300
301 **Fig 3. ROC curves of random forest models.** Abbreviation = ROC, receiver operating
302 characteristic curve

303
304 Table 4 presents the performance metrics for both the random forest and LASSO
305 models for each outcome. Relative to the LASSO models, the random forest models
306 had a higher performance in overall accuracy for each outcome. The models for
307 predicting mechanical ventilation and deaths had higher accuracy than predicting
308 pneumonia and ARDS. The random forest models for hospitalizations and death had a
309 strong recall, with a low number of false-negative predictions. Overall, we found that the
310 random forest models had a low precision with a high number of false positives. We

311 presented the relative feature importance for each model in Fig 4 to rank the patient
312 characteristics that had the most influence on each outcome. Age groups and presence
313 of a pre-existing condition were the top features for each outcome. Northeast region
314 was an important feature in predicting hospitalization or ventilation. This pattern draws
315 upon high hospitalization and ventilation rates in the Northeast during the first COVID-
316 19 wave in the spring. The youngest age group is likely a solid indicator for the model of
317 a lack of the outcome occurring.

318

319 **Table 4. Performance metrics.**

Outcome	Metric	Logistic	
		Random Forest	Regression
Hospitalization	Accuracy	0.818	0.808
	Precision	0.412	0.396
	Recall	0.804	0.804
	AUC	0.889	0.886
ICU Admission	Accuracy	0.786	0.777
	Precision	0.165	0.169
	Recall	0.741	0.817
	AUC	0.85	0.861
Mechanical Ventilation	Accuracy	0.863	0.814
	Precision	0.12	0.106
	Recall	0.663	0.821
	AUC	0.877	0.904

Pneumonia	Accuracy	0.738	0.756
	Precision	0.175	0.186
	Recall	0.746	0.743
	AUC	0.811	0.824
Acute Respiratory Distress	Accuracy	0.781	0.768
	Precision	0.053	0.055
	Recall	0.713	0.787
	AUC	0.809	0.853
Death	Accuracy	0.879	0.872
	Precision	0.275	0.267
	Recall	0.896	0.914
	AUC	0.952	0.953

320 AUC, Area under the receiver operating characteristic (ROC) curve.

321

322 **Fig 4. Figure importance by the outcome.** Abbreviation = HC Worker, Healthcare
323 Worker.

324

325 We aggregated the predictive models into the Severe COVID-19 Calculator web
326 application and published it online ([https://methodsconsultants-
327 apps.shinyapps.io/guthrie-cdc-covid-prediction/](https://methodsconsultants-apps.shinyapps.io/guthrie-cdc-covid-prediction/)). The web application accounts for all
328 predictive variables and provides a predicted estimate of the probability of severe
329 COVID-19 outcomes (Fig 5).

330

331 **Fig 5. Example of the Severe COVID-19 Calculator.**

332

333 **Discussion**

334 Among 3,798,261 patients with COVID-19 from the pre-vaccine era in the United
335 States, older patients had an increased risk of all six severe outcomes. The odds ratios
336 are significantly higher in the latter decades of life as some outcomes, such as death,
337 have a 300-fold increase. We discovered that White subjects had a decreased
338 occurrence of all severe outcomes than Non-Whites, except Native Hawaiian/Other
339 Pacific Islander patients who had similar mortality risks. The models also demonstrated
340 that the rates of severe COVID-19 effects were unequal between all five of the CDC
341 Census Regions. In both the LASSO selected and all-variable models, we found that
342 healthcare workers had a decreased risk of severe COVID-19 outcomes. Additionally,
343 we discovered the patients with pre-existing health conditions all had an increased
344 chance of severe outcomes. When we adjusted for the diagnosis time, we uncovered
345 that severe COVID-19 outcomes decreased by month as the pandemic progressed.

346

347 Overall, our results confirm the current COVID-19 literature on severe outcomes. We
348 found that age has the most significant effect measure compared to our other
349 covariables. However, instead of using threshold ages, our age groups provide more
350 details on risk stratification in special populations such as school children. A
351 correspondence published in the New England Journal of Medicine examining child and
352 teacher morbidity in Sweden without school closure also confirmed the low incidence of
353 severe COVID-19 among schoolchildren and preschool-age children during the

354 pandemic [21]. While most researchers confirmed the effects of chronic diseases on
355 severe COVID-19, our results demonstrate that pre-existing conditions have the second
356 most significant influence on outcomes after age [22]. These chronic conditions like
357 obesity and diabetes demonstrate an increased clinical severity in disease
358 presentations, while infrastructure shutdown delayed treatments in conditions such as
359 mental health disorders, chronic kidney disease, and cancer [22]. Additionally, our
360 results also demonstrate that race/ethnicity has a significant influence on severe
361 outcomes. A report from PNAS presented a disproportionate mortality rate in the Black
362 and Latino populations compared to White patients [16]. However, we discovered that
363 all Non-White populations, in general, had higher rates of severe outcomes.
364 Additionally, we found that while healthcare workers consisted of 3% of our sample,
365 they had a decreased risk of severe COVID-19 outcomes or had no influence on their
366 outcomes in this pre-vaccine era. Our results differ from the current idea that healthcare
367 workers mirror the general population rates of severe outcomes [23]. Because a few
368 models excluded the variable healthcare workers, we suspect that the decreased
369 occurrence of severe outcomes reflects a younger and healthier subgroup compared to
370 the general population. A systematic review of COVID-19 in healthcare workers
371 worldwide mirrors our results as providers over 70 or males had a higher risk of death
372 [23].

373

374 We present the first report on the CDC COVID-19 case surveillance system. This
375 restricted database presents individual-level data from autonomous reporting entities
376 from all U.S. states and territories. The large sample size and standardized data

377 dictionary allowed us to perform multiple logistic regression models to control each
378 predictive variable. Therefore, we are confident in our adjusted risk calculations as they
379 already account for any potential interactions. To the best of our knowledge, this is the
380 first study to stratify by date of COVID-19 diagnosis to account for any treatment biases
381 as COVID-19 management has changed over time. Additionally, as states have
382 different COVID-19 protocols and strains occupy different regions of the U.S., our study
383 is the first to account for these factors by accounting for geographical regions.

384
385 Our web app for our Severe COVID-19 Calculator allows anyone to estimate the risk of
386 the severe outcome if they have or if they were to contract COVID-19. As mentioned
387 throughout our report, all of the predictive variables discussed significantly influence
388 severe COVID-19 outcomes. Thereby, our calculator can provide accurate numeric
389 predictions using the subject's objective data without arbitrary thresholds or scoring
390 systems. For example, Fig 5 illustrates the user interface and outcome estimates for a
391 hypothetical patient.

392
393 While our study accounted for socioeconomic factors such as sex, race, and region, the
394 available data limited our ability to account for other important social determinants of
395 health, such as income and occupation, which may affect severe COVID-19 outcomes
396 and vaccine distribution plans [24]. Additionally, to preserve patient anonymity, we could
397 not distinguish between types of pre-existing conditions. This limitation also questions
398 the possibility of the number of pre-existing conditions correlating with disease severity.
399 If so, there may be a benefit in making this distinction in future predictive tools.

400 Like all prediction modeling projects, the exclusion of potentially essential variables
401 could reduce our models' performance. Also, within the available data, a few predictive
402 variables, such as race/ethnicity, healthcare worker status, pre-existing medical
403 conditions, contained unknown variables, which could have affected our calculator's
404 accuracy. However, the high AUC value and performance metrics for all of our
405 outcomes suggest that our predictive models are significantly better than random
406 chance.

407
408 As vaccine distribution continues to roll out in phases and trends of severe COVID-19
409 outcomes change, we will need to continuously recalculate our results to determine the
410 patients at the highest risk for severe COVID-19 outcomes. With the continuously
411 updated CDC data and our publicly available web app, we can recreate our Severe
412 COVID-19 Calculator to reflect the most up-to-date data. For example, in a new data
413 release, we may find a race/ethnicity, sex, or age group with a disproportionate chance
414 of severe COVID-19 outcomes in the vaccination era, and we would like to account for
415 that in future calculator iterations. Additionally, using some of the variables already
416 available in the CDC report, we plan to aggregate data sources to account for variables
417 of interest, such as estimated patient income.

418
419 Since the start of the pandemic, researchers have published dozens of studies that
420 identified risk factors for severe outcomes. However, the literature lacked any detail on
421 each risk factor's relative importance or adjusted for potential covariable interactions.
422 Our study performed these adjustments and found that predicting severe COVID-19

423 outcomes is a multifactorial problem and quantified that variables such as age, pre-
424 existing health conditions, and race/ethnicity are more important than others like health
425 care worker status.

426
427 While we provide an objective tool to assess a subject's risk of developing severe
428 COVID-19 outcomes, we have no comment on current vaccine distribution plans. The
429 Severe COVID-19 Calculator is an objective tool that mirrors other models already used
430 in medicine to help providers and researchers stratify patients in resource scarcity, such
431 as hospital beds, ventilators, masks, or vaccines. This study's sole purpose was to
432 identify independent risk factors and quantify these effects to understand the apparent
433 outcome disparities between different patient groups. After adjusting for covariates,
434 patients that are older, male, Non-White, non-healthcare workers, or possess at least
435 one pre-existing condition have an increased risk of severe COVID-19 outcomes.
436 Additionally, we found that our Severe COVID-19 Calculator accurately predicts the
437 chance of hospitalization, ICU hospitalization, mechanical ventilation/intubation,
438 pneumonia, ARDS, and death in citizens that may contract or have COVID-19.
439

440 **Acknowledgments**

441 We thank Vicky Hickey, Dr. Sri Harshavardhan Senapathi MD, Dr. Seungwoo Chai MD,
442 and Dr. Karen Williams, PharmD., BCPS-AQ ID (Guthrie Robert Packer Hospital) for
443 helping us develop our study design and comments on the manuscript.
444

445 **References**

- 446 1. Rivett L, Sridhar S, Sparkes D, Routledge M, Jones NK, Forrest S, et al.
447 Screening of healthcare workers for SARS-CoV-2 highlights the role of
448 asymptomatic carriage in COVID-19 transmission. *Elife*. 2020;9.
449 doi:10.7554/eLife.58728
- 450 2. Garcia-Vidal C, Sanjuan G, Moreno-García E, Puerta-Alcalde P, Garcia-Pouton
451 N, Chumbita M, et al. Incidence of co-infections and superinfections in
452 hospitalized patients with COVID-19: a retrospective cohort study. *Clin Microbiol*
453 *Infect*. 2021;27: 83–88. doi:10.1016/j.cmi.2020.07.041
- 454 3. Elkrief A, Desilets A, Papneja N, Cvetkovic L, Groleau C, Lakehal YA, et al. High
455 mortality among hospital-acquired COVID-19 infection in patients with cancer: A
456 multicentre observational cohort study. *Eur J Cancer*. 2020;139: 181–187.
457 doi:10.1016/j.ejca.2020.08.017
- 458 4. American Academy of Family Physicians. State COVID-19 Vaccine Distribution
459 Plans. Leawood, USA; 2020. Available: [AAFP.ORG](https://www.aafp.org)
- 460 5. Centers for Disease Control and Prevention. State & Territorial Health
461 Department Websites. In: Centers for Disease Control and Prevention (CDC)
462 [Internet]. 2021 [cited 1 Feb 2021]. Available:
463 [https://www.cdc.gov/publichealthgateway/healthdirectories/healthdepartments.ht](https://www.cdc.gov/publichealthgateway/healthdirectories/healthdepartments.html)
464 [ml](https://www.cdc.gov/publichealthgateway/healthdirectories/healthdepartments.html)
- 465 6. Assistant Secretary for Public Affairs. COVID-19 Vaccine Distribution□: The
466 Process. In: U.S. Department of Health and Human Services [Internet]. 2020
467 [cited 1 Feb 2021]. Available: <https://www.hhs.gov/coronavirus/covid-19->

- 468 [vaccines/distribution/index.html](#)
- 469 7. Dooling K. ACIP COVID-19 Vaccines Work Group Phased Allocation of COVID-
470 19 Vaccines Policy Questionnaire: 2020. Available:
471 [https://www.cdc.gov/vaccines/acip/meetings/downloads/slides-2020-11/COVID-](https://www.cdc.gov/vaccines/acip/meetings/downloads/slides-2020-11/COVID-04-Dooling.pdf)
472 [04-Dooling.pdf](https://www.cdc.gov/vaccines/acip/meetings/downloads/slides-2020-11/COVID-04-Dooling.pdf)
- 473 8. Reiss CS. Coronavirus Pandemic. *DNA Cell Biol.* 2020;39: 919–919.
474 doi:10.1089/dna.2020.29015.csr
- 475 9. Huang P, Audrey C. How Is The COVID-19 Vaccination Campaign Going In Your
476 State? In: NPR [Internet]. 2021 [cited 8 Mar 2021]. Available:
477 [https://www.npr.org/sections/health-shots/2021/01/28/960901166/how-is-the-](https://www.npr.org/sections/health-shots/2021/01/28/960901166/how-is-the-covid-19-vaccination-campaign-going-in-your-state)
478 [covid-19-vaccination-campaign-going-in-your-state](https://www.npr.org/sections/health-shots/2021/01/28/960901166/how-is-the-covid-19-vaccination-campaign-going-in-your-state)
- 479 10. Centers for Disease Control and Prevention. COVID-19 Vaccinations in the
480 United States. In: Centers for Disease Control and Prevention [Internet]. 2021
481 [cited 1 Feb 2021] p. 1. Available: [https://www.cdc.gov/coronavirus/2019-](https://www.cdc.gov/coronavirus/2019-ncov/vaccines/index.html)
482 [ncov/vaccines/index.html](https://www.cdc.gov/coronavirus/2019-ncov/vaccines/index.html)
- 483 11. Kincade K. Struggles people are having booking COVID-19 vaccine appointment.
484 WWLP. 2021. Available: [https://www.wwlp.com/news/health/coronavirus-local-](https://www.wwlp.com/news/health/coronavirus-local-impact/struggles-people-are-having-booking-covid-19-vaccine-appointment/)
485 [impact/struggles-people-are-having-booking-covid-19-vaccine-appointment/](https://www.wwlp.com/news/health/coronavirus-local-impact/struggles-people-are-having-booking-covid-19-vaccine-appointment/)
- 486 12. Scheible S. 5 seniors, 5 stories: Local residents struggle to secure vaccination
487 appointments. *The Patriot Ledger.* 2021. Available:
488 [https://www.patriotledger.com/story/news/2021/01/28/75-plus-trying-get-vaccine-](https://www.patriotledger.com/story/news/2021/01/28/75-plus-trying-get-vaccine-depends-luck-and-kindness/4285777001/)
489 [depends-luck-and-kindness/4285777001/](https://www.patriotledger.com/story/news/2021/01/28/75-plus-trying-get-vaccine-depends-luck-and-kindness/4285777001/)
- 490 13. Otterman S. The Maddening Red Tape Facing Older People Who Want the

- 491 Covid-19 Vaccine. The New York Times. 2021. Available:
492 [https://www.nytimes.com/2021/01/14/nyregion/covid-vaccine-older-people-senior-](https://www.nytimes.com/2021/01/14/nyregion/covid-vaccine-older-people-senior-citizens.html)
493 [citizens.html](https://www.nytimes.com/2021/01/14/nyregion/covid-vaccine-older-people-senior-citizens.html)
- 494 14. Paremoer L, Nandi S, Serag H, Baum F. Covid-19 pandemic and the social
495 determinants of health. *BMJ*. 2021;372: n129. doi:10.1136/bmj.n129
- 496 15. Thomas WC, Grabenstein H. First up for COVID vaccine: 75 and older. The
497 average Black person in Shelby County doesn't live that long. *MLK50*. 12 Feb
498 2021. Available: [https://mlk50.com/2021/02/12/first-up-for-covid-vaccine-75-and-](https://mlk50.com/2021/02/12/first-up-for-covid-vaccine-75-and-older-the-average-black-person-doesnt-live-that-long/)
499 [older-the-average-black-person-doesnt-live-that-long/](https://mlk50.com/2021/02/12/first-up-for-covid-vaccine-75-and-older-the-average-black-person-doesnt-live-that-long/)
- 500 16. Andrasfay T, Goldman N. Reductions in 2020 US life expectancy due to COVID-
501 19 and the disproportionate impact on the Black and Latino populations. *Proc Natl*
502 *Acad Sci*. 2021;118: e2014746118. doi:10.1073/pnas.2014746118
- 503 17. Wiesner R, Edwards E, Freeman R, Harper A, Kim R, Kamath P, et al. Model for
504 end-stage liver disease (MELD) and allocation of donor livers. *Gastroenterology*.
505 2003;124: 91–96. doi:10.1053/gast.2003.50016
- 506 18. Centers for Disease Control and Prevention. COVID-19 Case Surveillance Data
507 Access, Summary, and Limitations. Atlanta, GA; 2020.
- 508 19. U.S. Department of Commerce. Census regions and divisions of the United
509 States. 2010; 1–2. Available: [https://www2.census.gov/geo/pdfs/maps-](https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us_regdiv.pdf)
510 [data/maps/reference/us_regdiv.pdf](https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us_regdiv.pdf)
- 511 20. Core R Team. A Language and Environment for Statistical Computing. R
512 Foundation for Statistical Computing. Vienna, Austria: R Foundation for Statistical
513 Computing; 2019. p. [https://www.R--project.org](https://www.R-project.org).

- 514 21. Ludvigsson JF, Engerström L, Nordenhäll C, Larsson E. Open Schools, Covid-19,
515 and Child and Teacher Morbidity in Sweden. *N Engl J Med*. 2021;384: 669–671.
516 doi:10.1056/NEJMc2026670
- 517 22. Flaherty GT, Hession P, Liew CH, Lim BCW, Leong TK, Lim V, et al. COVID-19 in
518 adult patients with pre-existing chronic cardiac, respiratory and metabolic disease:
519 a critical literature review with clinical recommendations. *Trop Dis Travel Med*
520 *Vaccines*. 2020;6: 16. doi:10.1186/s40794-020-00118-y
- 521 23. Bandyopadhyay S, Baticulon RE, Kadhum M, Alser M, Ojuka DK, Badereddin Y,
522 et al. Infection and mortality of healthcare workers worldwide from COVID-19: a
523 systematic review. *BMJ Glob Heal*. 2020;5: e003097. doi:10.1136/bmjgh-2020-
524 003097
- 525 24. Raifman MA, Raifman JR. Disparities in the Population at Risk of Severe Illness
526 From COVID-19 by Race/Ethnicity and Income. *Am J Prev Med*. 2020;59: 137–
527 139. doi:10.1016/j.amepre.2020.04.003
- 528
- 529

530 **Supporting information**

531 **S1 Fig. Forest plot for mechanical ventilation logistic regression model with all**
532 **variables.** Male is the reference group for sex. Age 20-29 Years is the reference group
533 for age. Non-Hispanic White is the reference group for race/ethnicity. Northeast is the
534 reference group for the region. The reference group for the month of the positive test is
535 Jan-March.

536

537 **S2 Fig. Forest plot for pneumonia logistic regression model with all variables.**
538 Male is the reference group for sex. Age 20-29 Years is the reference group for age.
539 Non-Hispanic White is the reference group for race/ethnicity. Northeast is the reference
540 group for the region. The reference group for the month of the positive test is Jan-
541 March.

542

543 **S3 Fig. Forest plot for acute respiratory distress logistic regression model – all**
544 **variables.** Male is the reference group for sex. Age 20-29 Years is the reference group
545 for age. Non-Hispanic White is the reference group for race/ethnicity. Northeast is the
546 reference group for the region. The reference group for the month of the positive test is
547 Jan-March.

548

549 **S1 Table. Model fit statistics for hospitalization logistic regression model.**

Metric	Full Model with All Variables
AIC	921,131.2

BIC	921,589.4
AUC	0.121
McFadden's R ²	0.341
Nagelkerke's R ²	0.433

550 Abbreviations = AIC, Akaike information criterion; BIC, Bayesian information criterion;

551 AUC; area under the curve.

552

553 **S2 Table. Model fit statistics for ICU admission logistic regression models.**

Metric	Full Model with All Variables
AIC	237,255.7
BIC	237,684.0
AUC	0.139
McFadden's R ²	0.238
Nagelkerke's R ²	0.275

554 Abbreviations = AIC, Akaike information criterion; BIC, Bayesian information criterion;

555 AUC; area under the curve.

556

557 **S3 Table. Model fit statistics for mechanical ventilation logistic regression**
 558 **models.**

Metric	Full Model with All Variables	Model with LASSO Selected Variables
AIC	112,975.3	118,124.7
BIC	113,393.8	118,520.6

AUC	0.103	0.118
McFadden's R ²	0.310	0.278
Nagelkerke's R ²	0.339	0.306

559 Abbreviations = AIC, Akaike information criterion; BIC, Bayesian information criterion;

560 AUC; area under the curve; LASSO, least absolute shrinkage and selection operator.

561

562 **S4 Table. Model fit statistics for pneumonia logistic regression models.**

Metric	Full Model with All Variables	Model with LASSO Selected
		Variables
AIC	294,942.5	295,800.1
BIC	295,371.1	296,205.5
AUC	0.167	0.169
McFadden's R ²	0.214	0.212
Nagelkerke's R ²	0.255	0.253

563 Abbreviations = AIC, Akaike information criterion; BIC, Bayesian information criterion;

564 AUC; area under the curve; LASSO, least absolute shrinkage and selection operator.

565

566 **S5 Table. Model fit statistics for acute respiratory distress logistic regression**
 567 **models.**

Metric	Full Model with All Variables	Model with LASSO Selected
		Variables
AIC	115,744.8	116,143.1
BIC	116,173.3	116,548.5

AUC	0.164	0.167
McFadden's R^2	0.180	0.178
Nagelkerke's R^2	0.194	0.191

568 Abbreviations = AIC, Akaike information criterion; BIC, Bayesian information criterion;

569 AUC; area under the curve; LASSO, least absolute shrinkage and selection operator.

570

571 **S6 Table. Model fit statistics for mortality logistic regression model.**

Metric	Full Model with All Variables
AIC	346,641.7
BIC	347,098.9
AUC	0.048
McFadden's R^2	0.483
Nagelkerke's R^2	0.531

572 Abbreviations = AIC, Akaike information criterion; BIC, Bayesian information criterion;

573 AUC; area under the curve; LASSO, least absolute shrinkage and selection operator.

574

575 **S7 Table. Hyperparameter tuning results.**

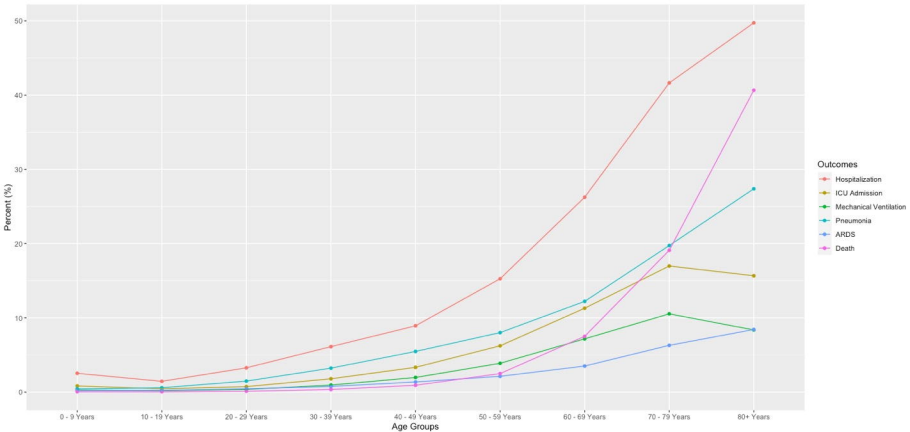
Outcome	mtry	min_n
Hospitalization	3	17
ICU Admission	3	17
Mechanical Ventilation	3	17
Pneumonia	1	15
Acute Respiratory Distress	1	20

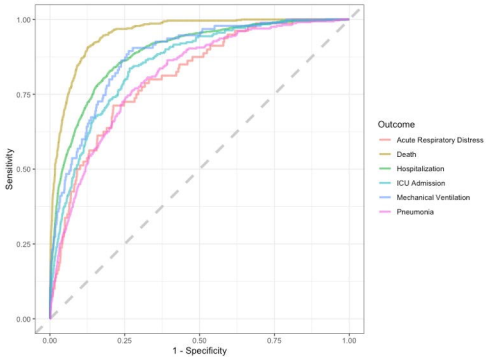
Death

3

17

576 Abbreviations = mtry, number of variables randomly sampled as candidates at each
577 split; min_n, the minimum number of data points in a node that is required for the node
578 to be split further.





Select Gender

Male

Select Age Group

20 - 29 Years

Select Race/Ethnicity

Asian, Non-Hispanic

Select the Region

Northeast

Is the patient a health care worker in the U.S.?

Yes

No

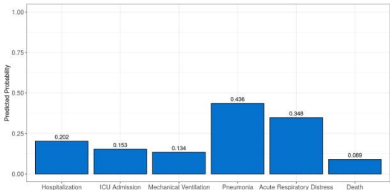
Unknown

Pre-existing medical conditions?

Yes

No

Unknown



Note: Probabilities range from zero (no risk of outcome) to one (certainty of outcome)

