



# Unbiased identification of clinical characteristics predictive of COVID-19 severity

Elliot H. Akama-Garren<sup>1</sup> · Jonathan X. Li<sup>1,2</sup>

Received: 12 April 2021 / Accepted: 28 May 2021 / Published online: 5 June 2021  
© The Author(s), under exclusive licence to Springer Nature Switzerland AG 2021

## Abstract

There is currently limited clinical ability to identify COVID-19 patients at risk for severe outcomes. To unbiasedly identify metrics associated with severe outcomes in COVID-19 patients, we conducted a retrospective study of 835 COVID-19 positive patients at a single academic medical center between March 10, 2020 and October 13, 2020. As of December 1, 2020, 656 (79%) patients required hospitalization and 149 (18%) died. Unbiased comparisons of all clinical characteristics and mortality revealed that abnormal pH (OR 8.54, 95% CI 5.34–13.6), abnormal creatinine (OR 6.94, 95% CI 4.22–11.4), and abnormal PTT (OR 4.78, 95% CI 3.11–7.33) were most significantly associated with mortality. Correlation with ordinal severity scores confirmed these associations, in addition to associations between respiratory rate (Spearman's rho = -0.56), absolute neutrophil count (Spearman's rho = -0.5), and C-reactive protein (Spearman's rho = 0.59) with disease severity. Unsupervised principal component analysis and machine learning model classification of patient demographics, laboratory results, medications, comorbidities, signs and symptoms, and vitals are capable of separating patients on the basis of COVID-19 mortality (AUC 0.82). This retrospective analysis identifies laboratory and clinical metrics most relevant to predict COVID-19 severity.

**Keywords** COVID-19 · Machine learning · Laboratory results · Prediction

## Introduction

As the number of COVID-19 deaths approaches 3.5 million worldwide as of May 11, 2021, there is increasing need to better understand what disease mechanisms and clinical correlates lead to poor outcomes. SARS-CoV-2 infection may result in a spectrum of severity ranging from asymptomatic disease to hospitalization requiring mechanical ventilation [1–7], making identification of patients at risk for severe COVID-19 at initial presentation imperative yet complex. Case series of hospitalized COVID-19 patients during the early pandemic identified key risk groups of severe COVID-19 [8–17], including patients with diabetes, obesity, chronic kidney disease, liver disease, and patients above 65 years old. Cytokine profiling [18] and multi-dimensional flow

cytometry [19–22] have identified hematologic profiles associated with severe COVID-19. Over the course of the pandemic, these advances along with improvements in supportive care such as prone positioning [23–25] have led to reductions in disease mortality [26, 27].

Despite these advances, clinical prediction of COVID-19 prognosis at the time of initial presentation remains imperfect [28]. A better understanding of the clinical correlates of COVID-19 severity would improve prognostic and therapeutic approaches to disease assessment. With an accumulating number of SARS-CoV-2 positive patients with a range of clinical outcomes, we are increasingly able to perform unbiased analyses across more diverse multi-dimensional clinical metrics, in order to identify novel associations with COVID-19 severity. We sought to leverage these data to determine which clinical characteristics are most useful to predict COVID-19 severity. Here, we perform analyses of over 1,700 clinical metrics including laboratory results, vitals, demographics, medications, and disease outcomes in 835 COVID-19 positive patients to identify correlates of disease severity.

✉ Elliot H. Akama-Garren  
elliott\_akama-garren@hms.harvard.edu

<sup>1</sup> Harvard Medical School, Boston, MA 02115, USA

<sup>2</sup> Division of General Medicine, Beth Israel Deaconess Medical Center, Boston, MA 02115, USA

## Methods

### Study design

This study was conducted at the Beth Israel Deaconess Medical Center (BIDMC) in Boston. The BIDMC Institutional Review Board approved this retrospective cohort study (2020P000699) as minimal risk using data collected during routine clinical care and waived the requirement for informed consent. BIDMC patients who presented for care and with confirmed SARS-CoV-2 infection by positive result of nasopharyngeal sample polymerase chain reaction between March 10, 2020 and October 13, 2020, and who had available past medical history, were included.

Data were obtained from the BIDMC COVID-19 Observational Research Effort (CORE) Data Registry REDCap database and BIDMC InSIGHT CORE service. Laboratory values were obtained from inpatient data acquired over the course of an individual patient's admission. When multiple laboratory draws were present over the course of a patient's admission, mean, maximum, and minimum laboratory values for each test collected were calculated for each patient. Time to follow-up was determined by the number of days between the earliest COVID-19 test date and date of death or December 1, 2020, the final date of follow-up, if still alive. COVID-19 severity was graded by the NIH Ordinal Severity Scale. Patients were stratified into eight groups with lower scores corresponding to greater severity: (1) death, (2) invasive mechanical ventilation, (3) noninvasive ventilation, (4) supplemental oxygen, (5) no supplemental oxygen but requiring medical care, (6) no supplemental oxygen and not requiring medical care, (7) limitation in activities, or (8) no limitation in activities.

### Principal component analysis (PCA)

Outcome metrics including mortality, hospitalization length and status, ICU length and status, ventilation and renal replacement therapy requirement, NIH Ordinal Severity Score, pathology results, and medications prescribed after COVID-19 diagnosis were excluded to allow for unsupervised PCA. Patients and metrics with missing data were excluded from analysis, and categorical factor variables were converted to dummy numerical variables. Data were scaled to unit variance and principal component analysis was performed using factextra (version 1.0.7). The top two principal components were used for two-dimensional mapping of patient data and variable eigenvectors.

### Machine learning classification

Mortality status was added to the data set used for PCA to allow for construction of a supervised machine learning

classifier. All machine learning analyses were performed in R (version 3.6.1). Training and test data sets were created using the createDataPartition function in caret (version 6.0), with 75% of patients allocated to the training data set. Training data were preprocessed by centering and scaling and training was performed using ten separate tenfold repeated cross-validations for resampling. A gradient boosting machine model [29, 30] was built using 100 trees, a tree complexity of 2, and a learning rate of 0.1 using the train function in caret. Training performance was measured using area under the ROC curve, and variable importance was calculated using the varImp function in caret. Model performance was tested on the test data set and evaluated using MLevel (version 0.3).

### Statistical analysis

All statistical analyses were performed in R (version 3.6.1). Bar graphs and violin plots were created using ggpubr (version 0.4.0), correlation plots were created using corrplot (version 0.84), Kaplan–Meier plots were created using survminer (version 0.4.8) and survival (version 3.2–7), and scatter plots and forest plots were created using ggplot2 (version 3.3.0). Heatmaps and hierarchical clustering were performed using pheatmap (version 1.0.12). Volcano plots were generated using EnhancedVolcano (version 1.4.0), and significant differences (absolute logFC > 0.2 and P-val < 0.05) were highlighted in red. When data were missing, these patients were not included in a given univariate analysis, eliminating potential confounding due to the presence or absence of a given clinical metric. When multiple comparisons were made, p values were corrected by the Benjamini–Hochberg procedure and a false discovery rate < 0.05 was considered significant.

## Results

### Demographics, comorbidities, and outcomes of COVID-19 patients

A total of 835 patients with PCR confirmed SARS-CoV-2 infection were included (Table 1). The median age was 64 years (IQR, 50–76 years; range, 17–102 years) and 438 (52%) were female. Of these patients, 363 (43%) were white and 253 (30%) were black. Past medical history was available for 549 patients and among these patients, common comorbidities included hypertension (347; 63%), diabetes (224; 41%), obesity (157; 30%), chronic kidney disease (144; 26%), and cancer (131; 24%). Active prescriptions at time of COVID-19 diagnosis were available for 697 patients, and among these the most common categories of prescribed drugs included

**Table 1** Demographics, comorbidities, and outcomes of COVID-19 patients

	Overall (N = 835)	Alive (N = 686)	Dead (N = 149)	P value
Gender				0.117
Female	438 (52%)	369 (54%)	69 (46%)	
Male	397 (48%)	317 (46%)	80 (54%)	
Age	64 (50–76)	61 (47–73)	73 (63–84)	<0.001
Race				0.0919
Native American	1 (0%)	1 (0%)	0 (0%)	
Asian	31 (4%)	23 (3%)	8 (5%)	
Black	253 (30%)	206 (30%)	47 (32%)	
Declined	1 (0%)	1 (0%)	0 (0%)	
Native Hawaiian	2 (0%)	1 (0%)	1 (1%)	
Other	90 (11%)	84 (12%)	6 (4%)	
Unknown	94 (11%)	73 (11%)	21 (14%)	
White	363 (43%)	297 (43%)	66 (44%)	
ABO Type				0.487
A	71 (9%)	43 (6%)	28 (19%)	
AB	11 (1%)	8 (1%)	3 (2%)	
B	40 (5%)	29 (4%)	11 (7%)	
O	104 (12%)	63 (9%)	41 (28%)	
Missing	609 (72.9%)	543 (79.2%)	66 (44.3%)	
BMI	29 (25–34)	29 (25–34)	30 (24–36)	0.905
Comorbidities available	549 (66%)	459 (67%)	90 (60%)	
Hypertension	347 (63%)	278 (61%)	69 (77%)	0.00549
Chronic kidney disease	144 (26%)	107 (23%)	37 (41%)	<0.001
Diabetes	224 (41%)	172 (37%)	52 (58%)	<0.001
Obesity	167 (30%)	136 (30%)	31 (34%)	0.434
Rheumatologic disease	127 (23%)	100 (22%)	27 (30%)	0.12
Autoimmune disease	49 (9%)	43 (9%)	6 (7%)	0.535
Cancer	131 (24%)	98 (21%)	33 (37%)	0.00287
Immunosuppressive Disease	128 (23%)	103 (22%)	25 (28%)	0.338
COPD	72 (13%)	54 (12%)	18 (20%)	0.0517
Asthma	81 (15%)	66 (14%)	15 (17%)	0.691
Coronary artery disease	130 (24%)	97 (21%)	33 (37%)	0.00241
Cerebrovascular disease	67 (12%)	46 (10%)	21 (23%)	<0.001
Medications available	697 (83%)	568 (83%)	129 (87%)	
Corticosteroid	179 (26%)	140 (25%)	39 (30%)	0.231
Calcineurin inhibitors	16 (2%)	12 (2%)	4 (3%)	0.726
Antirheumatic therapy	9 (1%)	7 (1%)	2 (2%)	1
Immunosuppressive therapy	46 (7%)	32 (6%)	14 (11%)	0.0361
Chemotherapy	26 (4%)	19 (3%)	7 (5%)	0.385
Antiglycemic therapy	241 (35%)	192 (34%)	49 (38%)	0.424
Asthma therapy	227 (33%)	178 (31%)	49 (38%)	0.177
Biologics	1 (0%)	1 (0%)	0 (0%)	1
Osteoporosis therapy	13 (2%)	9 (2%)	4 (3%)	0.43
Antihypertensive therapy	500 (72%)	392 (69%)	108 (84%)	0.00119
Labs				
Absolute lymphocyte count (10 <sup>6</sup> /mL)	1.2 (0.83–1.6)	1.2 (0.89–1.6)	0.97 (0.67–1.3)	<0.001
C-Reactive protein (mg/L)	94 (52–150)	82 (42–130)	140 (100–180)	<0.001
Creatinine (mg/dL)	1.0 (0.73–1.7)	0.91 (0.70–1.3)	1.8 (1.1–2.8)	<0.001
Ferritin (ng/mL)	680 (300–1500)	570 (260–1200)	1400 (570–2900)	<0.001

**Table 1** (continued)

	Overall (N = 835)	Alive (N = 686)	Dead (N = 149)	P value
D-Dimer (ng/mL FEU)	1300 (720–2700)	1100 (650–2200)	2400 (1200–4800)	< 0.001
Creatine kinase (IU/L)	150 (69–380)	140 (65–360)	170 (82–530)	0.0455
INR	1.2 (1.1–1.4)	1.2 (1.1–1.3)	1.3 (1.2–1.5)	< 0.001
Lactate dehydrogenase (IU/L)	330 (260–430)	320 (240–400)	420 (310–560)	< 0.001
pH	7.1 (6.7–7.3)	7.0 (6.5–7.3)	7.2 (7.0–7.3)	0.00105
Platelet count (10 <sup>6</sup> /mL)	230 (180–310)	250 (190–320)	190 (140–260)	< 0.001
PT (s)	13 (12–15)	13 (12–15)	14 (13–17)	< 0.001
PTT (s)	35 (30–55)	33 (29–47)	53 (35–70)	< 0.001
Absolute neutrophil count (10 <sup>6</sup> /mL)	5.5 (3.8–8.2)	5.0 (3.6–7.3)	7.8 (5.1–12)	< 0.001
A1c (%)	7.7 (6.4–9.3)	7.6 (6.3–9.3)	7.8 (7.2–8.8)	0.629
<b>Vitals</b>				
Respiratory rate	19 (18–21)	19 (18–20)	23 (20–25)	< 0.001
Heart rate	85 (76–94)	84 (74–93)	89 (82–98)	< 0.001
Tmax	100 (99–100)	100 (99–100)	100 (100–100)	< 0.001
SBP (minimum)	99 (91–110)	99 (92–110)	95 (84–110)	0.0275
DBP (minimum)	57 (49–65)	57 (50–65)	54 (44–63)	0.00543
<b>Status</b>				
Inpatient	656 (79%)	510 (74%)	146 (98%)	< 0.001
Outpatient	179 (21%)	176 (26%)	3 (2%)	
<b>Outcomes</b>				
Supplemental O <sub>2</sub>	336 (40%)	263 (38%)	73 (49%)	0.0208
Mechanical ventilation	196 (23%)	106 (15%)	90 (60%)	< 0.001
Total encounters	1.0 (1.0–1.0)	1.0 (1.0–1.0)	1.0 (1.0–1.0)	0.261
Length admission	9.0 (5.0–18)	8.0 (4.0–19)	12 (7.0–17)	< 0.001
Ordinal score	4.0 (2.0–5.0)	4.0 (4.0–5.0)	1.0 (1.0–1.0)	< 0.001
ICU admission	133 (16%)	87 (13%)	46 (31%)	< 0.001
ICU days	8.0 (3.0–17)	8.0 (2.0–19)	9.0 (4.0–14)	< 0.001

Continuous data presented as mean (95% CI). P values computed by Chi-squared test for categorical data and Wilcoxon signed-rank test for continuous data

antihypertensive drugs (500; 72%), antihistamines (324; 46%), and antiglycemic drugs (241; 35%). Most patients had an elevated temperature (median Tmax 100; IQR 99–100) and were tachypneic (median 19; IQR 18–21) but had normal heart rates (median 85; IQR 76–94). As of December 1, 2020, 656 (79%) patients required hospitalization, 336 (40%) required supplemental oxygen, 310 (37%) required intensive care unit (ICU) stays, and 196 (23%) required mechanical ventilation. Among patients who were hospitalized the median total length of stay was 9 days (IQR, 2–5 days) and among patients treated in the ICU the median length of stay in the ICU was 8 days (IQR, 3–17 days). NIH Ordinal Scoring was available for 322 patients, and mean ordinal score was 3.7 (SD 1.7). Overall, 149 (18%) patients died at the time of censoring.

### Clinical predictors of COVID-19 outcomes

To validate our ability to identify risk factors for COVID-19 severity, we compared mortality rates among currently recognized comorbidities for COVID-19 (Fig. 1A). In our cohort, hypertension (OR 2.14, 95% CI 1.27–3.60), chronic kidney disease (OR 2.30, 95% CI 1.44–3.68), cardiovascular disease (OR 2.73, 95% CI 1.54–4.84), diabetes (OR 2.28, 95% CI 1.44–3.60), coronary artery disease (OR 2.16, 95% CI 1.33–3.50), and cancer (OR 2.13, 95% CI 1.32–3.45) were associated with COVID-19 mortality. Risks for hospitalization included hypertension (OR 2.42, 95% CI 1.64–3.57), male gender (OR 1.69, 95% CI 1.21–2.38), diabetes (OR 2.17, 95% CI 1.43–3.29), chronic kidney disease (OR 2.42, 95% CI 1.44–3.68), coronary artery disease (OR 2.59, 95%

CI 1.51–4.42), and COPD (OR 3.63, 95% CI 1.65–7.96), whereas risks for ICU admission only included male gender (OR 2.17, 95% CI 1.42–3.31) and diabetes (OR 2.27, 95% CI 1.35–3.81). Notably, male gender was not significantly associated with mortality among COVID-19 patients in our cohort (OR 1.35, 95% CI 0.95–1.93).

In order to unbiasedly compare the relative association of clinical characteristics with COVID-19 outcomes, we calculated the odds ratios among binary categorical clinical metrics measured, including laboratory results, demographics, medications, comorbidities, and signs and symptoms (Fig. 1B). Mortality was most significantly associated with abnormal pH (OR 8.54, 95% CI 5.34–13.6), abnormal creatinine (OR 6.94, 95% CI 4.22–11.4), and abnormal PTT (OR 4.78, 95% CI 3.11–7.33). Hospitalization was most significantly associated with abnormal D-dimer (OR 8.87, 95% CI 4.18–18.8), NSAID use (OR 0.24, 95% CI 0.15–0.38), and abnormal C-reactive protein (OR 6.43, 95% CI 3.30–12.5), and ICU admission was associated with requiring supplemental oxygen at admission (OR 8.34, 95% CI 4.91–14.1), abnormal pH (OR 13.1, 95% CI 7.71–22.5), and abnormal PTT (OR 7.36, 95% CI 4.42–12.2).

We next sought to compare the relative association between continuous variables and COVID-19 outcomes. Mann–Whitney U tests between mortality and laboratory values and demographic information revealed that elevated creatinine was most significantly associated with mortality (average maximum creatinine 3.97 in dead vs 1.97 in alive, adjusted  $P$ -val  $< 2 \times 10^{-16}$ ) (Fig. 1C). Other significant associations with mortality included decreased albumin (average minimum albumin 2.50 in dead vs 3.22 in alive, adjusted  $P$ -val  $< 2 \times 10^{-16}$ ), decreased lymphocyte count (average minimum lymphocytes 7.53 in dead vs 13.56 in alive, adjusted  $P$ -val  $< 2 \times 10^{-16}$ ), elevated phosphate (average maximum phosphate 6.50 in dead vs 4.61 in alive, adjusted  $P$ -val  $< 2 \times 10^{-16}$ ), and older age (average age 71.9 years in dead vs 59.5 in alive, adjusted  $P$ -val  $= 8.6 \times 10^{-16}$ ) (Fig. 1D). Comparisons in hospitalization, ventilation, oxygen requirement, and ICU admission patient groups revealed similar associations between abnormal creatinine, albumin, lymphocytes, and phosphate and COVID-19 outcomes (Fig. 1C). These results suggest that laboratory abnormalities might be more informative in predicting outcomes from COVID-19 than patient demographic information including comorbidities.

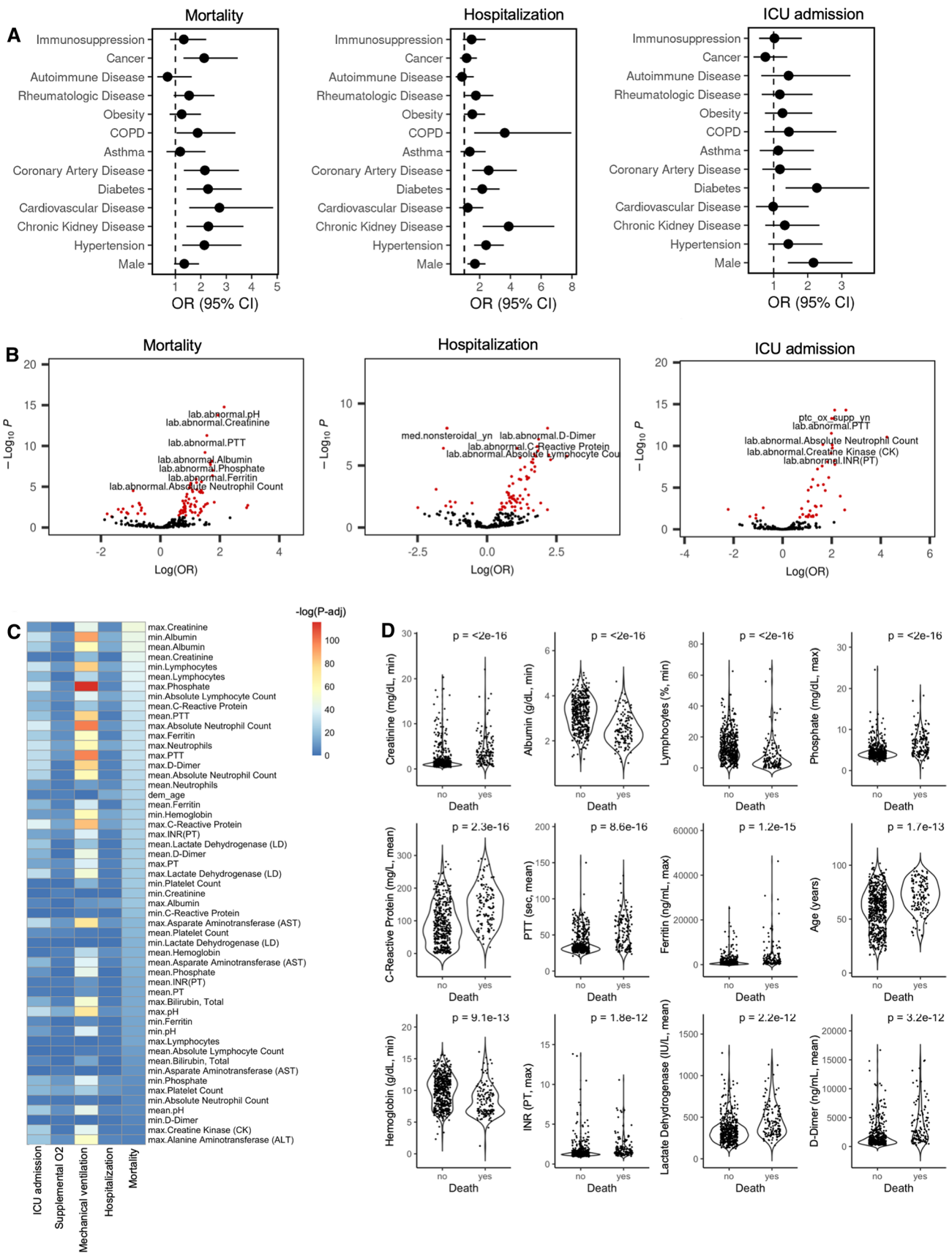
To quantify and rank the effects of clinical metrics on time to death following COVID-19 diagnosis, we performed Kaplan–Meier analysis of patient survival using positive COVID-19 test date and date of death. Among the 149 (18%) of patients that died, the median survival time after COVID-19 diagnosis was 13 days (IQR, 7–28 days) (Fig. 2A). Regression analysis of demographics, laboratory results, medications, comorbidities, and vitals against

survival probability revealed that abnormal pH (HR 6.5, 95% CI 4.2–10), stratified age groups (HR = 1.5, 95% CI 1.3–1.7), abnormal albumin (HR 3.6, 95% CI 2.4–5.5), and abnormal phosphate (HR 4.7, 95% CI 2.7–8.1) were most significantly associated with increased risk of COVID-19 death (Fig. 2B). These risks are greater than those associated with currently accepted comorbidities for severe COVID-19 in our cohort, such as hypertension (HR 2.0, 95% CI 1.2–3.3), diabetes (HR 2.1, 95% CI 1.4–3.3), and chronic kidney disease (HR 2.2, 95% CI 1.4–3.3) (Fig. 2C). Both race (HR 0.99, 95% CI 0.92–1.1) and gender (HR 1.3, 95% CI 0.91–1.7) were not significantly associated with decreased survival following COVID-19 diagnosis in our cohort.

### Clinical correlates of COVID-19 severity

To examine associations between clinical metrics and COVID-19 severity beyond binary categorical outcomes, we measured the correlation of each metric with NIH ordinal severity scores and total length of stay per patient (Fig. 3A). Ordinal score was most significantly correlated with maximum respiratory rate (Spearman's  $\rho = -0.56$ ), maximum absolute neutrophil count (Spearman's  $\rho = -0.5$ ), maximum C-reactive protein (Spearman's  $\rho = -0.52$ ), and minimum albumin (Spearman's  $\rho = 0.5$ ) (Fig. 3B). The total length of admission was most significantly correlated with maximum temperature (Spearman's  $\rho = 0.62$ ), maximum phosphate (Spearman's  $\rho = 0.60$ ), minimum hemoglobin (Spearman's  $\rho = -0.58$ ), and minimum systolic blood pressure (Spearman's  $\rho = -0.53$ ) (Fig. 3C). These results confirm our previous findings, suggesting that hematologic laboratory results are not only indicative of mortality in COVID-19 patients, but are also correlated with disease severity. These results also quantify the relative association of vitals such as respiratory rate and temperature with COVID-19 severity.

To determine relationships between multiple categorical and numerical outcomes and metrics, we performed correlation analysis across patient demographics, selected laboratory results, medications, comorbidities, vitals, and outcomes including continuous metrics of COVID-19 severity (Fig. 4). In addition to the associations noted previously, this analysis revealed significant correlations between COVID-19 outcomes and clinical interventions such as ICU admission and mechanical ventilation. As expected, comorbidities were highly correlated with prescriptions for appropriate medications (e.g., diabetes and antidiabetic drugs) as well as corresponding laboratory results (e.g., chronic kidney disease and mean creatinine). Notably, comorbidities were more closely associated with corresponding medications than COVID-19 outcomes, whereas laboratory values and vitals were more closely associated with COVID-19 outcomes than corresponding comorbidities. Overall, this



**Fig. 1** Univariate analyses identify key laboratory parameters associated with mortality in COVID-19 patients. **a** Forest plot comparing odds ratios of selected comorbidities with mortality, hospitalization, and ICU admission in COVID-19 patients. Horizontal lines indicate 95% CI. **b** Volcano plots of odds ratios of laboratory results, demographics, medications, comorbidities, and signs and symptoms with mortality, hospitalization, and ICU admission in COVID-19 patients. P values corrected for multiple comparisons by Benjamini–Hochberg procedure and significant metrics ( $P\text{-adj} < 0.05$ ) indicated in red. **c** Heatmap of adjusted p values from Mann–Whitney U tests for continuous laboratory values and demographic information between patients requiring or not requiring ICU admission, supplement oxygen, mechanical ventilation, hospitalization, and death. Metrics significantly altered between alive and dead patient cohorts are shown and arranged by increasing adjusted p value. **d** Violin plots of the most significantly altered clinical metrics alive and dead patient cohorts. Mann–Whitney U test p value shown

correlation analysis revealed the heterogeneity of COVID-19 patient presentation, and the relative utility of a spectrum of patient information in predicting COVID-19 severity.

### Principal component analysis and machine learning classification segregates COVID-19 patients by mortality

To determine whether COVID-19 patients can be stratified by severity based on clinical metrics typically present at admission to the emergency department, we performed unsupervised principal component analysis (PCA). We excluded metrics of COVID-19 outcomes and severity and metrics that would not be known at admission, such as pathology results and medications placed after COVID-19 diagnosis. Only patients for whom full demographic, laboratory, medication history, comorbidities, past medical history, and vitals were available were included, leaving 237 metrics across 209 patients. PCA distilled these 237 metrics into two dimensions, which were most defined by immunosuppression and anemia in Dimension 1, and by AST, LDH, ALT, and ferritin in Dimension 2 (Fig. 5A). The eigenvectors for mean AST and maximum ferritin were orthogonal to the eigenvector for immunosuppression (Fig. 5B), suggesting that these metrics capture independent meta-characteristics of COVID-19 patients.

We next plotted the 209 patients present in our PCA in two-dimensional space. There was no clear distribution of COVID-19 patients in PCA space on the basis of demographic information such as gender, race, and age (Fig. 5C). However, when we visualized mortality, which was not a variable included in our PCA, there was a separation among COVID-19 patients in PCA space. Similar trajectories could be appreciated for COVID-19 severity and outcomes metrics, such as length of stay, mechanical ventilation requirement, and ordinal score (Fig. 5C). Trajectories of COVID-19 severity in PCA space were orthogonal to the eigenvector for immunosuppression, suggesting

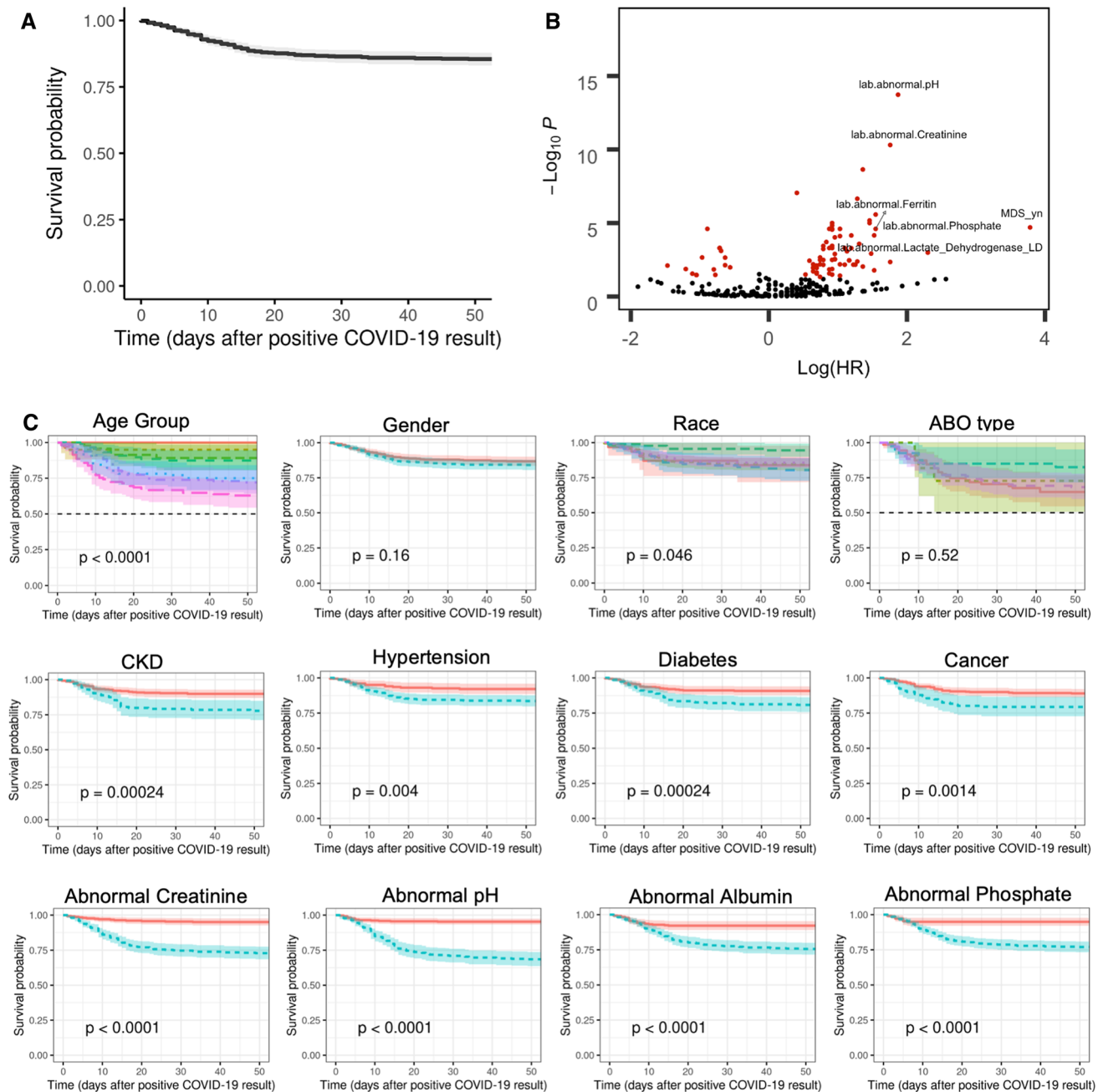
that although immunosuppression contributes to variability among COVID-19 patients, it likely does not contribute to disease severity.

Given our ability to segregate patients by COVID-19 severity using unsupervised PCA, we next sought to design a machine learning classifier to predict patient mortality. Using mortality in addition to the 237 variables used for PCA above, we partitioned our COVID-19 patient cohort into a training set of 157 patients and a test set of 52 patients. The training set of patients was used to build a supervised gradient boosting machine model to classify patient mortality. Our model achieved a sensitivity of 0.53 (95% CI 0.39–0.67), specificity of 0.88 (95% CI 0.81–0.93), and area under curve (AUC) for the ROC curve of 0.87 (95% CI 0.80–0.94) based on the training data (Fig. 5D). When applied to the test set, our model correctly identified 6 of 15 patients who died following COVID-19 diagnosis, achieving an accuracy of 0.77 (95% CI 0.63–0.87), a sensitivity of 0.92, specificity of 0.40, and AUC ROC of 0.82. Variable importance scores extracted from the gradient boosting machine model revealed that absolute neutrophil count, PTT, and patient age were the most contributory to model prediction (Fig. 5E). Together our PCA and machine learning classifier suggest that COVID-19 severity and outcomes can be correlated with clinical characteristics known at the time of admission and confirm the importance of laboratory data over demographic information in predicting disease outcome.

## Discussion

Here, we unbiasedly profile over 1700 unique clinical metrics in 835 COVID-19 patients to identify correlates of disease outcomes and severity. We observed similar odds ratios for COVID-19 mortality risk from comorbidities previously reported, such as increased age [11, 17, 31–33], hypertension [12], diabetes [8, 11–13], and chronic kidney disease [16]. Univariate, correlation, and multivariate analyses revealed strong associations between key laboratory parameters and COVID-19 severity. Several of these associations have been previously reported, such as elevated creatinine [34], decreased lymphocyte count [19, 20], elevated CRP [34], decreased hemoglobin [20], abnormal pH [35], decreased albumin [36], and elevated PTT [20]. Notably, through unbiased comparisons across all clinical metrics, we observed that these laboratory abnormalities are more strongly associated with mortality in COVID-19 patients than patient age, gender, comorbidities, or prescribed medications.

As this was a retrospective cohort study of associations with COVID-19 outcomes, it remains unclear whether the metrics identified here predispose patients to worse outcomes or are a consequence of severe COVID-19 itself.



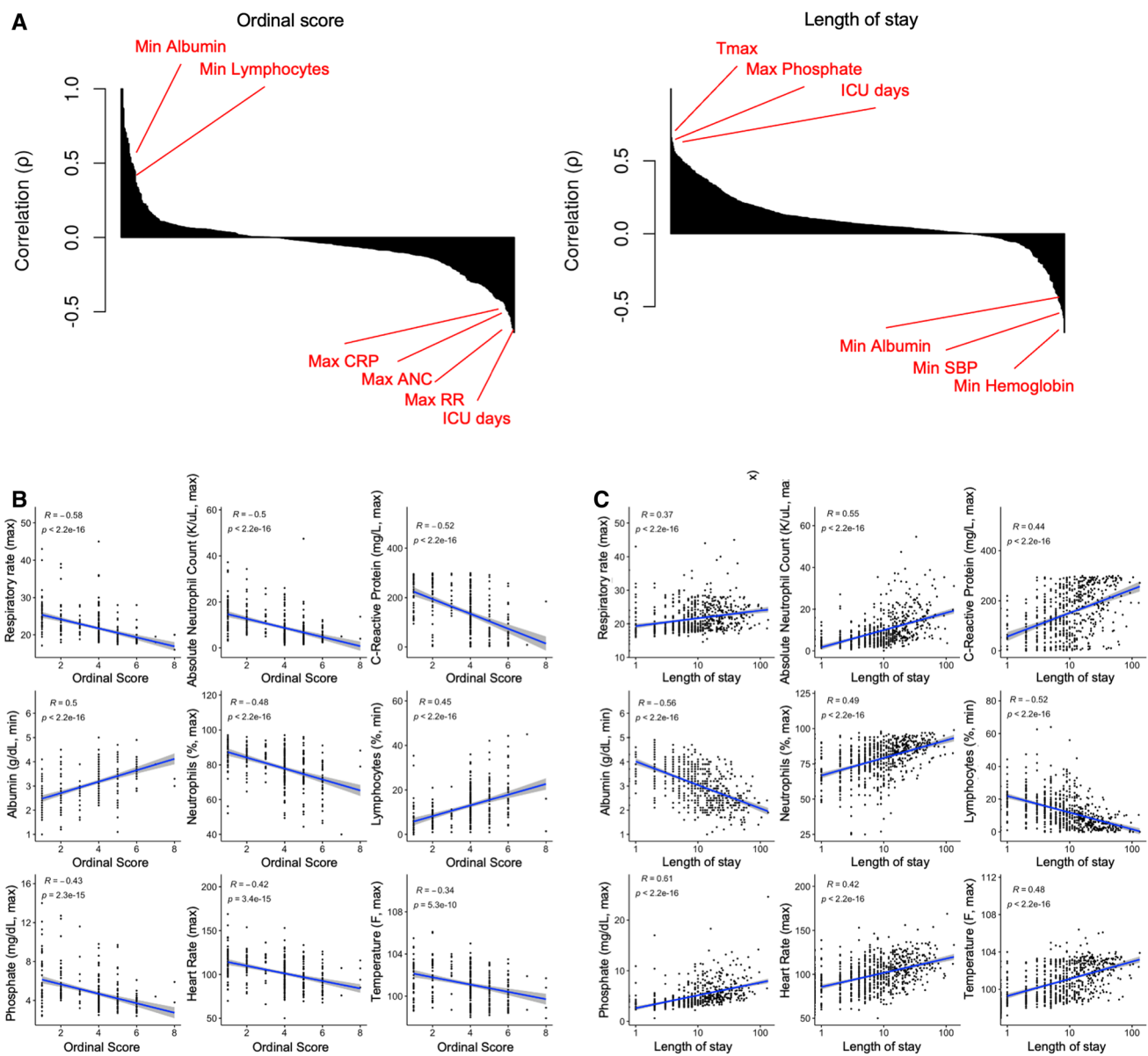
**Fig. 2** Unbiased identification of metrics most associated with increased risk of dying following COVID-19 diagnosis. **A** Kaplan–Meier plot of patient survival following COVID-19 diagnosis. **B** Volcano plot of hazard ratios (HR) calculated from unbiased Cox regression analysis between all measured patient metrics and patient survival following COVID-19 diagnosis. P values were calculated

using the Wald test statistic and corrected for multiple comparisons by Benjamini–Hochberg procedure. Significant metrics ( $P\text{-adj} < 0.05$ ) indicated in red. **C** Kaplan–Meier plots of patient survival following COVID-19 diagnosis stratified by indicated patient demographic or laboratory result. Log rank test p value indicated on plots and 95% CI indicated by shading

Abnormal pH and increased respiratory rate in patients with severe COVID-19 is likely reflective of the eventual acute respiratory distress syndrome and tissue malperfusion experienced by these patients [5], whereas the elevated inflammatory markers we observed are characteristic of the systemic inflammation observed in some case of severe COVID-19

[3, 37, 38]. Some laboratory perturbations such as prolonged PTT might reflect interventions employed preferentially in COVID-19 patients such as anticoagulants. Other laboratory parameters such as decreased lymphocytes and albumin might represent a unique inflammatory phenotype that predisposes patients to severe COVID-19 [19]. Regardless of





**Fig. 3** Correlation between continuous clinical metrics and COVID-19 severity. **A** Ranked order plots of Spearman correlation coefficients between all clinical metrics and NIH ordinal score and total length of admission. Selected significant associations indicated on

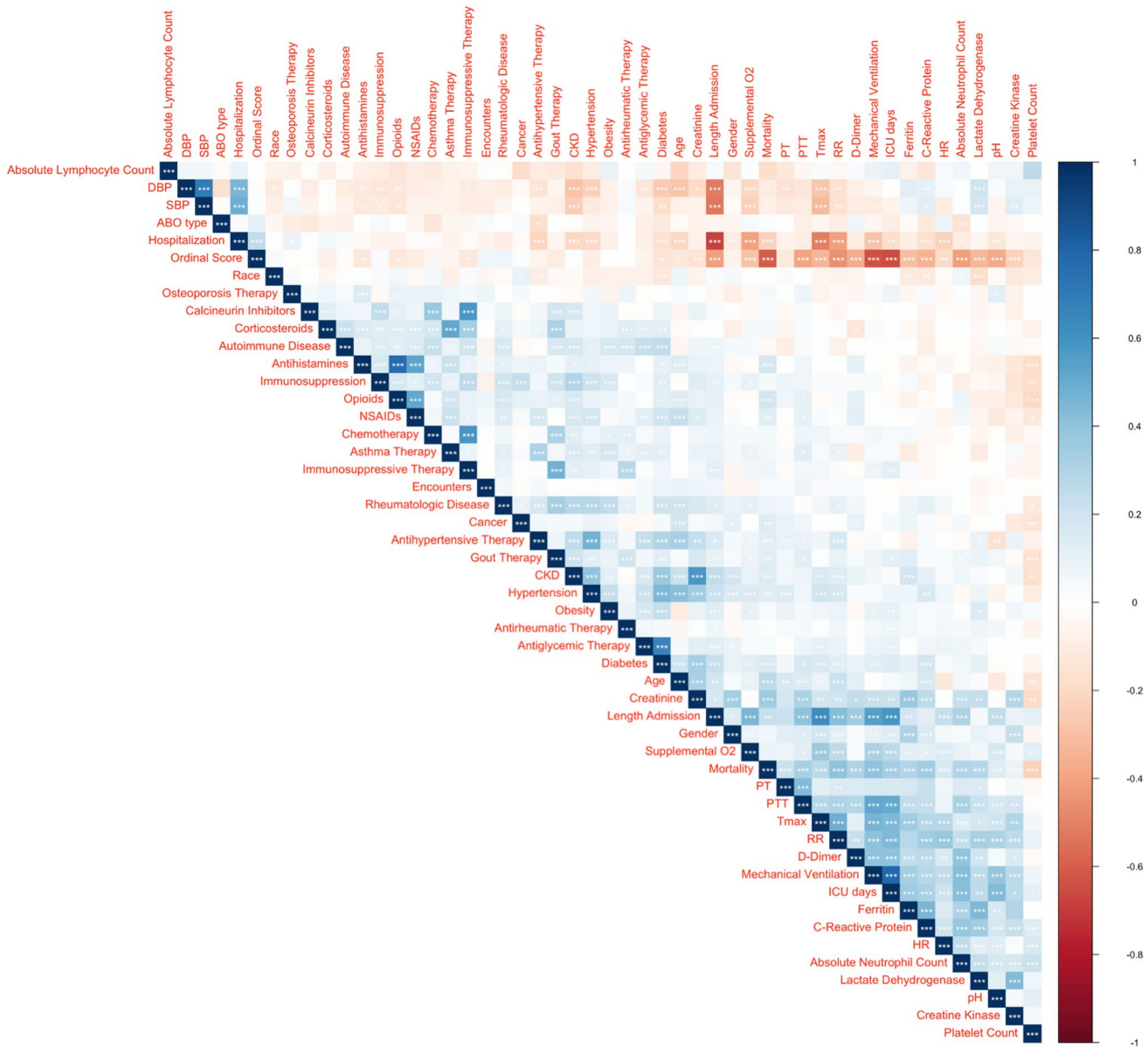
plot. **B–C** Scatter plots of correlation of selected clinical metrics and NIH ordinal score (**B**) or length of admission (**C**). Spearman correlation coefficient and p value indicated on plot, and regression line and 95% confidence interval indicated in blue

the root cause of the clinical associations we describe, we have identified key clinical metrics that may be obtained at emergency department admission to identify overall risk for COVID-19 mortality.

We observed a mortality rate of 18% and hospitalization rate of 79%, in contrast to currently estimated case fatality rates of 0.9–7.2% [17, 33, 39, 40] for SARS-CoV-2. This is likely due to sampling bias as only patients who sought care at an academic medical center, obtained a laboratory confirmed COVID-19 diagnosis, and had available medication or past medical history were included. Alternatively, this

might reflect the evolving mortality rate of the course of this pandemic, as our ability to diagnose and treat COVID-19 has improved the past year [41]. Nevertheless, a range of clinical presentations and disease severity scores are represented in our patient cohort, including outpatients and patients with asymptomatic disease.

COVID-19 remains a great threat to society relative to other respiratory viral diseases due to its case fatality rate and its striking range of clinical presentations and severity [17, 42, 43]. This study offers an unbiased retrospective approach to identify potential associations

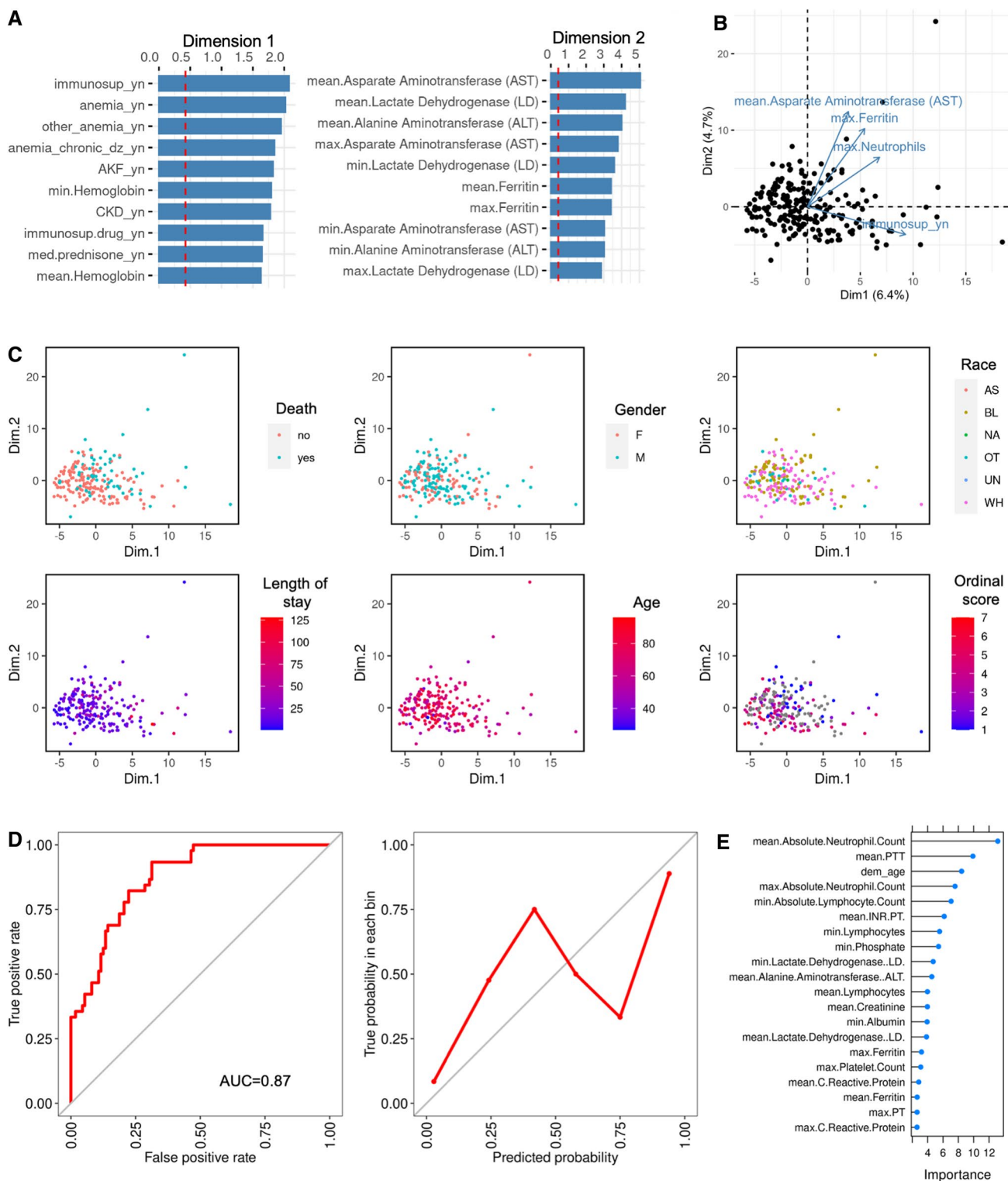


**Fig. 4** Correlation analysis reveals heterogeneity and associations among COVID-19 patient characteristics and outcomes. Correlation plot of Spearman correlation coefficients between indicated

clinical metrics and measures of disease outcomes among COVID-19 patients. Matrix display order was determined by angular order of eigenvectors. \* $P < 0.05$ , \*\* $P < 0.01$ , \*\*\* $P < 0.001$

with this fatality rate and spectrum of disease severity. Our data suggest that increased absolute neutrophil count, decreased albumin, and decreased lymphocytes are key correlates of severe COVID-19 and are clinical characteristics available at initial admission that might be informative of disease prognosis. By identifying which

COVID-19 patients are most at risk for severe disease, we may be better able to provide early and targeted therapeutic interventions, thereby combatting the current pandemic in an orthogonal but complementary approach to the preventative approaches currently being pursued across the world.



**Fig. 5** Multivariate analyses segregate COVID-19 patients by disease severity. **A** Bar plot indicating contributions of the top ten metrics to the top two principal components identified by unsupervised principal component analysis (PCA) of COVID-19 patients. **B** Biplot of principle component scores of COVID-19 patients (dots) and variable loadings (vectors). The top four metrics with the greatest contribution to variability are shown. **C** PCA plots of COVID-19 patients according to the top two principal components and colored according to the

indicated metric. **D** Receiver operator curve (left) and calibration plot (right) to assess ability of a supervised gradient boosting machine model to classify COVID-19 patient mortality using demographic, laboratory, medication history, comorbidities, past medical history, and vitals. Classification performance assessed by area under the curve (AUC). **E** Ranked plot of the importance scores of top 20 clinical metrics in the machine learning classifier constructed in (**D**)

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1007/s10238-021-00730-y>.

**Acknowledgements** We would like to thank the BIDMC COVID-19 Observational Research Effort (CORE) Data Registry REDCap database and BIDMC InSIGHT CORE service for help with data acquisition.

**Funding** This work was supported by the National Institutes of Health [T32GM007753 to E.A.G.].

## Declarations

**Conflicts of interest** The authors declare no conflicts of interest.

**Availability of data and material** De-identified patient data are available through GitHub ([https://github.com/egarren/COVID\\_ML](https://github.com/egarren/COVID_ML)).

**Code availability** Relevant code is available through GitHub ([https://github.com/egarren/COVID\\_ML](https://github.com/egarren/COVID_ML)).

**Ethics approval** The BIDMC Institutional Review Board approved this retrospective cohort study (2020P000699) as minimal risk using data collected during routine clinical care and waived the requirement for informed consent.

**Consent to participate** Not applicable.

**Consent for publication** Not applicable.

## References

- Chan JF-W, Yuan S, Kok K-H, et al. A familial cluster of pneumonia associated with the 2019 novel coronavirus indicating person-to-person transmission: a study of a family cluster. *Lancet Lond Engl.* 2020;395(10223):514–23. [https://doi.org/10.1016/S0140-6736\(20\)30154-9](https://doi.org/10.1016/S0140-6736(20)30154-9).
- Bajema KL, Oster AM, McGovern OL, et al. Persons evaluated for 2019 novel coronavirus—United States, January 2020. *MMWR Morb Mortal Wkly Rep.* 2020;69(6):166–70. <https://doi.org/10.15585/mmwr.mm6906e1>.
- Huang C, Wang Y, Li X, et al. Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China. *Lancet Lond Engl.* 2020;395(10223):497–506. [https://doi.org/10.1016/S0140-6736\(20\)30183-5](https://doi.org/10.1016/S0140-6736(20)30183-5).
- Chen N, Zhou M, Dong X, et al. Epidemiological and clinical characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. *Lancet Lond Engl.* 2020;395(10223):507–13. [https://doi.org/10.1016/S0140-6736\(20\)30211-7](https://doi.org/10.1016/S0140-6736(20)30211-7).
- Wang D, Hu B, Hu C, et al. Clinical characteristics of 138 hospitalized patients with 2019 novel coronavirus-infected pneumonia in Wuhan, China. *JAMA.* 2020;323(11):1061–9. <https://doi.org/10.1001/jama.2020.1585>.
- Liu K, Fang Y-Y, Deng Y, et al. Clinical characteristics of novel coronavirus cases in tertiary hospitals in Hubei Province. *Chin Med J (Engl).* 2020;133(9):1025–31. <https://doi.org/10.1097/CM9.0000000000000744>.
- Yang X, Yu Y, Xu J, et al. Clinical course and outcomes of critically ill patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered, retrospective, observational study. *Lancet Respir Med.* 2020;8(5):475–81. [https://doi.org/10.1016/S2213-2600\(20\)30079-5](https://doi.org/10.1016/S2213-2600(20)30079-5).
- Harrison SL, Fazio-Eynullayeva E, Lane DA, Underhill P, Lip GYH. Comorbidities associated with mortality in 31,461 adults with COVID-19 in the United States: a federated electronic medical record analysis. *PLoS Med.* 2020;17(9):e1003321. <https://doi.org/10.1371/journal.pmed.1003321>.
- Tartof SY, Qian L, Hong V, et al. Obesity and mortality among patients diagnosed with COVID-19: results from an integrated health care organization. *Ann Intern Med.* 2020;173(10):773–81. <https://doi.org/10.7326/M20-3742>.
- Lighter J, Phillips M, Hochman S, et al. Obesity in patients younger than 60 years is a risk factor for COVID-19 hospital admission. *Clin Infect Dis Off Publ Infect Dis Soc Am.* 2020;71(15):896–7. <https://doi.org/10.1093/cid/ciaa415>.
- Williamson EJ, Walker AJ, Bhaskaran K, et al. Factors associated with COVID-19-related death using OpenSAFELY. *Nature.* 2020;584(7821):430–6. <https://doi.org/10.1038/s41586-020-2521-4>.
- Zhou F, Yu T, Du R, et al. Clinical course and risk factors for mortality of adult inpatients with COVID-19 in Wuhan, China: a retrospective cohort study. *Lancet Lond Engl.* 2020;395(10229):1054–62. [https://doi.org/10.1016/S0140-6736\(20\)30566-3](https://doi.org/10.1016/S0140-6736(20)30566-3).
- CDC COVID-19 Response Team (2020) Preliminary estimates of the prevalence of selected underlying health conditions among patients with coronavirus disease 2019—United States, February 12–March 28, 2020. *MMWR Morb Mortal Wkly Rep.* 2020;69(13):382–386. <https://doi.org/10.15585/mmwr.mm6913e2>.
- Liang W, Guan W, Chen R, et al. Cancer patients in SARS-CoV-2 infection: a nationwide analysis in China. *Lancet Oncol.* 2020;21(3):335–7. [https://doi.org/10.1016/S1470-2045\(20\)30096-6](https://doi.org/10.1016/S1470-2045(20)30096-6).
- Cunningham JW, Vaduganathan M, Claggett BL, et al. Clinical outcomes in young US adults hospitalized with COVID-19. *JAMA Intern Med.* 2020. <https://doi.org/10.1001/jamainternmed.2020.5313>.
- Petrilli CM, Jones SA, Yang J, et al. Factors associated with hospital admission and critical illness among 5279 people with coronavirus disease 2019 in New York City: prospective cohort study. *BMJ.* 2020;369:m1966. <https://doi.org/10.1136/bmj.m1966>.
- Wu Z, McGoogan JM. Characteristics of and important lessons from the coronavirus disease 2019 (COVID-19) outbreak in China: summary of a report of 72 314 cases from the Chinese center for disease control and prevention. *JAMA.* 2020;323(13):1239–42. <https://doi.org/10.1001/jama.2020.2648>.
- Del Valle DM, Kim-Schulze S, Huang H-H, et al. An inflammatory cytokine signature predicts COVID-19 severity and survival. *Nat Med.* 2020;26(10):1636–43. <https://doi.org/10.1038/s41591-020-1051-9>.
- Mathew D, Giles JR, Baxter AE, et al. Deep immune profiling of COVID-19 patients reveals distinct immunotypes with therapeutic implications. *Science.* 2020. <https://doi.org/10.1126/science.abc8511>.
- Liao D, Zhou F, Luo L, et al. Haematological characteristics and risk factors in the classification and prognosis evaluation of COVID-19: a retrospective cohort study. *Lancet Haematol.* 2020;7(9):e671–8. [https://doi.org/10.1016/S2352-3026\(20\)30217-9](https://doi.org/10.1016/S2352-3026(20)30217-9).
- Garcia-Beltran WF, Lam EC, Astudillo MG, et al. COVID-19-neutralizing antibodies predict disease severity and survival. *Cell.* 2020. <https://doi.org/10.1016/j.cell.2020.12.015>.
- Zohar T, Loos C, Fischinger S, et al. Compromised humoral functional evolution tracks with SARS-CoV-2 mortality. *Cell.* 2020;183(6):1508–1519.e12. <https://doi.org/10.1016/j.cell.2020.10.052>.

23. Coppo A, Bellani G, Winterton D, et al. Feasibility and physiological effects of prone positioning in non-intubated patients with acute respiratory failure due to COVID-19 (PRON-COVID): a prospective cohort study. *Lancet Respir Med*. 2020;8(8):765–74. [https://doi.org/10.1016/S2213-2600\(20\)30268-X](https://doi.org/10.1016/S2213-2600(20)30268-X).
24. Koeckerling D, Barker J, Mudalige NL, et al. Awake prone positioning in COVID-19. *Thorax*. 2020;75(10):833–4. <https://doi.org/10.1136/thoraxjnl-2020-215133>.
25. Thompson AE, Ranard BL, Wei Y, Jelic S. Prone positioning in awake, nonintubated patients with COVID-19 hypoxemic respiratory failure. *JAMA Intern Med*. 2020;180(11):1537–9. <https://doi.org/10.1001/jamainternmed.2020.3030>.
26. Dennis JM, McGovern AP, Vollmer SJ, Mateen BA (2020) Improving survival of critical care patients with coronavirus disease in England: a national cohort study, March to June 2020. *Crit Care Med*. <https://doi.org/10.1097/CCM.0000000000004747>
27. Horwitz LI, Jones SA, Cerfolio RJ, et al. Trends in COVID-19 risk-adjusted mortality rates. *J Hosp Med*. Published online October 21, 2020. <https://doi.org/10.12788/jhm.3552>
28. Wynants L, Van Calster B, Collins GS, et al. Prediction models for diagnosis and prognosis of covid-19 infection: systematic review and critical appraisal. *BMJ*. 2020;369:m1328. <https://doi.org/10.1136/bmj.m1328>.
29. Friedman JH. Stochastic gradient boosting. *Comput Stat Data Anal*. 2002;38(4):367–78. [https://doi.org/10.1016/S0167-9473\(01\)00065-2](https://doi.org/10.1016/S0167-9473(01)00065-2).
30. Schapire RE. The boosting approach to machine learning: an overview. In: Denison DD, Hansen MH, Holmes CC, Mallick B, Yu B, eds. *Nonlinear estimation and classification*. Lecture notes in statistics. Springer; 2003:149–171. [https://doi.org/10.1007/978-0-387-21579-2\\_9](https://doi.org/10.1007/978-0-387-21579-2_9)
31. CDC COVID-19 Response Team. Severe outcomes among patients with coronavirus disease 2019 (COVID-19)—United States, February 12–March 16, 2020. *MMWR Morb Mortal Wkly Rep*. 2020;69(12):343–346. <https://doi.org/10.15585/mmwr.mm6912e2>
32. Richardson S, Hirsch JS, Narasimhan M, et al. Presenting characteristics, comorbidities, and outcomes among 5700 patients hospitalized with COVID-19 in the New York City area. *JAMA*. 2020;323(20):2052–9. <https://doi.org/10.1001/jama.2020.6775>.
33. Onder G, Rezza G, Brusaferro S. Case-fatality rate and characteristics of patients dying in relation to COVID-19 in Italy. *JAMA*. 2020;323(18):1775–6. <https://doi.org/10.1001/jama.2020.4683>.
34. Shi S, Qin M, Shen B, et al. Association of cardiac injury with mortality in hospitalized patients with COVID-19 in Wuhan, China *JAMA Cardiol*. 2020;5(7):802–10. <https://doi.org/10.1001/jamacardio.2020.0950>.
35. Rosenthal N, Cao Z, Gundrum J, Sianis J, Safo S. Risk factors associated with in-hospital mortality in a US national sample of patients with COVID-19. *JAMA Netw Open*. 2020;3(12):e2029058. <https://doi.org/10.1001/jamanetworkopen.2020.29058>.
36. Mesas AE, Caverro-Redondo I, Álvarez-Bueno C, et al. Predictors of in-hospital COVID-19 mortality: a comprehensive systematic review and meta-analysis exploring differences by age, sex and health conditions. *PLoS ONE*. 2020;15(11):e0241742. <https://doi.org/10.1371/journal.pone.0241742>.
37. Wang C, Kang K, Gao Y, et al. Cytokine levels in the body fluids of a patient with COVID-19 and acute respiratory distress syndrome: a case report. *Ann Intern Med*. 2020;173(6):499–501. <https://doi.org/10.7326/L20-0354>.
38. Mehta P, McAuley DF, Brown M, et al. COVID-19: consider cytokine storm syndromes and immunosuppression. *Lancet Lond Engl*. 2020;395(10229):1033–4. [https://doi.org/10.1016/S0140-6736\(20\)30628-0](https://doi.org/10.1016/S0140-6736(20)30628-0).
39. Weinberger DM, Chen J, Cohen T, et al. Estimation of excess deaths associated with the COVID-19 pandemic in the United States, March to May 2020. *JAMA Intern Med*. 2020;180(10):1336–44. <https://doi.org/10.1001/jamainternmed.2020.3391>.
40. Verity R, Okell LC, Dorigatti I, et al. Estimates of the severity of coronavirus disease 2019: a model-based analysis. *Lancet Infect Dis*. 2020;20(6):669–77. [https://doi.org/10.1016/S1473-3099\(20\)30243-7](https://doi.org/10.1016/S1473-3099(20)30243-7).
41. Nguyen NT, Chinn J, Nahmias J, et al. Outcomes and mortality among adults hospitalized with COVID-19 at US medical centers. *JAMA Netw Open*. 2021;4(3):e210417. <https://doi.org/10.1001/jamanetworkopen.2021.0417>.
42. Pollock AM, Lancaster J. Asymptomatic transmission of covid-19. *BMJ*. 2020;371:m4851. <https://doi.org/10.1136/bmj.m4851>.
43. Woolf SH, Chapman DA, Sabo RT, Weinberger DM, Hill L, Taylor DDH. Excess deaths from COVID-19 and other causes, March–July 2020. *JAMA*. 2020;324(15):1562–4. <https://doi.org/10.1001/jama.2020.19545>.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.