

Early risk assessment for COVID-19 patients from emergency department data using machine learning

Authors

Frank S. Heldt¹, Marcela P. Vizcaychipi^{2,3}, Sophie Peacock¹, Mattia Cinelli¹, Lachlan McLachlan¹, Fernando Andreotti¹, Stojan Jovanović¹, Robert Dürichen¹, Nadezda Lipunova¹, Robert A. Fletcher¹, Anne Hancock¹, Alex McCarthy², Richard A. Pointon², Alexander Brown², James Eaton², Roberto Liddi¹, Lucy Mackillop^{1,4,5}, Lionel Tarassenko^{1,6}, Rabia T. Khan^{1,*}

Affiliations

¹Sensyne Health plc, Schrodinger Building, Heatley Road, Oxford Science Park, Oxford, OX4 4GE.

²Chelsea and Westminster Hospital NHS Foundation Trust, 369 Fulham Road, London, SW10 9NH, UK.

³Academic Department of Anaesthesia & Intensive Care Medicine, Imperial College London, Chelsea & Westminster Campus, 369 Fulham Road, London, SW10 9NH, UK.

⁴Oxford University Hospitals NHS Foundation Trust, Women's Centre, John Radcliffe Hospital, Headley Way, Headington, Oxford, OX3 9DU, UK.

⁵Nuffield Department of Women's and Reproductive Health, University of Oxford, Women's Centre, John Radcliffe Hospital, Headley Way, Headington, Oxford, OX3 9DU, UK.

⁶Institute of Biomedical Engineering, Department of Engineering Science, University of Oxford, OX3 7DQ.

*Corresponding author

Correspondence address:

Rabia Tahir Khan

Sensyne Health plc, Schrodinger Building, Heatley Road, Oxford Science Park, Oxford, OX4 4GE

Email: rabia.khan@sensynehealth.com

Keywords

SARS-CoV-2, COVID-19, machine learning, electronic healthcare records, risk factors, critical care, mechanical ventilation, mortality

Running title

COVID-19 patient risk assessment using machine learning

33 Abstract

34 **Background** Since its emergence in late 2019, the severe acute respiratory syndrome
35 coronavirus 2 (SARS-CoV-2) has caused a pandemic, with more than 4.8 million reported
36 cases and 310 000 deaths worldwide. While epidemiological and clinical characteristics of
37 COVID-19 have been reported, risk factors underlying the transition from mild to severe
38 disease among patients remain poorly understood.

39
40 **Methods** In this retrospective study, we analysed data of 820 confirmed COVID-19 positive
41 patients admitted to a two-site NHS Trust hospital in London, England, between January 1st
42 and April 23rd, 2020, with a majority of cases occurring in March and April. We extracted
43 anonymised demographic data, physiological clinical variables and laboratory results from
44 electronic healthcare records (EHR) and applied multivariate logistic regression, random
45 forest and extreme gradient boosted trees. To evaluate the potential for early risk
46 assessment, we used data available during patients' initial presentation at the emergency
47 department (ED) to predict deterioration to one of three clinical endpoints in the remainder of
48 the hospital stay: A) admission to intensive care, B) need for mechanical ventilation and C)
49 mortality. Based on the trained models, we extracted the most informative clinical features in
50 determining these patient trajectories.

51
52 **Results** Considering our inclusion criteria, we have identified 126 of 820 (15%) patients that
53 required intensive care, 62 of 808 (8%) patients needing mechanical ventilation, and 170 of
54 630 (27%) cases of in-hospital mortality. Our models learned successfully from early clinical
55 data and predicted clinical endpoints with high accuracy, the best model achieving AUC-
56 ROC scores of 0.75 to 0.83 (F1 scores of 0.41 to 0.56). Younger patient age was associated
57 with an increased risk of receiving intensive care and ventilation, but lower risk of mortality.
58 Clinical indicators of a patient's oxygen supply and selected laboratory results were most
59 predictive of COVID-19 patient trajectories.

60
61 **Conclusion** Among COVID-19 patients machine learning can aid in the early identification of
62 those with a poor prognosis, using EHR data collected during a patient's first presentation at
63 ED. Patient age and measures of oxygenation status during ED stay are primary indicators of
64 poor patient outcomes.

65 Introduction

66 COVID-19, caused by the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), is
67 a novel infectious disease that leads to severe acute respiratory distress in humans. In March
68 2020, the World Health Organisation declared the outbreak a pandemic and, by May 19th, it
69 had caused more than 4 800 000 confirmed cases and 310 000 deaths worldwide [1].
70 Disease severity for COVID-19 appears to vary dramatically between patients, including
71 asymptomatic infection, mild upper respiratory tract illness and severe viral pneumonia with
72 acute respiratory distress, respiratory failure and thromboembolic events that can lead to
73 death [2–4]. Initial reports suggest that 6%-10% of infected patients are likely to become
74 critically ill, most of whom will require mechanical ventilation and intensive care [3,5].

75 Currently, few prognostic markers exist to forecast whether a COVID-19 patient may
76 deteriorate to a critical condition and require intensive care. In general, patients can be
77 grouped into three phenotypes, being at risk of thromboembolic disease, respiratory
78 deterioration and cytokine storm [6]. Early clinical reports find that age, sex and underlying
79 comorbidities, such as hypertension, cardiovascular disease and diabetes, can adversely
80 affect patient outcomes [7,8]. However, few studies have leveraged machine learning to
81 systematically explore risk factors for poor prognosis.

82 Increasingly, hospitals collate large amounts of patient data as electronic healthcare
83 records (EHRs). Combined with state-of-the-art machine learning algorithms, these data can
84 help to predict patient outcomes with greater accuracy than traditional methods [9,10].
85 However, EHR data for COVID-19 remains scarce in the public domain, prompting many
86 authors to focus on statistical analyses instead [11–14]. Where machine learning has been
87 applied to COVID-19, results have been promising, but most studies suffer from a lack of
88 statistical power owing to small sample size [15–18]. Jiang *et al.* applied predictive analytics
89 to data from two hospitals in Wenzhou, China, which included 53 hospitalised COVID-19
90 patients, to predict risk factors for acute respiratory distress syndrome (ARDS) [15]. Exploring
91 the risk factors for in-hospital deaths, Zhou and co-workers used univariate and multivariate
92 logistic regression on data of 191 patients in two hospitals in Wuhan, China [16]. Similarly, Xie
93 *et al.* used logistic regression to predict mortality, training a model on 299 patients and
94 validating it on 145 patients from a different hospital in Wuhan, China [18]. Gong *et al.* used a
95 logistic regression model to identify patients at risk of deterioration to severe COVID-19,
96 applied to the data of 189 patients in Wuhan and Guangdong, China [17].

97 A key factor that determines the success of risk prediction models is the quality and richness
98 of the available data. Studies to date have used a combination of demographics,
99 comorbidities, symptoms, and laboratory tests [15–17,19]. These data typically comprise the
100 patients' entire historical record, as well as observations collected during the current hospital
101 stay [16,18–20]. While the inclusion of a patient's full EHR history improves predictive
102 performance, such approaches may be limited in their clinical applicability to early risk-
103 assessment; at the point of presentation in hospital, the entire EHR of a patient is rarely
104 available.

105 In this work, we retrospectively apply machine learning to data of 820 confirmed COVID-19
106 patients from two tertiary referral urban hospitals in London to predict patients' risk of
107 deterioration to one of three clinical endpoints: A) admission to an adult intensive care unit
108 (AICU), B) need for mechanical ventilation, and C) in-hospital mortality. We restrict our
109 analysis to EHR data available during a patient's first presentation in the emergency
110 department (ED) as this more accurately resembles the hospital reality of early-risk

111 assessment and patient-stratification. Our analysis provides a proof of principle for COVID-19
112 risk assessment, with models achieving a high prediction performance, indicating that patient
113 age, oxygenation status and selected laboratory tests are prime indicators of patient
114 outcome.

115

116 Methods

117 Data collection and study design

118 Anonymised EHR data of patients admitted to two hospitals in London, England, between
119 January 1st, 2020 and April 23rd, 2020, were gathered by Chelsea & Westminster NHS
120 Foundation Trust (NHS Trust, hereafter). The data was supplied in accordance with internal
121 information governance review, NHS Trust information governance approval, and General
122 Data Protection Regulation (GDPR) procedures outlined under the Strategic Research
123 Agreement (SRA) and relative Data Sharing Agreements (DSAs) signed by the NHS Trust and
124 Sensyne Health plc on 25th July 2018.

125 Data encompasses clinical observations collated from inpatient encounters. The analysis was
126 restricted to adult patients aged between 18 and 100 years at the time of their most recent
127 hospital admission (assumed to be the COVID-19-related admission). Only confirmed SARS-
128 CoV-2 positive patients, as determined by quantitative reverse-transcription PCR (qRT-PCR),
129 were included. 65% of patients were male and 35% female (Table 1). The majority were white
130 British (28%) or did not state their ethnicity (24%) (see also Fig. S1). All clinical features and
131 their coverage in the data set are listed in Table S1. Features include patient demographics (3
132 in total), vital signs (4 in total), laboratory measurements and clinical observations (60 in total).
133 For vital signs and laboratory measurements, patients may have received multiple test results
134 during their stay. These values were aggregated for each feature to only retain the respective
135 minimum, maximum, mean and last observation value. Only clinical features with at least 5%
136 coverage in the patient population were considered. The data set covered the patient's entire
137 encounter history from their admission to the hospital's ED, with a median length of stay in
138 that department of 5 hours, to their discharge. The median length of in-hospital stay was 7.2
139 days.

140

141 Cohort definition

142 A total of 3229 patients fell within the observation time and study parameters. From these
143 patients, three cohorts were derived, one for each clinical endpoint, as follows (see Fig. S2
144 for flow diagram and patient numbers). Only confirmed COVID-19 positive patients were
145 considered. Patients who did not have information relating to an admission to any hospital
146 department in 2020 were excluded. Furthermore, the following exclusion criteria were applied
147 to each of the considered endpoints: for cohort A) patients without a documented ward
148 location were excluded; for cohort B) patients without information on oxygen supply were
149 excluded; for cohort C) patients without hospital discharge information were excluded.
150 Finally, since our models were trained on data available during a patient's stay in the ED, we
151 removed patients who did not have a documented ED visit.

152

153 Each cohort was divided into target and control groups (see Table 2). For AICU admission,
154 target patients comprise those that were admitted to an AICU at any time during their
155 hospital stay, while control patients are those that remained in any other ward for their entire
156 admission. Target patients in the ventilation cohort were defined as requiring invasive

157 mechanical ventilation, whereas control patient required no or only minimal breathing
 158 assistance. Both categories are based on clinical records of oxygen supply according to
 159 Table 3. Note that from clinical data the total number of mechanically ventilated patients was
 160 135, however only 62 were visible in our data. This results from staggered deployment of
 161 EHR data in the two hospitals such that one site is understood to lack certain data related to
 162 mechanical ventilation. Mortality data was based on the discharge destination (mortuary) in
 163 clinical records. All regularly discharged patients or patients remaining in hospital were
 164 considered alive.

165
 166 *Table 1. Composition of patient population.*

| Demographics | |
|-----------------------------------|-----------------|
| Patient age (years) | |
| Range | 18-100 169 |
| Overall mean (standard deviation) | 67.3 (16.8) 170 |
| Female mean (standard deviation) | 70.3 (17.2) 171 |
| Male mean (standard deviation) | 65.8 (16.4) 172 |
| Sex (number of patients) | |
| Female | 286 (34.9%) 174 |
| Male | 533 (65.0%) 175 |
| unknown | 1 (0.1%) 176 |
| Ethnicity (number of patients) | |
| White British | 230 (28%) 177 |
| Not Stated | 196 (23.9%) 179 |
| Ethnic Other | 97 (11.8%) 180 |
| White Other | 76 (9.3%) 181 |
| Asian Indian | 63 (7.7%) 182 |
| Asian Other | 39 (4.8%) 183 |
| Unknown | 29 (3.5%) 184 |
| Black African | 24 (2.9%) 185 |
| Black Caribbean | 23 (2.8%) 186 |
| Asian Pakistani | 11 (1.3%) 187 |
| Black Other | 10 (1.2%) 188 |
| Others | 22 (2.7%) 189 |

190
 191 *Table 2. Clinical endpoint cohorts.*

| | Cohort A (AICU admission) | Cohort B (ventilation) | Cohort C (mortality) |
|--------------------|------------------------------|---------------------------|-------------------------|
| Number of patients | 820 | 808 | 630 |
| Target patients | 126 (15%) | 62 (8%) | 170 (27%) |
| Control patients | 694 (85%) | 742 (92%) | 460 (73%) |

192
 193 *Table 3. Target and control definition for ventilation cohort.*

| Category | Clinical observation value |
|----------|--|
| Control | room air, air/none, nasal cannulae, high flow nasal cannulae, venturi mask, face mask, non-rebreather mask, simple face mask, swedish nose with, oxygen, mask, HFOV, face/tracheostomy mask, CPAP, BiPAP |
| Target | ventilator, tracheostomy, CMV, VC-CMV, t-piece, HELIOX, IPPV, SIMV, PC-BIPAP, APRV, CPAP / ASB_SPN / CPAP/PS |

194 Patient outcome prediction

195 Three machine-learning algorithms were benchmarked to predict patient outcomes from EHR
196 data: logistic regression, random forest and Extreme Gradient Boosted Trees (XGBoost).
197 Logistic regression, which predicts the probability of a clinical endpoint as a linear function of
198 the feature space, was used as a baseline algorithm. The model was regularised with elastic
199 net using equal weighting given to L_1 and L_2 penalties in order to account for the high
200 dimensionality of the data set relative to the number of observations. A random forest [21],
201 i.e., an ensemble of decision trees where each tree is trained on a slightly different subset of
202 data, was trained using 100 trees and splits were evaluated using Gini impurity. Classes were
203 inversely weighted to account for the class imbalance present in the data set. An XGBoost
204 algorithm [22] was trained with its hyperparameters set to 100 trees, max tree-depth of 6,
205 step-shrinkage of 0.3, no subsampling and L_2 regularisation, to minimize log-loss. This tree-
206 based algorithm trains decision trees sequentially, with each new tree being trained on the
207 residuals of previous trees.

209 Performance evaluation

210 All models were evaluated using a stratified 3-fold cross-validation strategy. Results are
211 reported as mean and standard deviation across these folds. Predictive performance was
212 measured in terms of area under curve (AUC) of the receiver operating characteristic (ROC)
213 as well as F1 score at each model's ideal classification threshold as derived from the ROC
214 curve. Given the presence of class-imbalance, precision-recall curves were also computed to
215 assess expected real-world performance relative to random classifiers.

216
217 In order to extract the clinical features most relevant to predictions, permutation feature
218 importance (PFI) was calculated for each model post-hoc [21,23]. Each feature was
219 individually randomised. The model's AUC-ROC on the validation sets was then compared to
220 the AUC-ROC before the feature had been randomised. PFI provides an estimate of the
221 extent to which a model relies on a feature for its predictive performance and generalisability.
222 The changes in performance were normalised by the sum of absolute changes over all
223 features. Averages and standard deviations over the validation sets have been reported.

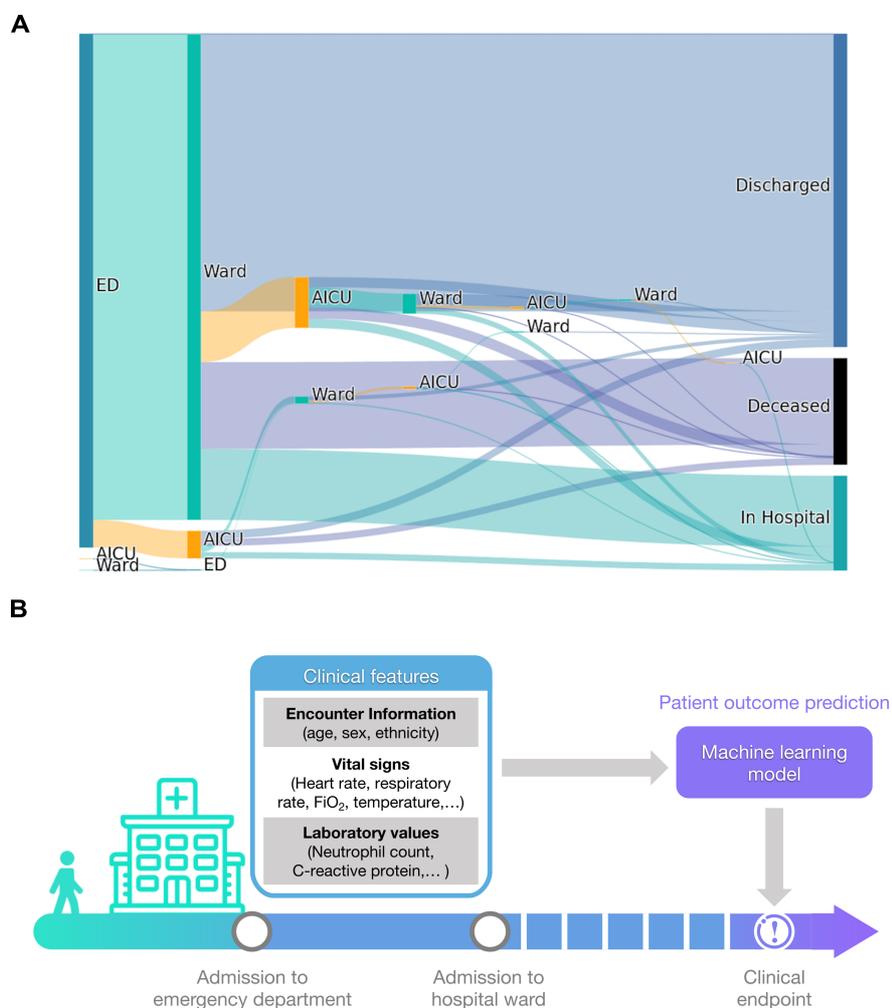
224
225 Accumulated local effects (ALE) were computed to determine the directionality of a feature's
226 effect on model predictions [24]. Specifically, the feature space was divided into ten
227 percentile bins and each feature's effect was calculated as the difference in predictions
228 between the upper and lower bounds of each bin, leaving all other features unchanged.
229 Binning features in this way can reduce the influence of correlated features often encountered
230 when trying to isolate the effect of a single feature.

231 Results

232 Patient pathways

233 A summary of observed patient in-hospital pathways is shown in Figure 1A. Of the 820
 234 patients in cohort A, which we present as an example, 818 (99.8%) entered the hospital via
 235 the ED, while 1 (0.1%) and 1 (0.1%) patients were admitted directly to a ward and the AICU,
 236 respectively. Upon leaving the ED, 775 (94.5%) patients transitioned to regular wards and 44
 237 (5.4%) to an AICU. Of the 775 patients in regular wards, 81 (10.5%) patients required
 238 subsequent admission to an AICU, 441 (57%) were discharged, 113 (14.5%) remained in
 239 hospital and 138 (18%) succumbed to the infection. From the 126 patients that have been
 240 admitted to an AICU, 57 (37%) were ultimately discharged, 32 (35%) did not survive and 37
 241 (29%) are still in hospital. Patients' median length of stay in ED was 5 hours (IQR 3.45 hours).
 242 During this time, demographic information, vitals and laboratory values were collected (Fig.
 243 1B). To aid an early patient stratification, our models use data collected during the ED stay
 244 only to predict whether a patient reached any of three clinical endpoints during their
 245 subsequent admission.

246



247

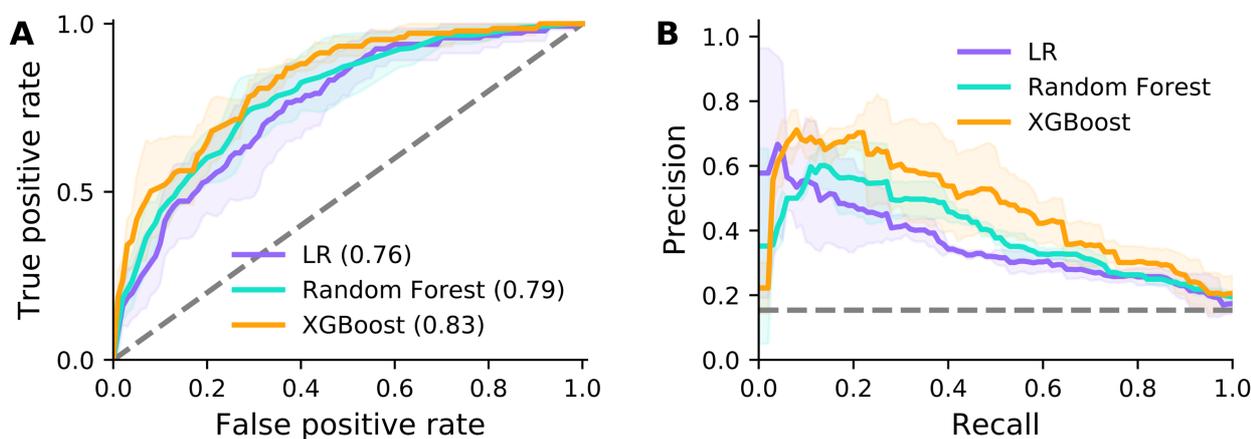
248 **Figure 1. Patient pathways and outcome prediction.** (A) Patient transitions between hospital departments are
 249 shown as bands proportional in size to patient numbers. Different departments are indicated by rectangles (ED,
 250 emergency department; Ward, regular hospital ward; AICU, adult intensive care unit). Patients who remain in
 251 hospital, are being discharged or die in hospital are indicated on the right. (B) Patient outcome prediction
 252 models use clinical data recorded within the ED stay of a patient to predict clinical endpoints during the
 253 remainder of the in-hospital stay.

254

255 **AICU admission**

256 First, we studied patients transitioning to critical care and requiring admission to an AICU. All
 257 three models reach good prediction performance on this endpoint, as measured by area
 258 under the curve (AUC) of the receiver operating characteristic (ROC) and precision-recall
 259 curves, significantly outperforming random classifiers (Fig. 2). The best performing model,
 260 XGBoost, reaches an AUC-ROC of 0.83 and an F1 score of 0.51. Both tree-based methods
 261 perform better than logistic regression (Table 4). This is to be expected since logistic
 262 regression cannot model interactions between features unless such interactions are explicitly
 263 encoded into the training data set through feature engineering. All models show a moderate
 264 amount of variability across cross-validation folds (notice standard deviations in Fig. 2 and
 265 Table 4), which can compromise subsequent analyses. This instability originates from the
 266 limited number of patients and high class imbalance between target and control patients (see
 267 Table 2). Specifically, in each of the three cross-validation folds the models are only trained
 268 and validated on two thirds and one third of the data set, respectively, leaving few target
 269 patients for these tasks.

270



271

272 *Figure 2. Prediction performance for AICU admission. Model performance for the logistic regression (LR),*
 273 *random forest and XGBoost models are shown as ROC (A) and precision-recall curves (B). AUC under ROC is*
 274 *provided in brackets. Solid lines and shaded areas indicate the mean and standard deviation across three*
 275 *cross-validation folds, respectively. Dashed lines indicate random classifiers.*

276

277 *Table 4. Model performance on clinical endpoint prediction (standard deviation shown in brackets).*

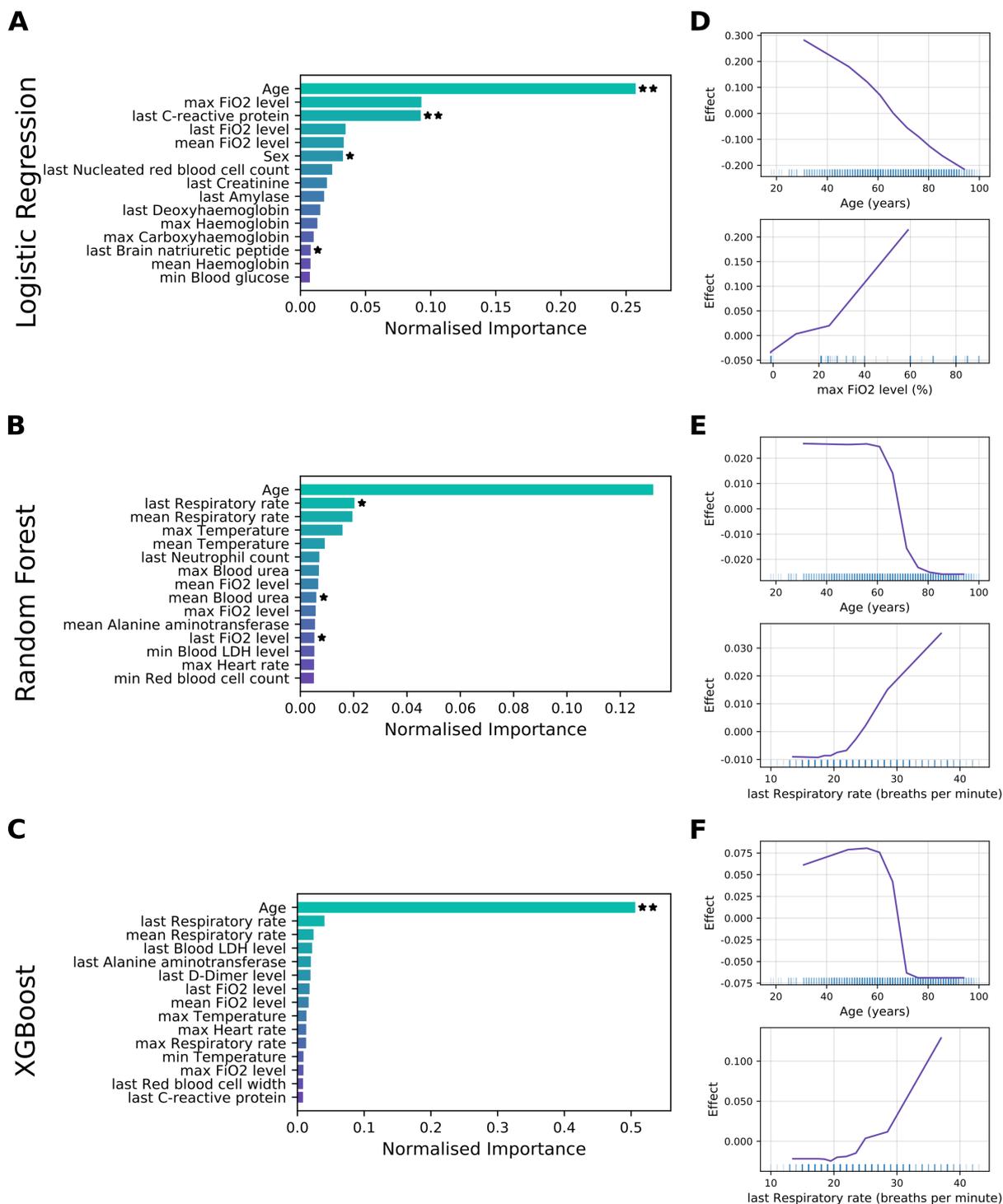
| Model | Endpoint A (AICU admission) | | Endpoint B (ventilation) | | Endpoint C (mortality) | |
|---------------------|--------------------------------|------------------------|-----------------------------|------------------------|---------------------------|------------------------|
| | AUC | F1 | AUC | F1 | AUC | F1 |
| Logistic regression | 0.76 (0.067) | 0.40 (0.029) | 0.79 (0.097) | 0.41 (0.083) | 0.66 (0.030) | 0.50 (0.035) |
| Random forest | 0.79 (0.058) | 0.41 (0.031) | 0.81 (0.045) | 0.37 (0.081) | 0.75 (0.016) | 0.55 (0.039) |
| XGBoost | 0.83 (0.045) | 0.51 (0.037) | 0.83 (0.083) | 0.41 (0.052) | 0.74 (0.011) | 0.56 (0.035) |

278

279 Next, we assessed which clinical variables contribute the most to model predictions by
 280 applying PFI. Figure 3A presents the 15 most important features for the logistic regression
 281 with elastic net regularisation. Note that clinical variables that can be recorded multiple times

282 during a patient's ED visit were aggregated to retain only the minimum, maximum, mean and
283 last observation value during the ED stay. Patient age, C-reactive protein and sex reached
284 high importance and significance over cross-validation folds for the logistic regression.
285 Moreover, the fraction of inspired oxygen (FiO₂) contributes to predictions, albeit without
286 being significant. The random forest (Fig. 3B) and XGBoost (Fig. 3C) models assign a higher
287 importance to patient age, with respiratory rate following thereafter. Intriguingly, ALE analyses
288 reveal that lower patient age increases the likelihood of AICU admission in all three
289 models (Figs. 3D-F). This agrees well with a bias towards younger patients when comparing
290 AICU-admitted patients with control patients (Fig. S3A). However, clinical indicators of
291 disease severity, such as C-reactive protein and ferritin levels, show no clear trend across
292 age groups (Fig. S4). We also find that the fraction of inspired oxygen (Fig. 3D) and
293 respiratory rate (Figs. 3E and F) exhibit a positive effect on AICU admission probability.

294 In summary, machine learning algorithms can predict those patients most likely to require
295 AICU admission in COVID-19 patients from EHR data available during the initial ED stay with
296 high precision. Patient age and indicators of oxygenation status are strong indicators of
297 patient outcome, with advanced age decreasing the probability of AICU admission.



298
299

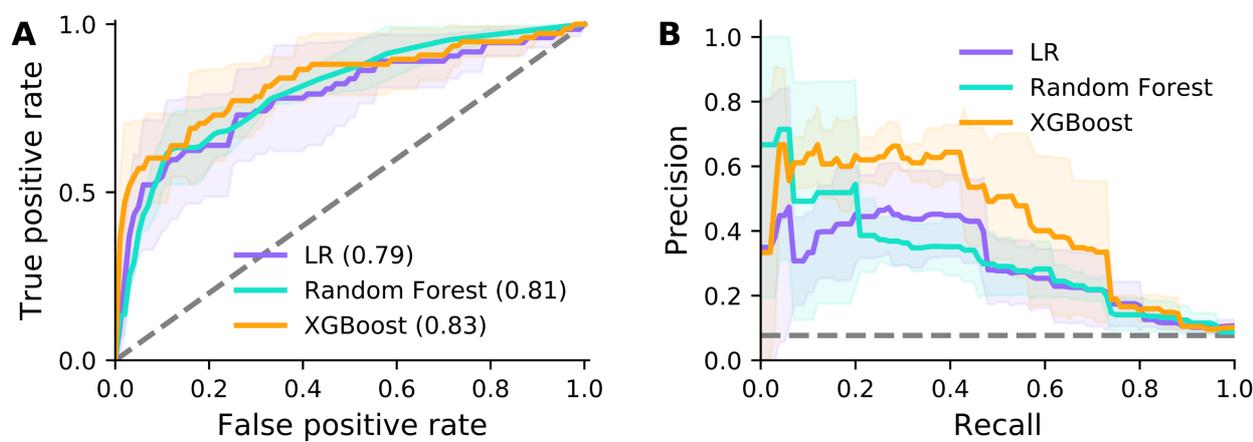
300 *Figure 3. Feature importance for AICU admission. (A-C) Permutation feature importance for the logistic*
 301 *regression (A), random forest (B) and XGBoost (C) models. Only the top 15 features are shown. Asterisks mark*
 302 *features with importance scores significantly different from zero across three cross-validation folds with t-test p-*
 303 *value thresholds of 5% (*) and 1% (**). (D-F) Accumulated local effects plots for the logistic regression (D),*
 304 *random forest (E) and XGBoost models (F). The top two features according to permutation feature importance*
 305 *are shown for each model. Vertical bars at the bottom indicate feature values observed in the data set.*

306
307

308 **Mechanical ventilation**

309 For mechanical ventilation prediction, we categorised patients into those that needed a
310 ventilator (e.g., patients receiving SIMV, BIPAP or APRV ventilation) and control patients that
311 either were able to breathe normally or required minimal assistance (e.g., those patients
312 receiving oxygen via nasal cannulae or face masks). Prediction performance on this endpoint
313 is comparable to prediction of AICU admission (Fig. 4). Specifically, XGBoost performs best,
314 reaching an AUC of 0.83, while logistic regression and random forest reach 0.79 and 0.81,
315 respectively (Table 4). This result is expected since most patients receive mechanical
316 ventilation in AICU, meaning the ventilation cohort is a subset of the critical care cohort (56 of
317 62 target patients in Cohort B are target patients in Cohort A). Notably, all models show a
318 decrease in stability in predicting this clinical endpoint. This is most likely due to a higher
319 class-imbalance and lower number of patients receiving ventilation.

320



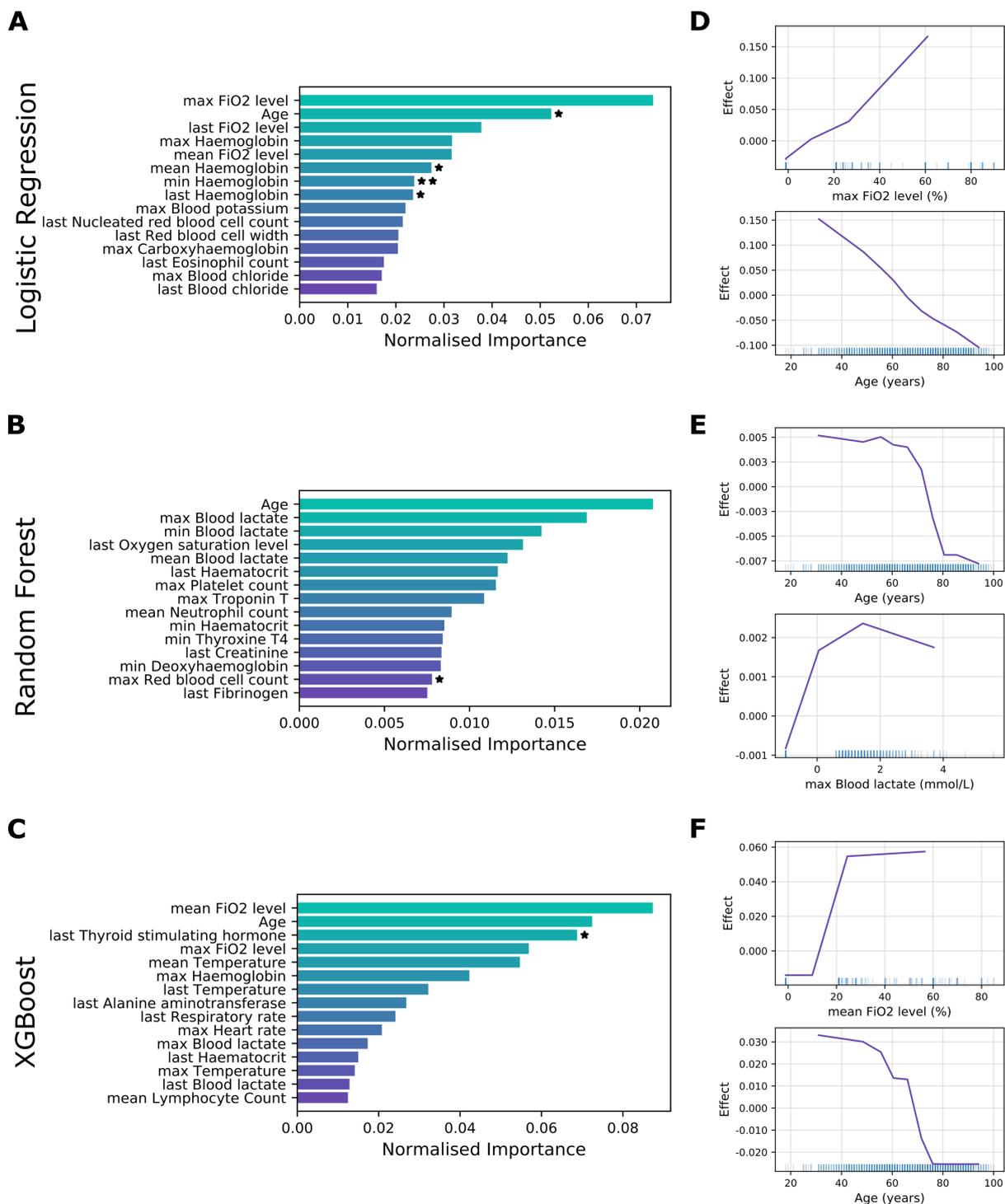
321

322 **Figure 4. Prediction performance for mechanical ventilation.** Model performance for the logistic regression (LR),
323 random forest and XGBoost models are shown as ROC (A) and precision-recall curves (B). AUC under ROC is
324 provided in brackets. Solid lines and shaded areas indicate the mean and standard deviation across three
325 cross-validation folds, respectively. Dashed lines indicate random classifiers.

326

327 Feature importance analysis for the logistic regression shows a large effect of the fraction of
328 inspired oxygen and patient age (Fig. 5A). This mirrors the results for AICU admission. We
329 also observe a significant influence of haemoglobin levels on model predictions. Both tree-
330 based methods rank age highly (Figs. 5B and C). In addition, blood lactate levels and oxygen
331 saturation are used by the random forest (Fig. 5B), while XGBoost relies on the fraction of
332 inspired oxygen and levels of thyroid stimulating hormone (Fig. 5C), although few values are
333 significant. In general, all models rely on a broader set of features for the ventilation endpoint.
334 ALE analysis shows younger patients had an increased probability of receiving
335 ventilation (Fig. 5D-F), which agrees with an inherent bias towards younger age when
336 comparing ventilated with non-ventilated patients (Fig.S4B). By contrast, a higher fraction of
337 inspired oxygen and higher blood lactate level were associated with a poor prognosis.

338 Taken together, models show good performance when predicting ventilation, albeit with a
339 decreased model stability (higher standard deviation). Patient age and oxygenation status are
340 most predictive of poor outcome, with additional contributions from blood test values, such
341 as lactate and haemoglobin levels.



342

343

344

345

346

347

348

349

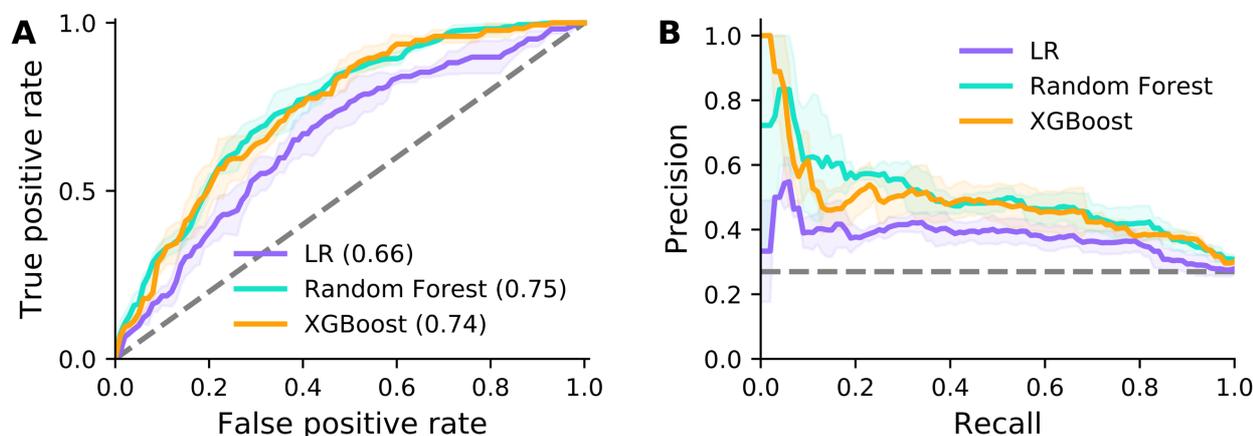
350

Figure 5. Feature importance for mechanical ventilation. Permutation feature importance for the random forest (A), logistic regression (B) and XGBoost (C) models. Only the top 15 features are shown. Asterisks mark features with importance scores significantly different from zero across three cross-validation folds with *t*-test *p*-value thresholds of 5% (*) and 1% (**). (D-F) Accumulated local effects plots for the logistic regression (D), random forest (E) and XGBoost models (F). The top two features according to permutation feature importance are shown for each model. Vertical bars at the bottom indicate feature values observed in the data set.

351 **Mortality**

352 The performance of all three models shows a marked decrease when predicting mortality
353 (Fig. 6). The logistic regression and XGBoost reach AUCs of 0.66 and 0.74, respectively, only
354 outperformed by random forest reaching an AUC of 0.75. However, model stability is
355 improved with standard deviations across cross-validation folds reaching their lowest levels
356 over all three clinical endpoints (Table 4).

357



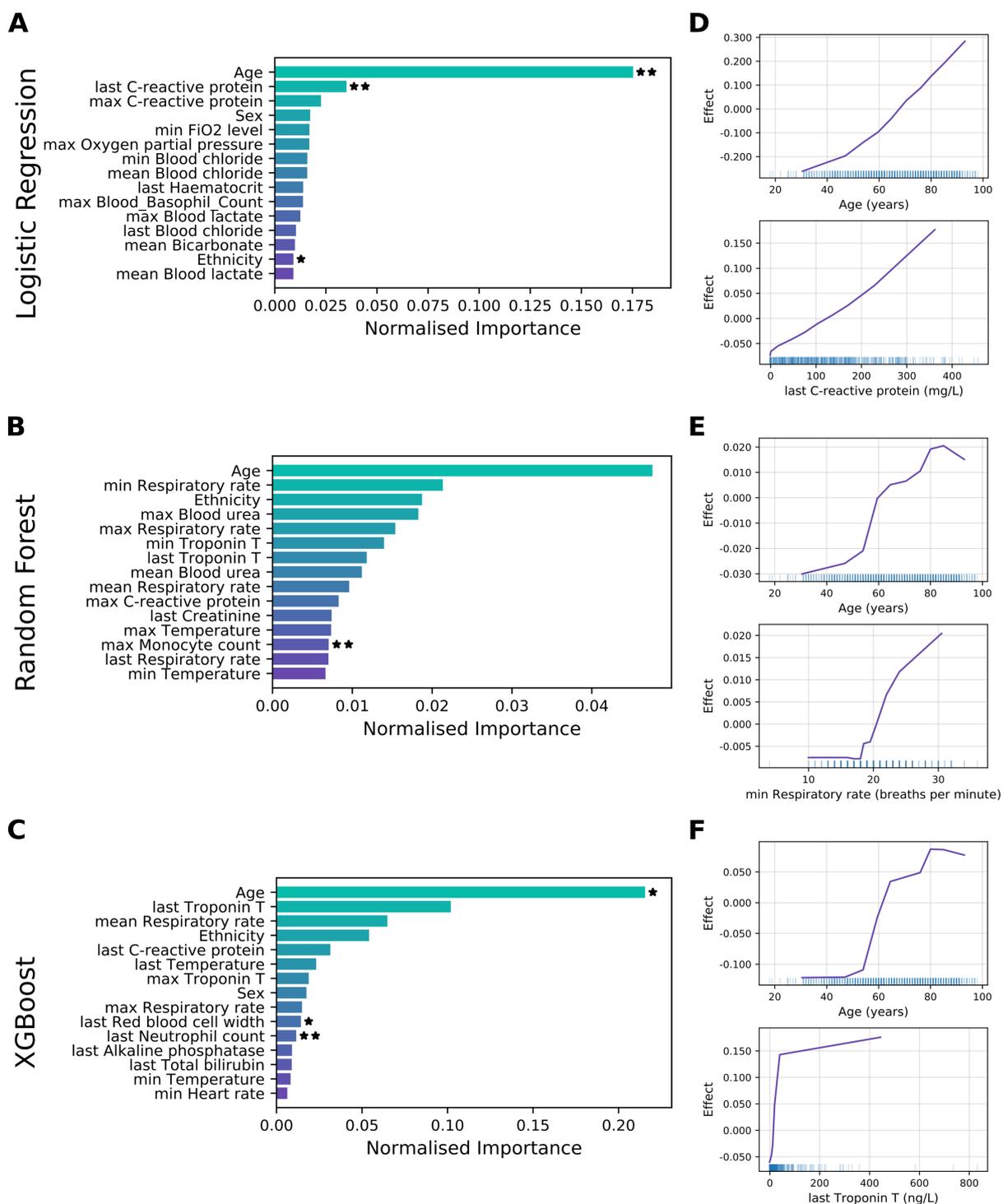
358

359 **Figure 6. Prediction performance for mortality.** Model performance for the logistic regression (LR), random forest
360 and XGBoost models are shown as ROC (A) and precision-recall curves (B). AUC under ROC is provided in
361 brackets. Solid lines and shaded areas indicate the mean and standard deviation across three cross-validation
362 folds, respectively. Dashed lines indicate random classifiers.

363

364 Predictions from the logistic regression model are dominated by patient age, with C-reactive
365 protein levels adding a small but significant contribution (Fig. 7A). Similarly, tree-based
366 methods rely heavily on age for their predictions, with smaller contributions of respiratory rate
367 and Troponin T levels (Figs. 7B and C). More generally, prediction of mortality relies more
368 strongly on blood tests as opposed to indicators of oxygen supply observed in other cohorts.
369 ALE analysis shows that advanced age is predictive of higher mortality (Fig. 7D-F). This
370 agrees with a bias towards older age in patients that die in hospital (Fig. S4C). Higher C-
371 reactive protein, respiratory rate and Troponin T levels increase the risk of mortality in our
372 models (Figs. 7D-F).

373 In summary, models show an increased stability but lower overall performance when
374 predicting mortality. Feature importance scores reveal a high and significant contribution of
375 patient age with advanced age contributing to poor patient outcomes.



376
377

378 *Figure 7. Feature importance for mortality. (A-C) Permutation feature importance for the logistic regression (A),*
 379 *random forest (B) and XGBoost (C) models. Only the top 15 features are shown. Asterisks mark features with*
 380 *importance scores significantly different from zero across three cross-validation folds with t-test p-value*
 381 *thresholds of 5% (*) and 1% (**). (D-F) Accumulated local effects plots for the logistic regression (D),*
 382 *random forest (E) and XGBoost models (F). The top two features according to permutation feature importance*
 383 *are shown for each model. Vertical bars at the bottom indicate feature values observed in the data set.*

384 Discussion

385 Disease severity can vary dramatically between COVID-19 patients, ranging from
386 asymptomatic infection to severe respiratory distress and failure. To evaluate the potential of
387 an early stratification of hospitalised patients into risk groups, we built machine learning
388 models from EHR care data of confirmed Covid-19 positive patients, aimed at predicting one
389 of three clinical endpoints: admission to AICU, the need for mechanical ventilation and
390 mortality. On all three cohorts, our models reach good performance with the best model
391 showing AUC-ROC between 0.75 and 0.83. Overall, mortality proved to be the most difficult
392 prediction task, presumably reflecting the complex interactions underlying in-hospital death.

393 The most predictive feature for all three endpoints was patient age, followed by indicators of
394 patients' oxygenation status, including fraction of inspired oxygen and respiratory rate. Given
395 that SARS-CoV-2 causes an infection of the respiratory tract, which can lead to severe
396 respiratory distress, these results were to be expected. Our findings are supported by similar
397 works, in which age is consistently found to be the most important feature [16–18]. However,
398 we note that other potential indicators for severe viral infection, like increased temperature
399 and markers of immune system activation, e.g. C-reactive protein, are less prominent in our
400 feature importance scores. Overall, prediction of mortality relies more strongly on blood tests
401 as opposed to indicators of oxygen supply observed in other cohorts. The reason for this
402 observation and its clinical significance is, as of yet, unclear. Our ALE analysis reveals that
403 lower patient age contributes to an increased probability of receiving mechanical ventilation
404 and critical care in AICU, while coinciding with lower mortality. We also note that Docherty *et al.*
405 find that 17% of COVID-19 patients require admission to a High Dependency or Intensive
406 Care Unit [25], which is similar to 15% of patients in our data.

407 Conversely, our findings concerning the importance of features relating to patients'
408 oxygenation status are not corroborated by other works. Specifically, other studies find that
409 one important predictor of patient outcome is the level of lactate dehydrogenase [17,18],
410 which, although present in our data set, does not significantly contribute to predictions.

411 A novel aspect of the present analysis is the use of data limited to a patient's first few hours
412 in ED. While this perhaps more accurately reflects the data available at the time of admission,
413 it may well come at the cost of missing important information, such as medical history or
414 primary care data, for predicting patient outcome. This may explain the comparative difficulty
415 in predicting mortality, since a patient's overall chance of surviving infection may depend
416 heavily on their medical history. Also note that, in our analysis, all patients were considered
417 together for mortality prediction and the cohort was not further split according to
418 confounding factors such as age or sex. In addition, mortality data for recent hospital
419 admissions are by their nature censored, with clinical endpoints for patients who remain in
420 hospital not yet fully known.

421 While we base our study on a comparatively large data set from two hospitals, longitudinal
422 information from additional treatment centres and geographic regions may improve a model's
423 ability to generalise. We note that such data is currently unavailable for COVID-19. However,
424 future studies may benefit from a multicentre approach. As a result of limited data and the
425 imbalanced cohorts, model stability remains a major challenge. While we use inverse class
426 weights and stratified 3-fold cross validation to mitigate this issue, large uncertainties in
427 model results persist, and many predictions do not reach statistical significance. Increased
428 patient numbers, in particular among target patients, may lead to more conclusive results.
429 Once such data is available, more complex models, such as deep neural networks, may

430 achieve higher prediction performance. A key aspect which should be considered in such
431 works is the prediction horizon, which impacts on how useful a model could be.

432 In conclusion, our models represent a first step towards the prediction of COVID-19 patient
433 pathways in hospital at the point of admission in the emergency department. While they
434 succeed in predicting patient outcomes and reveal critical clinical variables that may influence
435 patient trajectories, larger data sets and further analyses are required to draw clinically
436 relevant conclusions.

437

438 Acknowledgments

439 This work uses data provided by patients and collected by the NHS as part of their care and
440 support. We believe using patient data is vital to improve health and care for everyone and
441 would, thus, like to thank all those involved for their contribution.

442 Special thanks are due to the Chelsea and Westminster (CW) COVID-19 AICU Consortium,
443 comprising all critical care personnel who were part of the delivery of care during the COVID-
444 19 pandemic as follows: CW Anaesthetics Consultants, Critical Care Consultants, Trainees &
445 Fellows from ICU, Anaesthesia, and seconded to ICU from other specialities, Surgeons, the
446 supporting Respiratory and ED Physicians, Operating Department Practitioners and CW
447 Critical Care Nurses. This united approach to an unprecedented clinical condition was critical
448 not only to the management of the patients but also to our ability to document and collate the
449 key data in a timely manner to support this analysis.

450 Also a special thank you to Trystan Hawkin, Chris Chaney from CWplus, the Planned Care
451 Clinical Division managers, porters, domestic personnel and the CW local community who
452 without hesitation have supported the National Healthcare System throughout the COVID-19
453 pandemic.

454

455 Ethics statement

456 The data were extracted, anonymised, and supplied by the Trust in accordance with internal
457 information governance review, NHS Trust information governance approval, and General
458 Data Protection Regulation (GDPR) procedures outlined under the Strategic Research
459 Agreement (SRA) and relative Data Sharing Agreements (DSAs) signed by the Trust and
460 Sensyne Health plc on 25th July 2018.

461 References

- 462 1. Situation update worldwide, as of 1 May 2020. In: European Centre for Disease
463 Prevention and Control [Internet]. [cited 1 May 2020]. Available:
464 <https://www.ecdc.europa.eu/en/geographical-distribution-2019-ncov-cases>
- 465 2. Wu Z, McGoogan JM. Characteristics of and Important Lessons From the Coronavirus
466 Disease 2019 (COVID-19) Outbreak in China: Summary of a Report of 72 314 Cases
467 From the Chinese Center for Disease Control and Prevention. *JAMA*. 2020;323: 1239–
468 1242. doi:10.1001/jama.2020.2648
- 469 3. Yang X, Yu Y, Xu J, Shu H, Xia J, Liu H, et al. Clinical course and outcomes of critically ill
470 patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered,
471 retrospective, observational study. *Lancet Respir Med*. 2020 [cited 27 Apr 2020].
472 doi:10.1016/S2213-2600(20)30079-5
- 473 4. Klok F, Kruip M, Van der Meer N, Arbous M, Gommers D, Kant K, et al. Incidence of
474 thrombotic complications in critically ill ICU patients with COVID-19. *Thromb Res*. 2020.
475 doi:10.1016/j.thromres.2020.04.013
- 476 5. Anderson RM, Heesterbeek H, Klinkenberg D, Hollingsworth TD. How will country-based
477 mitigation measures influence the course of the COVID-19 epidemic? *The Lancet*.
478 2020;395: 931–934. doi:10.1016/S0140-6736(20)30567-5
- 479 6. Vizcaychipi MP, Shovlin CL, Hayes M, Singh S, Christie L, Sisson A, et al. Early detection
480 of severe COVID-19 disease patterns define near real-time personalised care,
481 bioseverity in males, and decelerating mortality rates. *medRxiv*. 2020;
482 2020.05.08.20088393. doi:10.1101/2020.05.08.20088393
- 483 7. Novel Coronavirus Pneumonia Emergency Response Epidemiology Team. The
484 epidemiological characteristics of an outbreak of 2019 novel coronavirus diseases
485 (COVID-19) in China. *Chin Cent Dis Control Prev*. 2020;41: 145–151.
486 doi:10.3760/cma.j.issn.0254-6450.2020.02.003
- 487 8. Chen T, Wu D, Chen H, Yan W, Yang D, Chen G, et al. Clinical characteristics of 113
488 deceased patients with coronavirus disease 2019: retrospective study. *BMJ*. 2020;368.
489 doi:10.1136/bmj.m1091
- 490 9. Goldstein BA, Navar AM, Pencina MJ, Ioannidis JPA. Opportunities and challenges in
491 developing risk prediction models with electronic health records data: a systematic
492 review. *J Am Med Inform Assoc*. 2017;24: 198–208. doi:10.1093/jamia/ocw042
- 493 10. Wynants L, Van Calster B, Bonten MM, Collins GS, Debray TP, De Vos M, et al.
494 Prediction models for diagnosis and prognosis of covid-19 infection: systematic review
495 and critical appraisal. *bmj*. 2020;369. doi:10.1136/bmj.m1328
- 496 11. Wang D, Hu B, Hu C, Zhu F, Liu X, Zhang J, et al. Clinical Characteristics of 138
497 Hospitalized Patients With 2019 Novel Coronavirus–Infected Pneumonia in Wuhan,
498 China. *JAMA*. 2020;323: 1061–1069. doi:10.1001/jama.2020.1585
- 499 12. Yang X, Yu Y, Xu J, Shu H, Liu H, Wu Y, et al. Clinical course and outcomes of critically
500 ill patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered,
501 retrospective, observational study. *Lancet Respir Med*. 2020. doi:10.1016/S2213-
502 2600(20)30079-5
- 503 13. Arentz M, Yim E, Klaff L, Lokhandwala S, Riedo FX, Chong M, et al. Characteristics and
504 outcomes of 21 critically ill patients with COVID-19 in Washington State. *Jama*. 2020.
505 doi:10.1001/jama.2020.4326
- 506 14. Hu L, Chen S, Fu Y, Gao Z, Long H, Ren H, et al. Risk Factors Associated with Clinical
507 Outcomes in 323 COVID-19 Patients in Wuhan, China. *medRxiv*. 2020.
508 doi:10.1101/2020.03.25.20037721

- 509 15. Jiang X, Coffee M, Bari A, Wang J, Jiang X, Huang J, et al. Towards an artificial
510 intelligence framework for data-driven prediction of coronavirus clinical severity. *CMC-*
511 *Comput Mater Contin.* 2020;63: 537–51. doi:10.32604/cmc.2020.010691
- 512 16. Zhou F, Yu T, Du R, Fan G, Liu Y, Liu Z, et al. Clinical course and risk factors for
513 mortality of adult inpatients with COVID-19 in Wuhan, China: a retrospective cohort
514 study. *The Lancet.* 2020. doi:10.1016/S0140-6736(20)30566-3
- 515 17. Gong J, Ou J, Qiu X, Jie Y, Chen Y, Yuan L, et al. A Tool to Early Predict Severe 2019-
516 Novel Coronavirus Pneumonia (COVID-19): A Multicenter Study using the Risk
517 Nomogram in Wuhan and Guangdong, China. *medRxiv.* 2020.
518 doi:10.1101/2020.03.17.20037515
- 519 18. Xie J, Hungerford D, Chen H, Abrams ST, Li S, Wang G, et al. Development and
520 external validation of a prognostic multivariable model on admission for hospitalized
521 patients with COVID-19. 2020. doi:10.2139/ssrn.3562456
- 522 19. Pourhomayoun M, Shakibi M. Predicting Mortality Risk in Patients with COVID-19 Using
523 Artificial Intelligence to Help Medical Decision-Making. *medRxiv.* 2020.
524 doi:10.1101/2020.03.30.20047308
- 525 20. Yan L, Zhang H-T, Xiao Y, Wang M, Sun C, Liang J, et al. Prediction of criticality in
526 patients with severe Covid-19 infection using three clinical features: a machine learning-
527 based prognostic model with clinical data in Wuhan. *medRxiv.* 2020.
528 doi:10.1101/2020.02.27.20028027
- 529 21. Breiman L. Random Forests. *Mach Learn.* 2001;45: 5–32.
530 doi:10.1023/A:1010933404324
- 531 22. Chen T, Guestrin C. Xgboost: A scalable tree boosting system. *Proceedings of the 22nd*
532 *acm sigkdd international conference on knowledge discovery and data mining.* 2016.
533 pp. 785–794. doi:10.1145/2939672.2939785
- 534 23. Fisher A, Rudin C, Dominici F. All Models are Wrong, but Many are Useful: Learning a
535 Variable’s Importance by Studying an Entire Class of Prediction Models Simultaneously.
536 *J Mach Learn Res.* 2019;20: 1–81.
- 537 24. Apley DW, Zhu J. Visualizing the Effects of Predictor Variables in Black Box Supervised
538 Learning Models. *ArXiv161208468 Stat.* 2019 [cited 15 Jan 2020]. Available:
539 <http://arxiv.org/abs/1612.08468>
- 540 25. Docherty AB, Harrison EM, Green CA, Hardwick HE, Pius R, Norman L, et al. Features
541 of 16,749 hospitalised UK patients with COVID-19 using the ISARIC WHO Clinical
542 Characterisation Protocol. *medRxiv.* 2020. doi:10.1101/2020.04.23.20076042
543