

# Development, evaluation, and validation of machine learning models for COVID-19 detection based on routine blood tests

Cabitza Federico<sup>1</sup>, Campagner Andrea<sup>2</sup>, Ferrari Davide<sup>3</sup>, Di Resta Chiara<sup>4</sup>, Ceriotti Daniele<sup>5</sup>, Sabetta Eleonora<sup>5</sup>, Colombini Alessandra<sup>2</sup>, De Vecchi Elena<sup>2</sup>, Banfi Giuseppe<sup>2</sup>, Locatelli Massimo<sup>5</sup>, Carobene Anna<sup>5\*</sup>

<sup>1</sup> DISCO, Università degli Studi di Milano-Bicocca, Viale Sarca 336, Milano, 20126, Italy

<sup>2</sup> IRCCS Istituto Ortopedico Galeazzi, Orthopaedic Biotechnology Lab, Via Riccardo Galeazzi, 4, 20161, Milano, Italy

<sup>3</sup> SCVSA Department, University of Parma, Parco Area delle Scienze 11/a, 43124, Parma, Italy.

<sup>4</sup> Vita-Salute San Raffaele University; Unit of Genomics for Human Disease Diagnosis, Division of Genetics and Cell Biology., Via Olgettina 58, 20132, Milan, Italy

<sup>5</sup> Laboratory Medicine, IRCCS San Raffaele Scientific Institute, Via Olgettina 60, 20132, Milan, Italy

**Running title:** Machine learning for COVID-19 detection using blood test

**Keywords:** SARS-CoV-2, COVID-19, machine learning, gradient boosted decision tree, complete blood count, blood laboratory tests

## \*Corresponding author

Anna Carobene, Laboratory Medicine, IRCCS San Raffaele Scientific Institute, Via Olgettina 60, 20132 Milan, Italy

Telephone: +39 02 26432850

E-mail address: [Carobene.anna@hsr.it](mailto:Carobene.anna@hsr.it)

## List of abbreviations

ALP	Alkaline Phosphatase
ALT	Alanine Aminotransferase
ANGPOC	Anion gap
AST	Aspartate Aminotransferase
AUC	Area under the curve
BA	Basophils count (%)
BAT	Basophils count
BEEPOC	Actual Base Excess
BEPOC	Base Excess
BICPOC	Bicarbonates
BILD	Direct Bilirubin
BILIN	Indirect Bilirubin
BILT	Total Bilirubin
BISPOC	Standard Calculated Bicarbonates
BO2POC	Bound O2 Maximum Concentration
CA	Calcium
CAPOC	Ionized Calcium (POC)
CASPOC	Standard Ionized Calcium (POC)
CBC	complete blood count
CK	Creatine kinase
CLPOC	Chloride (POC)
CO2POC	Carbonic Anhydride (pCO2)
CREA	Creatinine
CRP	C-reactive Protein
CT	computed tomography
CTOPOC	Total Oxygen
ED	emergency department
EO	Eosinophils count (%)
EOT	Eosinophils count
FCOPOC	Carboxyhemoglobin
FG	Fibrinogen
FIOPOC	Inspired Oxygen Fraction
FO2POC	Oxyhemoglobin / Total Hemoglobin
GGT	Gamma Glutamyltransferase
GLU	Glucose
GLUEM O	Glucose Blood Gas
HCT	Hematocrit
HCTPOC	Hematocrit (POC)
HGB	Hemoglobin
HHBPOC	Deoxyhemoglobin
IL6	Interleukin 6
IOG	Istituto Ortopedico Galeazzi
K	Potassium

1		
2		
3		
4	KNN	k-nearest neighbors
5	KPOC	Potassium (POC)
6	LATPOC	Lactate (POC)
7	LDH	Lactate Dehydrogenase
8	LR	logistic regression
9	LY	Lymphocytes count (%)
10	LYT	Lymphocytes count
11	MCH	Mean Corpuscular Hemoglobin
12	MCHC	Mean Corpuscular Hemoglobin Concentration
13	MCV	Average Globular Volume
14	METPOC	Methemoglobin
15	ML	machine learning
16	MO	Monocytes count (%)
17	MOT	Monocytes count
18	MPV	Average Platelet Volume
19	NA	Sodium
20	NAPOC	Sodium (POC)
21	NB	Naive Bayes
22	NE	Neutrophils count (%)
23	NET	Neutrophils count
24	NPV	Negative predictive Value
25	OFIPOC	Inspired O <sub>2</sub> / O <sub>2</sub> ratio
26	OSR	Ospedale San Raffaele
27	PCR	Polymerase Chain Reaction
28	PHPOC	pH
29	PLT	Platelets
30	PO2POC	Oxygen (pO <sub>2</sub> )
31	PPTR	Activated partial thromboplastin time ( R )
32	PPV	positive predictive value
33	PROBNP	NT-proB-type Natriuretic Peptide
34	PTINR	Prothrombin Time (INR)
35	rRT-PCR	reverse transcription polymerase chain reaction
36	RBC	Red Blood Cells
37	RDW	Erythrocyte distribution width
38	RF	random forest
39	ROC	receiver operating characteristic
40	RT-PCR	reverse transcriptase-PCR
41	SO2POC	O <sub>2</sub> Saturation
42	SVM	support vector machine
43	THBPOC	Total Oxyhemoglobin
44	TROPOT	Troponin T
45	UREA	Urea
46	WBC	White blood cells
47	XDP	D-Dimer
48		
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		

# Abstract

**Background** The rRT-PCR test, the current gold standard for the detection of coronavirus disease (COVID-19), presents with known shortcomings, such as long turnaround time, potential shortage of reagents, false-negative rates around 15–20%, and expensive equipment. The hematochemical values of routine blood exams could represent a faster and less expensive alternative.

**Methods** Three different training data set of hematochemical values from 1,624 patients (52% COVID-19 positive), admitted at San Raphael Hospital (OSR) from February to May 2020, were used for developing machine learning (ML) models: the complete OSR dataset (72 features: complete blood count (CBC), biochemical, coagulation, hemogasanalysis and CO-Oxymetry values, age, sex and specific symptoms at triage) and two sub-datasets (COVID-specific and CBC dataset, 32 and 21 features respectively). 58 cases (50% COVID-19 positive) from another hospital, and 54 negative patients collected in 2018 at OSR, were used for internal-external and external validation.

**Results** We developed five ML models: for the complete OSR dataset, the area under the receiver operating characteristic curve (AUC) for the algorithms ranged from 0.83 to 0.90; for the COVID-specific dataset from 0.83 to 0.87; and for the CBC dataset from 0.74 to 0.86. The validations also achieved good results: respectively, AUC from 0.75 to 0.78; and specificity from 0.92 to 0.96.

**Conclusions** ML can be applied to blood tests as both an adjunct and alternative method to rRT-PCR for the fast and cost-effective identification of COVID-19-positive patients. This is especially useful in developing countries, or in countries facing an increase in contagions.

# 1 Introduction

2 To date, at eight months post-outbreak, the coronavirus disease (COVID-19) caused by the SARS-CoV-2  
3 coronavirus has infected more than 20 million people and has resulted in approximately one million deaths  
4 worldwide. To manage this unprecedented pandemic emergency, the early identification of patients and of  
5 infectious people is extremely important due to the fact that this disease, unlike others caused by coronaviruses  
6 (e.g. SARS, MERS), can coexist in a host organism without causing any symptoms, or it can produce very mild  
7 and non-characteristic symptoms in—nevertheless—infectious subjects (1). To identify SARS-CoV-2 infections,  
8 the instrument of choice, or the gold standard, is the molecular test performed using the reverse polymerase chain  
9 reaction (PCR) or the reverse transcriptase–PCR (RT-PCR) technique. However, the execution of the test is time-  
10 consuming (at no less than 4–5 hours under optimal conditions), requires the use of special equipment and reagents,  
11 the involvement of specialized and trained personnel for the collection of the samples, and relies on the proper  
12 genetic conservation of the RNA sequences that are selected for annealing the primers (2). In addition, for these  
13 pre-analytical vulnerabilities (3), the RT-PCR test’s accuracy, and especially its sensitivity (i.e. its ability to avoid  
14 false negatives), is far from ideal. A recently published article in the *New England Journal of Medicine* suggests  
15 that a reasonable estimate for the sensitivity of this test is 70% (4).

16 To improve our diagnostic capabilities, in order to contain the spread of the pandemic, the data science community  
17 has proposed several machine learning (ML) models, recently reviewed in (5). Most of these models are based on  
18 computed tomography (CT) scans or chest x-rays (5–9). Despite the reported promising results, some concerns  
19 have been raised regarding these and other works, especially in regard to solutions based on chest x-rays, which  
20 have been associated with high rates of false-negative results (10). On the other hand, solutions based on CT  
21 imaging, although accurate, are affected by the characteristics of this modality: CTs are costly, time-consuming,  
22 and require specialized equipment; thus, approaches based on this imaging technique cannot reasonably be applied  
23 for screening exams. Although various clinical studies (11–13) have highlighted how blood test-based diagnostics  
24 might provide an effective and low-cost alternative for the early detection of COVID-19 cases, relatively few ML  
25 models have been applied to hematological parameters (14–18).

26 To overcome the above limitations, and following a successful feasibility study performed on a smaller dataset  
27 (19), we developed different classification models by applying ML techniques to blood-test results that are  
28 generally available in clinical practice within minutes (under emergency conditions, at even less than 15 minutes)  
29 and are only a fraction of the cost of the RT-PCR test and CT imaging (i.e. a few euros). As we will show, routine  
30 blood tests can be exploited by our method to diagnose COVID-19 patients in low-resource settings, in particular,  
31 where there is a shortage of RT-PCR reagents, such as during a pandemic peak. On the other hand, the developed  
32 method can also be used as a complement to the RT-PCR test in order to increase the sensitivity of the latter or to

1 provide its interpreters a sort of pre-test probability to compute NPV and PPV. Furthermore, the rapid blood-test  
2 results could be a valuable—although non-conclusive—indication for the early identification of COVID-19  
3 patients, resulting in their better management/isolation while waiting for the gold standard results.

## 4 5 **Methods**

6 In this section, we describe the datasets and statistical methods used to train and validate the ML models. The  
7 reporting follows the TRIPOD Guideline for Model Development and Validation (20). The study protocol  
8 (BIGDATA-COVID19) was approved by the Institutional Ethical Review Board in agreement with the World  
9 Medical Association Declaration of Helsinki.

## 10 **Data Description**

### 11 OSR dataset

12 The main dataset used for this study (the OSR dataset) consisted of routine blood-test results performed on 1,925  
13 patients on admission to the ED at the San Raffaele Hospital (OSR) from February 19, 2020, to May 31, 2020. In  
14 order to control for potentially confounding pathologies and other sources of bias, such as insufficient data  
15 availability, in ML development, 301 (15.6%) patients, admitted between February and April, were excluded from  
16 further analysis. All patients admitted during May 2020, on the other hand, were considered for the study, to have  
17 a balanced number of patients also from the late portions of the time frame considered.

18 For each case, COVID-19 positivity was determined based on the result of the molecular test for SARS-CoV-2  
19 performed by RT-PCR on nasopharyngeal swabs. On a set of 165 uncertain cases, we also used the result of chest  
20 radiography and x-rays to improve over the sensitivity of the RT-PCR test (21–24). Uncertain cases were identified  
21 through two different methods: either patients who resulted positive within 72 hours after a first negative test  
22 and were admitted as inpatients despite this test result; or patients who, despite having a negative test, had an  
23 hematochemical profile more similar to positive patients, as determined through multi-variate clustering based  
24 on a set of COVID-19 characteristic biomarkers (12) (AST, lymphocytes, calcium, LDH, PCR, WBC, XDP,  
25 Fibrinogen). Of the 165 uncertain cases, only 52 of them have been considered as positive after comparison with  
26 the radiologic gold standard, while the remaining 113 were considered as negative (having a double negative  
27 test from both the RT-PCR and the radiologic gold standard): this results in an estimate of 93% sensitivity of the  
28 RT-PCR with respect to the composite ground truth.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59

1 Therefore, the OSR dataset consisted of a total of 1,624 cases: 786 of them received a positive diagnosis (48%)  
2 and 838 were negative cases (52%).

3 As covariate features, for each case, the patient's age and gender, the presence of COVID-19 related  
4 symptomatology at admission (dyspnea, pneumonia, pyrexia, sore throat, influenza, cough, pharyngitis, bronchitis,  
5 generalized illness), and a set of 69 hematochemical values from laboratory tests were considered. The list of the  
6 analytes and instruments are reported in Table 1. The laboratory blood tests were performed according to the  
7 International Federation of Clinical Chemistry and Laboratory Medicine (IFCC) recommendations (25).

8 The demographic and clinical characteristics of the two different groups of COVID-19 patients are summarized in  
9 Figure 1 and Figure 2.

10 The missing data rate for each of the examined features is reported in Table 1. In order to reduce the bias due to  
11 imputation, we discarded all the features with a missing data rate greater than 75%. Thus, among the 1,624 cases  
12 in the OSR dataset, 1,189 (73%) cases had at least 75% of the attributes; while 1,324 (82%) cases had complete  
13 data for the CBC features.

14 From the complete *OSR dataset*, we obtained two other datasets by selecting two relevant subsets of the features  
15 (thus, the three datasets share the same set of patients):

- 16 1 A dataset consisting of the 34 features under the column header “COVID-specific features” (see Table 1),  
17 denoted as the *COVID-specific dataset*.
- 18 2 A dataset consisting of the 21 features under the column header “CBC features” (see Table 1), denoted as  
19 the *CBC dataset*.

## 20 External datasets

21 In addition to the previously described datasets, all obtained from the OSR dataset, we considered two external  
22 datasets for the internal–external validation and for the external validation of the models.

23 The first dataset, the Istituto Ortopedico Galeazzi (*IOG dataset*), was obtained from blood samples collected at the  
24 ED of the IOG of Milan between March 5, 2020, and May 26, 2020, and encompassed the parameters under the  
25 “COVID-specific features” column header (see Table 1). Notably, this hospital specializes in the diagnosis and  
26 treatment of musculoskeletal disorders and was not considered a destination of choice during the acute phase of  
27 the pandemic in the Milan area. Therefore, the patients were presumably of a different severity and were admitted  
28 for other reasons than pulmonary conditions with respect to OSR. The IOG dataset consisted of a total of 58 cases,  
29 29 with negative swab results and 29 with positive swab results, and with the same features as the *COVID-specific*  
30 *dataset*. For the IOG and OSR, different instruments were used for the CBC and biochemical parameters; in  
31 particular, the iSystem XN-2000 system was used instead of the Sysmex XE 2100 system for CBC counts and the

1  
2  
3 1 Atellica® CH Analyzer (Siemens Healthineers) was used instead of the Roche COBAS 6000 system for the  
4  
5 2 biochemical parameters.

6  
7 3 The second dataset (the *2018 dataset*) was obtained from blood samples collected at the OSR in November 2018  
8  
9 4 from 54 randomly chosen patients. These were obviously negative for COVID-19: 20 (37%) of them were  
10  
11 5 specifically chosen to act as confounding cases, as they exhibited pneumonia-like symptoms.

## 12 6 Machine Learning Experimental Design

13  
14  
15 7 We implemented a four-step pipeline for ML model development encompassing imputation, data normalization,  
16  
17 8 feature selection, and classification. The data analysis pipeline was implemented in Python (version 3.7), using  
18  
19 9 the numpy (version 1.19), pandas (version 1.1) and scikit-learn (version 0.23) libraries. For imputation, the  
20  
21 10 multivariate k-nearest neighbors algorithm was used (26), with  $k = 5$ . For feature-selection, the recursive feature-  
22  
23 11 elimination algorithm was used (27). The optimal features to select were determined through hyper-parameter  
24  
25 12 optimization. For classification we evaluated five different algorithms: Random Forest (RF), naive Bayes (NB),  
26  
27 13 logistic regression (LR), support vector machine (SVM), and k-nearest neighbors (KNN). We specifically  
28  
29 14 evaluated these algorithms as all have been shown to achieve state-of-the-art performance on tabular data (28)  
30  
31 15 and, at least to some degree (for example using feature-attribution methods), interpretable (29). The hyper-  
32  
33 16 parameters of the different classification algorithms are reported in Suppl. Table 1. All hyper-parameters were  
34  
35 17 optimized automatically using a grid search approach.

36  
37 18 In regard to model selection, training and evaluation, we performed a two-step procedure to minimize the risk of  
38  
39 19 over-fitting: first, the dataset was split into a training set (80% of the instances) and a hold-out test set (20% of the  
40  
41 20 instances), using a stratified procedure; second, hyper-parameter optimization was performed (on the training set)  
42  
43 21 through 5-fold stratified cross-validation grid search and using AUC as reference measure; third, the models were  
44  
45 22 trained and calibrated on the whole training set; finally, the calibrated models were evaluated on the hold-out test  
46  
47 23 set in terms of accuracy, sensitivity, specificity, AUC, and the Brier score (30) (a standard metric to measure the  
48  
49 24 models' calibration, with a lower score being better). In all stages of model development, the randomization was  
50  
51 25 controlled in order to ensure repeatability of the experiments.

52  
53 26 For each model class, we considered two versions: a standard one, and the three-way version (a model that abstains  
54  
55 27 from prediction when the confidence score is below 75%) (31). For each of these two versions, the model selection,  
56  
57 28 training and evaluation pipeline was implemented for each of the three datasets mentioned above (the *OSR dataset*,  
58  
59 29 *COVID-specific dataset*, and *CBC dataset*).

30  
31 30 The *IOG dataset* and the *2018 dataset* were used, respectively, for the internal–external and external validation of  
the models developed for the *COVID-specific dataset* and the *CBC dataset*.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59

1 The internal–external validation procedure—the purpose of which was to evaluate the models’ ability to generalize  
2 to a new setting when provided with a limited quantity of new data—was implemented using a bootstrap-based  
3 approach (see Supplementary Materials - Implementation of the Internal-External Validation).

4 The external validation procedure—the purpose of which was to test both the specificity of the developed models  
5 and their ability to identify potential suspect cases—was implemented by training the best models found for the  
6 *COVID-specific dataset* (respectively, the *CBC dataset*) on the combined dataset that also consisted of the *IOG*  
7 *dataset* and then evaluating the trained models against the *2018 dataset*.

8 The combined dataset consisting of the *COVID-specific* (respectively, the *CBC dataset*) and the *IOG* datasets were  
9 also used to evaluate the sensitivity and specificity for symptomatic and asymptomatic patients separately: in this  
10 case, the models were retrained after deletion of the Suspect feature (to avoid bias) and the re-trained models were  
11 then evaluated on symptomatic and asymptomatic patients (both from the test set) separately.

## 12 13 Results

14 The results of the ML models on the three datasets (OSR, COVID-specific, CBC) are reported in Table 2.

15  
16 The receiver operating characteristic (ROC) curves of the best model (in terms of the highest AUC) for each of  
17 the three datasets is reported in Figure 3. The ROC curves for all models (on each of the three datasets) are reported  
18 in Suppl. Figures 1, 2 and 3. The feature importance scores—which were computed in order to enable the  
19 interpretability of the developed models—are reported in Figure 4 and in Suppl. Figure 13 for the best model of  
20 each of the three datasets. The positive predictive value (PPV)-sensitivity curves are reported in Suppl. Figures 4,  
21 5 and 6, while the calibration curves are reported in Suppl. Figures 7, 8, and 9, and the PPV/NPV prevalence  
22 curves are reported in Suppl. Figures 10, 11 and 12.

23  
24 The results for the internal–external validation and the external validation (specificity only) procedures are  
25 reported in Table 3. In this table we highlight the results of the models that obtained the best performance in the  
26 internal validation (KNN for the COVID-specific dataset; and both KNN and RF for the CBC dataset).  
27 Specifically, in the first 4 columns we report the results of the internal-external validation (in terms of accuracy,  
28 sensitivity, specificity and AUC), while in the last column we report the results of the external validation (in terms

1  
2  
3 of specificity). In Table 3, we also report on the performance of the models for asymptomatic patients and  
4 symptomatic patients, as described in the Methods section.  
5  
6  
7  
8

## 9 10 Discussion

11  
12  
13 The unprecedented worldwide public health emergency caused by the COVID-19 pandemic has motivated  
14 different research groups to develop ML applications with the aim of automating—at least partially—the diagnosis  
15 or screening of COVID-19.  
16  
17

18  
19 Nonetheless, only a few studies have focused on the development of ML models based on routine blood exams.  
20 Formica et al.(13) developed a CBC-based ML model, reporting 83% sensitivity and 82% specificity; however,  
21 the analysis was based on a small sample (171 patients) collected in a limited time frame (between March 7 and  
22 March 19, 2020). Banerjee et al. (32) applied ML methods to a public dataset of CBC data encompassing 598  
23 cases of which only 39 cases were COVID-19 positive; the authors report good specificity (91%) but very low  
24 sensitivity (43%), thus making the proposed model unsuitable for early detection tasks. Further, this work presents  
25 some major limitations affecting replicability and generalizability, as the authors do not provide any information  
26 regarding how the values of the considered features were measured (analytical instruments, analytical principle,  
27 and units of measurement). Avila et al. (33) used the same dataset considered in (32) to develop a Bayesian model,  
28 reporting 76.7% sensitivity and specificity. Notably, the authors report a number of complete instances (510)  
29 which is different from that reported in (32). Joshi et al. (34) developed a logistic regression model trained using  
30 CBC data on a dataset of 380 cases, reporting good sensitivity (93%) but low specificity (43%).  
31  
32  
33  
34  
35  
36  
37

38  
39 More in general, a recent critical survey (5) raised some concerns about these and other evaluated studies (most  
40 of which have not yet undergone peer-review), noting the possibility of high rates of bias and over-fitting, and  
41 little compliance with reporting and replication guidelines (18).  
42

43  
44 Finally, a recent study Yang et al. (35), considered the development of a Gradient Boosting model on a set of 3,356  
45 patients (42% COVID-19 positive) using a set of 27 parameters encompassing both blood count and biochemical  
46 parameters, achieving 0.85 AUC, and also reporting a comparable result (AUC 0.84) for validation on an external  
47 dataset. This work can be viewed as similar but complementary with respect to the results that we report, both in  
48 terms of considered features and used laboratory instrumentation (the authors used the UniCel DXH 800 analyzer  
49 for the CBC features, and Siemens ADVIA XPT analyzers for biochemical parameters). Indeed, compared with  
50 the parameters considered in this study, the authors of (35) considered albumin, total protein, magnesium, ferritin  
51 and globulin; but lacked a set of parameters (some of which known to be significantly altered in COVID-19  
52 patients), such as creatinine (CREA), aspartate aminotransferase (AST), alanine aminotransferase (ALT), Gamma-  
53  
54  
55  
56  
57  
58  
59

1 glutamyl transferase (GGT), Creatin-kinase (CK), Potassium (K), Interleukin 6 (IL6), NT-proB-type Natriuretic  
2 Peptide (ProBNP), total (BilT) and direct (BilD) Bilirubin, all coagulation tests, hemogasanalysis and CO-  
3 Oxymetry parameters. We think that this complementarity in the two studies could lend support to the usefulness  
4 of blood tests as an alternative approach for COVID-19 diagnosis.

5 To overcome the limitations of the above models, we applied the ML methodology to routine blood examination  
6 outcomes, which are usually available for inpatients and for patients admitted to the ED in shorter time frames and  
7 at much lower cost than both molecular tests and radiological exams. In this endeavor, we addressed three subtasks:  
8 1) detecting COVID-19 from a full battery of hematochemical tests, commonly collected from suspected  
9 respiratory tract disease patients (OSR dataset); 2) detecting COVID-19 from only a restricted subset of parameters  
10 known to be altered in COVID-19 patients (COVID-specific dataset); and 3) detecting COVID-19 from a very  
11 small subset of hematological parameters (CBC and the WBC differential) representing the basic routine blood  
12 examinations, usually also available in low-resource settings (CBC dataset). For each of the datasets described  
13 above, we applied five different models that were selected from among those that are more frequently adopted in  
14 medical ML.

15 These models achieve COVID-19 detection in different ways and exhibit good performance, although they are  
16 associated with different sensitivities and specificities. This makes them good candidates for embedding in an  
17 online service (which we are currently developing) in which doctors can specify their preferences (with respect to  
18 greater sensitivity, greater specificity, or a balanced performance (36) and needs according to their diagnostic  
19 purpose (i.e., screening, triage, or a secondary diagnosis), and thus gain an indication from the optimal model. In  
20 addition, the users can decide whether they want an indication from the system—irrespective of the confidence in  
21 the advice given—or if they would prefer only to be advised about high-confidence indications, as the three-way  
22 approach allows for. This approach was specifically developed to mitigate the risk of automation bias and the odds  
23 of machine-induced errors (31).

24 With respect to the patterns used to discriminate between positive and negative cases, the ML models identified,  
25 as the most predictive features, those parameters that are known to be significantly altered in COVID-19 patients  
26 (5,37-38) (see Figure 4 and Suppl. Figure 13). For instance, when applied to the OSR and COVID-specific datasets,  
27 the models identify lactate dehydrogenase (LDH), AST, C-reactive protein (CRP), and calcium as the most  
28 important features, while all the models also reported WBC and its corresponding differential as important. Also  
29 the patients' age was reported by the models to be a significant predictor, which is consistent with the literature  
30 (5), where it was also found to be a significant predictor not only for prognostic, but also for diagnostic tasks.  
31 Notably, fibrinogen and cross-linked fibrin degradation products (XDPs), known to be associated with COVID-  
32 19 severity (39), were also considered by the model as being among the most important features when applied to  
33 the OSR dataset.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59

1 With respect to the calibration of the developed models, the good internal calibration of the models can be  
2 confirmed by the calibration curves in Suppl. Figures 7–9.

3 As can be seen in Table 3, the internal–external validation and the external validation procedures also achieved  
4 good results. In this respect, it is important to note that the validation procedures involved blood tests performed  
5 on different types of analytical instrumentation for the clinical chemistry tests (Siemens instead of Roche),  
6 although the CBC standardization was less problematic than the other tests. For this reason, as ML models exhibit  
7 poor performance when considering out-of-distribution samples (40), the goal of the internal–external validation  
8 process was to assess the capability of the models to generalize across different settings. All of the models showed  
9 good performance and, more specifically, good specificity. The models achieved good performances in  
10 symptomatic patients (with both the sensitivity and specificity at approximately 80%) and they performed even  
11 better in terms of specificity in asymptomatic patients (100% specificity), although the sensitivity was as low as  
12 50% (see Table 3). Nevertheless, considering that the developed ML-based tests were based on low-cost and rapid  
13 blood-test examinations, the reported values can be considered good enough, specifically in regard to screening  
14 (16).

15 The external validation procedure also achieved very good results (at around 95% for all models in the standard  
16 version), but it should be noted that this only relates to COVID-19-negative patients, and hence, to specificity.  
17 Notably, in the external validation process, all five patients for which the models failed had symptoms that were  
18 compatible with COVID-19 disease.

19 As hinted at above, the outputs from our models can be used in different scenarios. They could be used together  
20 and combined with the molecular test to obtain a compound test with higher accuracy, and, most importantly,  
21 higher sensitivity regarding suspected cases, thus allowing for the identification of a larger number of COVID-19-  
22 positive patients so that they can be isolated and treated in a timely manner. Indeed, we can see in Table 3 that the  
23 sensitivity in symptomatic patients is adequate for this type of use. In the same vein, the models' outputs could be  
24 used while waiting for the results from other tests, allowing for the timely and prudent management of suspected  
25 COVID patients, or in screening and pool-testing scenarios (41) where low accuracy is not a critical problem if a  
26 test such as a CBC can be performed frequently (42). In Table 3 we can see that the model that was defined based  
27 on the smallest dataset (the CBC dataset) reaches 100% specificity in asymptomatic patients. Consequently, we  
28 are planning to use our model for epidemiological purposes on the blood donor population to estimate the  
29 prevalence of the condition in the asymptomatic population. On the other hand, the scenarios in which the results  
30 from our models replace those of the molecular tests address an emergency need, especially if the time to obtain  
31 the molecular test results is too long (due to a high demand for such tests in an outbreak area), if there is a shortage  
32 of materials (swabs or reagents) for any supply problem, or in poor health contexts or in contexts where there are  
33 serious structural deficiencies (such as in some developing countries or in a geographical area that is, in the

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59

1 meantime, affected by other socio-sanitary and humanitarian emergencies). In these situations where resources are  
2 limited and population-wide testing cannot be performed, CBC-based scores may help to pre-evaluate patients and  
3 activate COVID-19-specific pathways and molecular testing for patients with high scores independent of symptom  
4 severity. In the presence of suspected COVID-19 cases and high scores, logistical management can promptly  
5 activate isolation procedures (43).

## 6 7 Conclusion

8 All things considered, the ML models that we presented in this article achieved a performance that is comparable,  
9 although inferior, to RT-PCR (4), which is the current gold standard for COVID-19 diagnosis. Nevertheless,  
10 although our models are less accurate, they aim to be an additional tool available among those that, being much  
11 faster and cheaper than the current diagnostic reference tests, can be used for the screening of whole populations.  
12 This use can facilitate the shift in testing strategy that, grounding on a faster, although less accurate, identification  
13 of infected individuals, is said to have a positive potential in slowing the virus' spread and contributing for the safe  
14 reopening of schools and workplaces (44).

## **Data Availability**

The datasets collected and used in this study are available from the corresponding author on reasonable request.

## **Author Contributions (in alphabetical order):**

Giuseppe Banfi: conceived and designed the study, revised, and approved the manuscript

Federico Cabitza conceived and designed the study, analyzed, and interpreted the data, wrote and revised the manuscript, approved the manuscript.

Andrea Campagner: analyzed and interpreted the data, wrote and revised the manuscript, approved the manuscript.

Anna Carobene: conceived and designed the study, provided study materials or patients, collected and organized the data, wrote and revised the manuscript, approved the manuscript

Daniele Ceriotti: provided study materials or patients, revised and approved the manuscript

Alessandra Colombini: provided study materials or patients, revised and approved the manuscript

Elena De Vecchi: provided study materials or patients, revised and approved the manuscript

Chiara Di Resta: collected and organized the data, wrote and revised the manuscript, approved the manuscript

Davide Ferrari: conceived and designed the study, revised and approved the manuscript

Massimo Locatelli: conceived and designed the study, revised and approved the manuscript

Eleonora Sabetta: provided study materials or patients, revised and approved the manuscript

## **Competing Interests statement:**

The authors declare no competing interests.

1  
2  
3 **Figure captions**

4  
5  
6 **Figure 1.** Demographic feature (gender and age) distributions for positive and negative cases. The blue and orange  
7 areas correspond to negative and positive cases, respectively.

8  
9  
10 **Figure 2.** Violin plots depicting the distributions of eight relevant features in the *OSR dataset* (selected for their  
11 predictivity toward COVID-19). The blue and orange areas correspond to negative and positive cases, respectively.

12  
13  
14 **Figure 3.** Receiver operating characteristic curves for the best models (in terms of AUC), for each of the three  
15 considered datasets (*OSR*, *COVID-specific*, *CBC*). For the *CBC* dataset we report the ROC curve for both RF and  
16 KNN as they had equal AUC (see Table 2).

17  
18  
19 **Figure 4.** Feature importance scores for the random forest algorithm trained using the *OSR dataset* (on the left)  
20 and the k-nearest neighbors' algorithm trained using the *COVID-specific dataset* (on the right).

## References

1. Oran DP, Topol EJ. Prevalence of Asymptomatic SARS-CoV-2 Infection: A Narrative Review. [Published online June 3, 2020]. *Ann intern med.* doi:10.7326/M20-3012
2. Vogels CBF, Brito AF, Wyllie AL, Fauver JR, Ott IM, Kalinich CC, et al. Analytical sensitivity and efficiency comparisons of SARS-CoV-2 RT-qPCR primer-probe sets. [Published online July 10, 2020]. *Nat Microbiol.* doi:10.1038/s41564-020-0761-6
3. Lippi G, Simundic A-M, Plebani M. Potential preanalytical and analytical vulnerabilities in the laboratory diagnosis of coronavirus disease 2019 (COVID-19). *Clin Chem Lab Med.* 2020;58:1070-6
4. Woloshin S, Patel N, Kesselheim AS. False Negative Tests for SARS-CoV-2 Infection — Challenges and Implications. *N Engl J Med.* 2020;383:e38. doi: 10.1056/NEJMp2015897.
5. Wynants L, Van Calster B, Collins GS, Riley RD, Heinze G, Schuit E, et al. Prediction models for diagnosis and prognosis of covid-19: Systematic review and critical appraisal. *BMJ.* 2020; 369:m1328. doi: 10.1136/bmj.m1328
6. Li L, Qin L, Xu Z, Yin Y, Wang X, Kong B, et al. Artificial Intelligence Distinguishes COVID-19 from Community Acquired Pneumonia on Chest CT. [Published online April 3, 2020]. *Radiology.* doi:10.1148/radiol.2020200905
7. Gozes O, Frid-Adar M, Greenspan H, Browning PD, Zhang H, Ji W, et al. Rapid AI Development Cycle for the Coronavirus (COVID-19) Pandemic: Initial Results for Automated Detection & Patient Monitoring using Deep Learning CT Image Analysis. [Published online March 24, 2020]. *arXiv Prepr arXiv* <http://arxiv.org/abs/2003.05037>
8. Ozturk T, Talo M, Yildirim EA, Baloglu UB, Yildirim O, Rajendra Acharya U. Automated detection of COVID-19 cases using deep neural networks with X-ray images. *Computers in Biology and Medicine* 121 (2020) 103792. doi:10.1016/j.combiomed.2020.103792
9. Mei X, Lee HC, Diao K, Huang M, Lin B, Liu C, et al. Artificial intelligence-enabled rapid diagnosis of patients with COVID-19. *Nat Med* 26, 1224–1228 (2020).
10. Weinstock MB, Echenique A, Russell JW, Leib A, Miller J, Cohen DJ, et al. Chest X-Ray Findings in 636 Ambulatory Patients with COVID-19 Presenting to an Urgent Care Center: A Normal Chest X-Ray Is no Guarantee. [Published online May, 2020]. *JUCM.* <https://www.jucm.com/documents/jucm-covid-19-studypub-april-2020.pdf>. Accessed August 17, 2020
11. Fan BE, Chong VCL, Chan SSW, Lim GH, Tan GB, Mucheli SS et al. Hematologic parameters in patients with COVID-19 infection. *Am J Hematol.* 2020;95:E131-E134.
12. Ferrari D, Motta A, Strollo M, Banfi G, Locatelli M. Routine blood tests as a potential diagnostic tool for COVID-19. *Clin Chem Lab Med.* 2020;58:1095-9.
13. Formica V, Minieri M, Bernardini S, Ciotti M, D'Agostini C, Roselli M, et al. Complete blood count might help to identify subjects with high probability of testing positive to SARS-CoV-2. *Clin Med (Lond).* 2020;20:e114-e119.



14. Wu J, Zhang P, Zhang L, Meng W, Li J, Tong C, et al. Rapid and accurate identification of COVID-19 infection through machine learning based on clinical available blood test results. [Published online 2020].medRxiv. doi:10.1101/2020.04.02.20051136
15. Soares F. A novel specific artificial intelligence-based method to identify {COVID}-19 cases using simple blood exams. [Published online 2020]. medRxiv. <https://www.medrxiv.org/content/10.1101/2020.04.10.20061036v2>
16. Soltan AAS, Kouchaki S, Zhu T, Kiyasseh D, Taylor T, Hussain ZB, et al. Artificial intelligence driven assessment of routinely collected healthcare data is an effective screening test for COVID-19 in patients presenting to hospital. [Published online 2020]. medRxiv.doi:10.1101/2020.07.07.20148361
17. Kukar M, Gunčar G, Vovko T, Podnar S, Černelč P, Brvar M, et al. COVID-19 diagnosis by routine blood tests using machine learning. [Published online June 2020]. arXiv Prepr arXiv.<http://arxiv.org/abs/2006.03476>. Accessed August 17, 2020
18. Collins GS, Moons KGM. Reporting of artificial intelligence prediction models. *Lancet*. 2019;393:1577-9.
19. Brinati D, Campagner A, Ferrari D, Locatelli M, Banfi G, Cabitza F. Detection of COVID-19 Infection from Routine Blood Exams with Machine Learning: A Feasibility Study. *J Med Syst*. 2020;44:135.
20. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The TRIPOD Statement. *Eur Urol*. 2015;131:211-9.
21. Watson J, Whiting PF, Brush JE. Interpreting a covid-19 test result. [Published online May 12, 2020].*BMJ*. doi:10.1136/bmj.m1808
22. Zitek T. The appropriate use of testing for Covid-19. *West J Emerg Med*. 2020 Apr 13;21:470-2.
23. Fang Y, Zhang H, Xie J, Lin M, Ying L, Pang P, et al. Sensitivity of Chest CT for COVID-19: Comparison to RT-PCR. *Radiology*. 2020;296:E115-E117.
24. Liu J, Yu H, Zhang S. The indispensable role of chest CT in the detection of coronavirus disease 2019 (COVID-19). *Eur J Nucl Med Mol Imaging*. 2020;47:1638-9.
25. Bohn MK, Lippi G, Horvath A, Sethi S, Koch D, Ferrari M, et al. Molecular, serological, and biochemical diagnosis and monitoring of COVID-19: IFCC taskforce evaluation of the latest evidence. *Clin Chem Lab Med*. 2020 25;58:1037-52.
26. Jadhav A, Pramod D, Ramanathan K. Comparison of Performance of Data Imputation Methods for Numeric Dataset. *Appl Artif Intell*. 2019;10:913-33
27. Guyon I, Weston J, Barnhill S, Vapnik V. Gene selection for cancer classification using support vector machines. *Mach Learn*. 2002;46:389–422.
28. Caruana R, Karampatziakis N, Yessensalina, A. An empirical evaluation of supervised learning in high dimensions. *Proceedings of the 25th ICML*. 2008;96-103
29. Du M, Liu N, Hu X. Techniques for interpretable machine learning. *Communications of the ACM*. 2019;63:68-77.

30. Brier GW. Verification of Forecasts Expressed in Terms of Probability. *Mon Weather Rev.* 1950;78:1-3.
31. Campagner A, Cabitza F, Ciucci D. The three-way-in and three-way-out framework to treat and exploit ambiguity in data. 2020;119: 292-312.
32. Banerjee A, Ray S, Vorselaars B, Kitson J, Mamalakis M, Weeks S, et al. Use of Machine Learning and Artificial Intelligence to predict SARS-CoV-2 infection from Full Blood Counts in a population. [Published online June 16, 2020]. *Int Immunopharmacol.* doi:10.1016/j.intimp.2020.106705
33. Avila E, Kahmann A, Alho C, Dorn M. Hemogram data as a tool for decision-making in COVID-19 management: applications to resource scarcity scenarios. [Published online June 29, 2020]. *PeerJ.* doi:10.7717/peerj.9482
34. Joshi RP, Pejaver V, Hammarlund NE, Sung H, Lee SK, Furmanchuk A, et al. A predictive tool for identification of SARS-CoV-2 PCR-negative emergency department patients using routine test results. *J Clin Virol.* 2020;129:104502. doi: 10.1016/j.jcv.2020.104502.
35. Yang HS, Vasovic L V, Steel P, Chadburn A, Hou Y, Racine-Brzostek SE, et al. Routine laboratory blood tests predict SARS-CoV-2 infection using machine learning. *Clin Chem.* [Published online August 21, 2020]. doi:10.1093/clinchem/hvaa200
36. Cabitza F, Campagner A, Ciucci D, Seveso A. Programmed Inefficiencies in DSS-Supported Human Decision Making. In: *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. 2019. doi:10.1007/978-3-030-26773-5\_18
37. Rodriguez-Morales AJ, Cardona-Ospina JA, Gutiérrez-Ocampo E, Villamizar-Peña R, Holguin-Rivera Y, Escalera-Antezana JP, et al. Clinical, laboratory and imaging features of COVID-19: A systematic review and meta-analysis. *Travel Med Infect Dis.* 2020;34:101623. doi:10.1016/j.tmaid.2020.101623
38. Zhang ZL, Hou YL, Li DT, Li FZ. Laboratory findings of COVID-19: a systematic review and meta-analysis. [Published online May 23, 2020] *Scand J Clin Lab Invest.*;1-7. doi: 10.1080/00365513.2020.1768587
39. Connors JM, Levy JH. COVID-19 and its implications for thrombosis and anticoagulation. *Blood.* 2020;135:2033-40.
40. Rabanser S, Günemann S, Lipton ZC. Failing Loudly: An Empirical Study of Methods for Detecting Dataset Shift. 2018;(NeurIPS). <http://arxiv.org/abs/1810.11953>
41. Augenblick N, Kolstad JT, Obermeyer Z, Wang A. Group Testing in a Pandemic: The Role of Frequent Testing, Correlated Risk, and Machine Learning. *Natl Bur Econ Res.* 2020; <http://www.nber.org/papers/w27457.pdf>
42. Larremore DB, Wilder B, Lester E, Shehata S, Burke JM, Hay JA, et al. Test sensitivity is secondary to frequency and turnaround time for COVID-19 surveillance. [Published online 2020]. medRxiv. doi:10.1101/2020.06.22.20136309
43. Song JY, Yun JG, Noh JY, Cheong HJ, Kim WJ. Covid-19 in South Korea - Challenges of subclinical manifestations. *N Engl J Med.* 2020; 382:1858-9
44. Service R. Fast, cheap tests could enable safer reopening. *Science.* 2020;369:608-9.

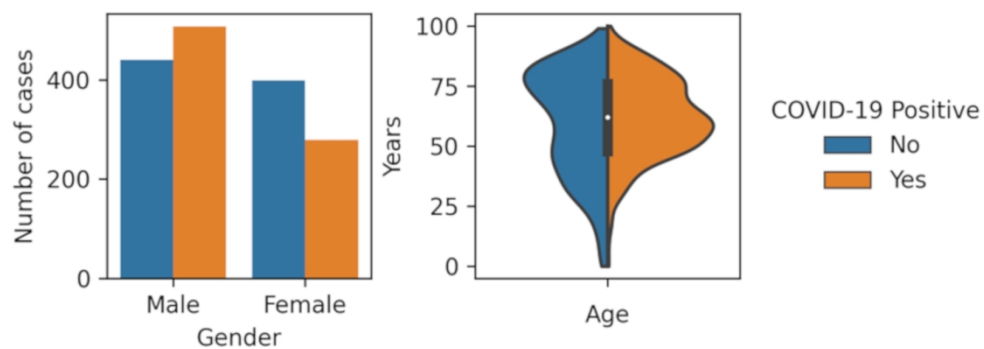


Figure 1. Demographic feature (gender and age) distributions for positive and negative cases. The blue and orange areas correspond to negative and positive cases, respectively

127x48mm (300 x 300 DPI)

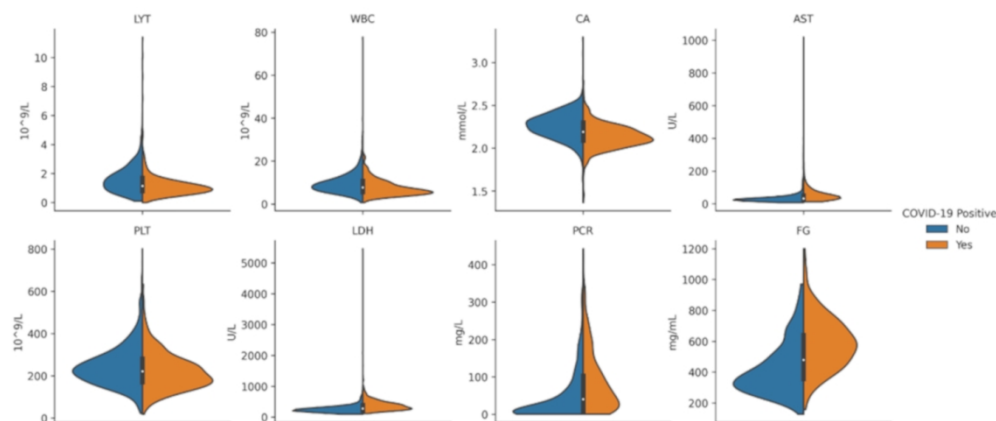


Figure 2. Violin plots depicting the distributions of eight relevant features in the OSR dataset (selected for their predictivity toward COVID-19). The blue and orange areas correspond to negative and positive cases, respectively.

127x56mm (300 x 300 DPI)

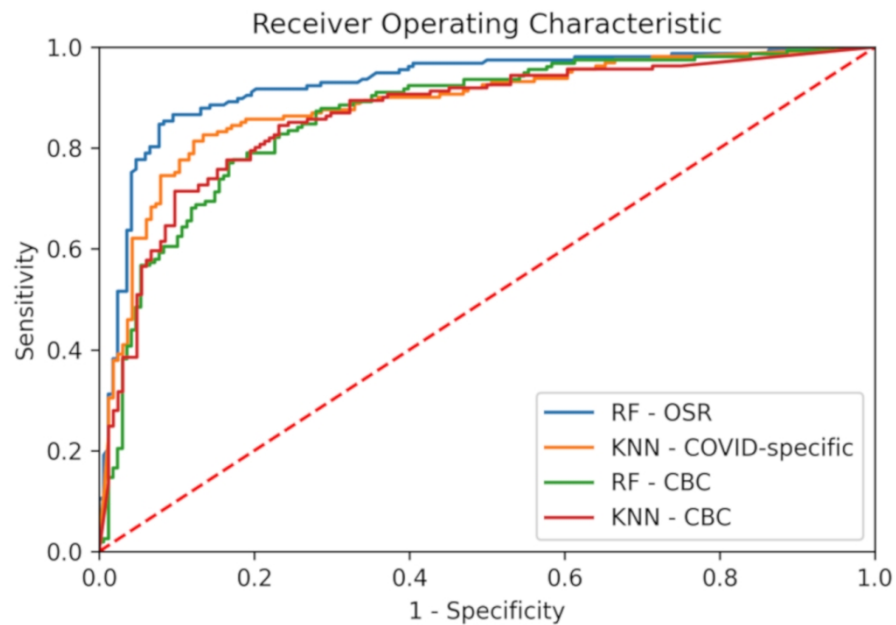


Figure 3. Receiver operating characteristic curves for the best models (in terms of AUC), for each of the three considered datasets (OSR, COVID-specific, CBC). For the CBC dataset we report the ROC curve for both RF and KNN as they had equal AUC (see Table 2).

127x84mm (300 x 300 DPI)

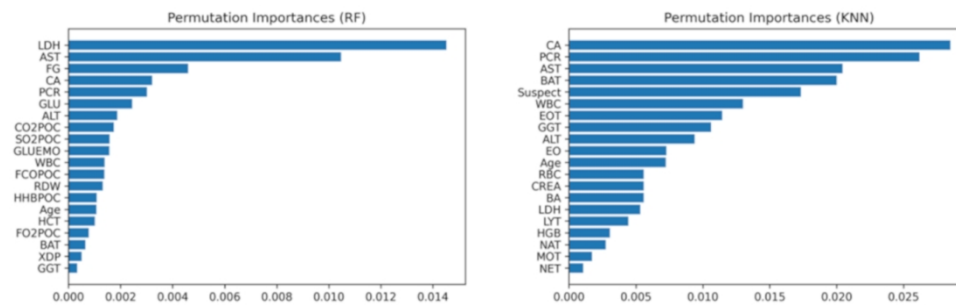


Figure 4. Feature importance scores for the random forest algorithm trained using the OSR dataset (on the left) and the k-nearest neighbors' algorithm trained using the COVID-specific dataset (on the right).

127x40mm (300 x 300 DPI)

**Table 1.** Complete list of the analyzed features in the *OSR dataset*  
 medRxiv preprint doi: <https://doi.org/10.1101/2020.10.02.20205070>; this version posted October 4, 2020. The copyright holder for this preprint  
 (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.  
 It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

Category	Instrument / sample	Parameter	Acronym	Unit of Measure	COVID-specific features	CBC features	Missing rate (%)
Hematological	Sysmex XE 2100 / whole blood	White blood cells	WBC	10 <sup>9</sup> /L	X	X	2.4
		Red Blood Cells	RBC	10 <sup>12</sup> /L	X	X	3.6
		Hemoglobin	HGB	g/dL	X	X	2.4
		Hematocrit	HCT	%	X	X	2.4
		Average Globular Volume	MCV	fL	X	X	3.6
		Mean Corpuscular Hemoglobin	MCH	pg/Cell	X	X	3.6
		Mean Corpuscular Hemoglobin Concentration	MCHC	g Hb/dL	X	X	2.4
		Erythrocyte distribution width	RDW	CV%	X	X	3.7
		Platelets	PLT	10 <sup>9</sup> /L	X	X	3.6
		Average Platelet Volume	MPV	fL	X	X	5.9
		Neutrophils count (%)	NE	%	X	X	18.9
		Lymphocytes count (%)	LY	%	X	X	15.2
		Monocytes count (%)	MO	%	X	X	15.2
		Eosinophils count (%)	EO	%	X	X	15.2
		Basophils count (%)	BA	%	X	X	15.2
		Neutrophils count	NET	10 <sup>9</sup> /L	X	X	15.2
		Lymphocytes count	LYT	10 <sup>9</sup> /L	X	X	15.2
		Monocytes count	MOT	10 <sup>9</sup> /L	X	X	18.9
Eosinophils count	EOT	10 <sup>9</sup> /L	X	X	15.2		
Basophils count	BAT	10 <sup>9</sup> /L	X	X	18.9		
Coagulation	STA - R MAX/ Plasma sample	Prothrombin Time (INR)	PTINR	INR			31.0
		Activated partial thromboplastin time ( R )	PPTR	Ratio			31.5

Fibrinogen	FG	mg/dL		70.2
D-Dimer	XDP	µg/mL		70.4

		Biochemical			
Cobas 6000 Roche/ serum sample	Glucose	GLU	mg/dL	X	3.4
	Creatinine	CREA	mg/dL	X	2.4
	Urea	UREA	mg/dL	X	37.0
	Direct Bilirubin	BILD	mg/dL		23.3
	Indirect Bilirubin	BILIN	mg/dL		23.3
	Total Bilirubin	BILT	mg/dL		25.3
	Alanine Aminotransferase	ALT	U/L	X	3.1
	Aspartate Aminotransferase	AST	U/L	X	3.2
	Alkaline Phosphatase	ALP	U/L	X	23.7
	Gamma Glutamyltransferase	GGT	U/L	X	24.5
	Lactate Dehydrogenase	LDH	U/L	X	13.2
	Creatine kinase	CK	U/L	X	60.3
	Sodium	NA	mmol/L	X	3.9
	Potassium	K	mmol/L	X	2.7
	Calcium	CA	mmol/L	X	3.8
	C-reactive Protein	CRP	mg/L	X	5.5
	NT-proB-type Natriuretic Peptide	PROBNP	pg/mL		91.1
	Troponin T	TROPOT	ng/L		62.8
Interleukin 6	IL6	pg/mL		92.2	
Rapidpoint 500 (Siemens Healthcare)	Hemogasanalysis, venous blood gas	pH	PHPOC	U	18.5
		Carbonic Anhydride (pCO2)	CO2POC	mmHg	22.4
		Oxygen (pO2)	PO2POC	mmHg	22.4



CO-Oxymetry	Bicarbonates	BICPOC	mmol/L	18.7
	Standard Calculated Bicarbonates	BISPOC	mmol/L	23.0
	Base Excess	BEPOC	mmol/L	22.8
	Actual Base Excess	BEEPOC	mmol/L	18.9
	Hematocrit (POC)	HCTPOC	%	22.7
	Total Oxyhemoglobin	THBPOC	g/dL	22.8
	O2 Saturation	SO2POC	%	18.3
	Oxyhemoglobin / Total Hemoglobin	FO2POC	%	18.6
	Carboxyhemoglobin	FCOPOC	%	18.8
	Methemoglobin	METPOC	%	22.5
	Deoxyhemoglobin	HHBPOC	%	18.8
	Oxygenation	Bound O2 Maximum Concentration	BO2POC	mL/dL
Total Oxygen		CTOPOC	mL/dL	20.9
Inspired Oxygen Fraction		FIOPOC	mL/dL	67.4
Inspired O2 / O2 ratio		OFIPOC	ratio	64.0
Sodium (POC)		NAPOC	mmol/L	22.5
Electrolytes POC	Potassium (POC)	KPOC	mmol/L	22.4
	Chloride (POC)	CLPOC	mmol/L	22.7
	Ionized Calcium (POC)	CAPOC	mmol/L	23.1
	Standard Ionized Calcium (POC)	CASPOC	mmol/L	23.2
	Anion gap	ANGPOC	mmol/L	19.6
	Glucose Blood Gas	GLUEMO	mg/dL	18.6
	Lactate (POC)	LATPOC	mmol/L	18.5

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Additional Information	Age	Age	Years	X	X	0
	Gender	Sex	Male/Female	X	X	0
	COVID-19 Suspect (Patient suffers from COVID-19 specific symptoms at triage)	Suspect	Yes / No	X	X	0
Target	COVID-19 positivity	Target	Positive/Negative	X	X	0

For Review Only

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

medRxiv preprint doi: <https://doi.org/10.1101/2020.10.02.20205070>; this version posted October 4, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

Table 2. Results for the models trained using the *OSR dataset*, the *COVID-specific dataset* and the *complete CBC dataset*. The first column reports the model used, the second column reports the accuracy, the third column reports the sensitivity, the fourth column reports the specificity, the fifth column reports the AUC, the sixth column reports the Brier score, and the seventh column reports the coverage for the 3-way version of the classifier with at least a 75% confidence score.

Dataset	Model	Accuracy	Sensitivity	Specificity	AUC	Brier <sup>1</sup>	Coverage for the 3-way version (75% confidence)
OSR dataset	Logistic Regression	0.86/0.92	0.88/0.95	0.84/0.90	0.86/0.95	0.13	0.76
	Naive Bayes	0.85/0.87	0.82/0.83	0.88/0.90	0.85/0.91	0.12	0.94
	KNN	0.83/0.89	0.76/0.82	0.90/0.95	0.83/0.90	0.12	0.72
	Random Forest	<b>0.88/0.93<sup>2</sup></b>	0.86/0.92	<b>0.91/0.94</b>	<b>0.90/0.94</b>	<b>0.10</b>	<b>0.70</b>
	SVM	<b>0.88/0.91</b>	<b>0.89/0.92</b>	0.87/0.90	0.88/0.94	0.11	0.77
COVID – specific dataset	Logistic Regression	0.83/0.87	<b>0.85/0.89</b>	0.82/0.85	0.83/0.88	0.14	0.70
	Naive Bayes	0.83/0.88	0.84/0.85	0.83/0.91	0.83/0.91	0.13	0.76
	KNN	<b>0.86/0.90</b>	0.80/0.85	<b>0.92/0.94</b>	<b>0.87/0.94</b>	<b>0.11</b>	<b>0.81</b>
	Random Forest	0.84/0.89	0.84/0.92	0.84/0.87	0.84/0.92	0.12	0.82
	SVM	<b>0.86/0.87</b>	0.83/0.83	0.89/0.91	0.86/0.93	0.12	0.74
CBC dataset	Logistic Regression	0.74/0.80	0.70/0.78	0.79/0.83	0.74/0.85	0.18	0.60
	Naive Bayes	0.78/0.83	0.74/0.79	0.82/0.87	0.78/0.88	0.16	0.69
	KNN	<b>0.86/0.90</b>	0.82/0.84	<b>0.89/0.95</b>	<b>0.86/0.89</b>	<b>0.13</b>	0.76
	Random Forest	0.83/0.90	<b>0.84/0.92</b>	0.82/0.87	<b>0.86/0.91</b>	<b>0.13</b>	0.68
	SVM	0.77/0.91	0.70/0.90	0.82/0.92	0.76/0.92	0.14	0.70

<sup>1</sup> Brier score, the lower it is, the better it is.

<sup>2</sup> The best value, for each score, is denoted in bold

medRxiv preprint doi: <https://doi.org/10.1101/2020.10.02.20205070>; this version posted October 4, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

**Table 3.** Results for the best models for the internal–external and external validation procedures. For each of the two features sets, we report the performances of the best models on the internal validation (namely, KNN for COVID-specific; and both Random Forest and KNN for CBC). The first four columns report on the results on the internal-external validation, while the last column reports on the results of the external validation.

Dataset	Accuracy	Sensitivity	Specificity	AUC	External Validation (specificity)
<i>COVID specific dataset</i> (KNN)	0.78	0.74	0.81	0.78	0.94
<i>CBC dataset</i> (RF)	0.76	0.70	0.82	0.76	0.96
<i>CBC dataset</i> (KNN)	0.75	0.72	0.78	0.75	0.92

For Review Only

# Development, evaluation, and validation of machine learning models for COVID-19 detection based on routine blood tests

Cabitza Federico<sup>1</sup>, Campagner Andrea<sup>2</sup>, Ferrari Davide<sup>3</sup>, Di Resta Chiara<sup>4</sup>, Ceriotti Daniele<sup>5</sup>, Sabetta Eleonora<sup>5</sup>, Colombini Alessandra<sup>2</sup>, De Vecchi Elena<sup>2</sup>, Banfi Giuseppe<sup>2</sup>, Locatelli Massimo<sup>5</sup>, Carobene Anna<sup>5</sup>

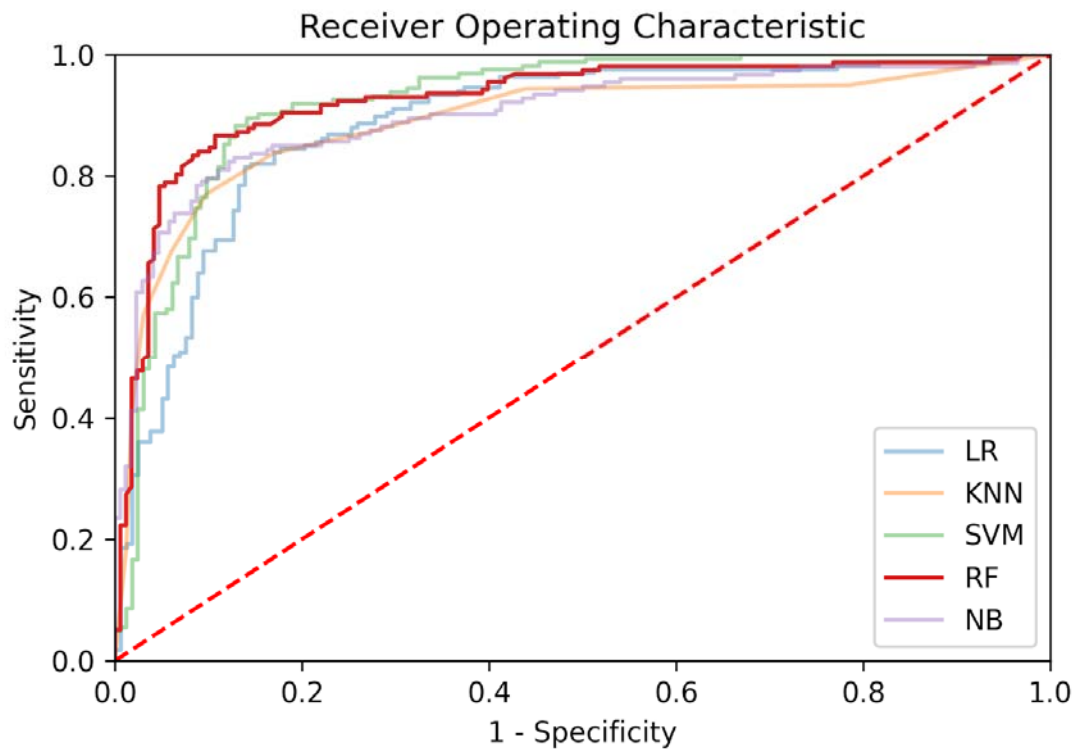
## SUPPLEMENTAL MATERIAL

**Suppl. Table 1.** Hyper-parameters of the machine learning models under consideration.

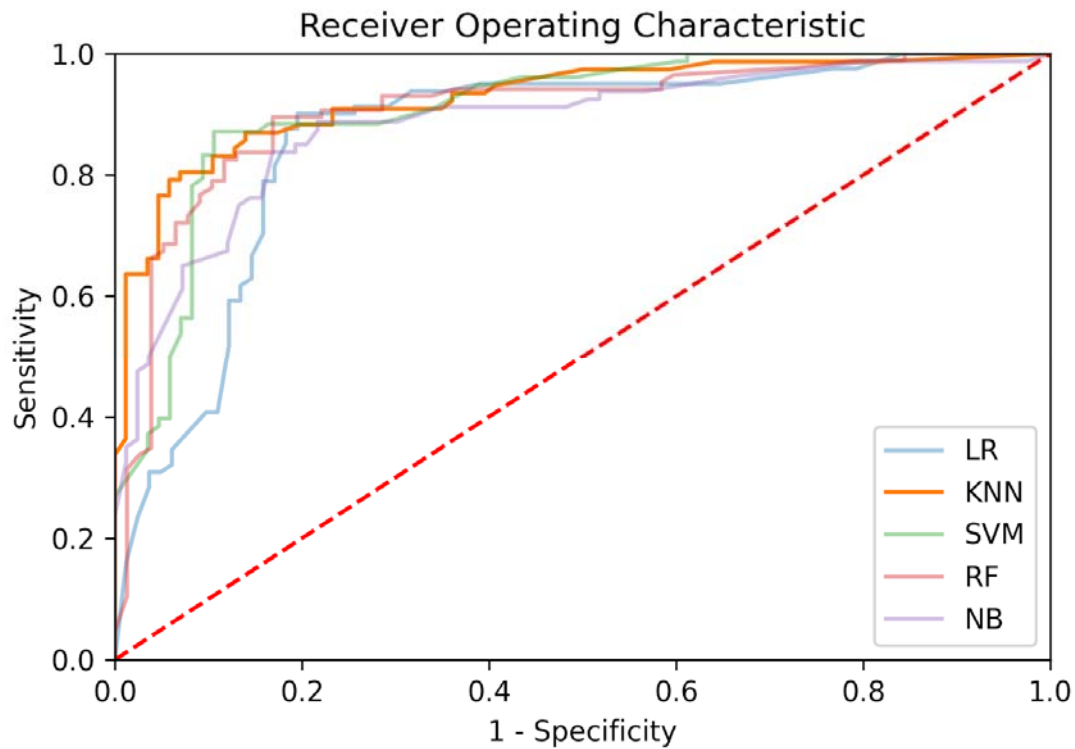
Model	Random Forest	Naive Bayes	Support Vector Machine	Logistic Regression	k-Nearest Neighbor
Hyper-parameters	Number of estimators, Maximum tree depth, Split criterion, Maximum number of features	/	Kernel, Maximum polynomial degree, Kernel coefficient, Regularization parameter	Regularization penalty, Regularization parameter	Distance weighting, Nearest Neighbor algorithm, Number of neighbors

# Supplementary Figures

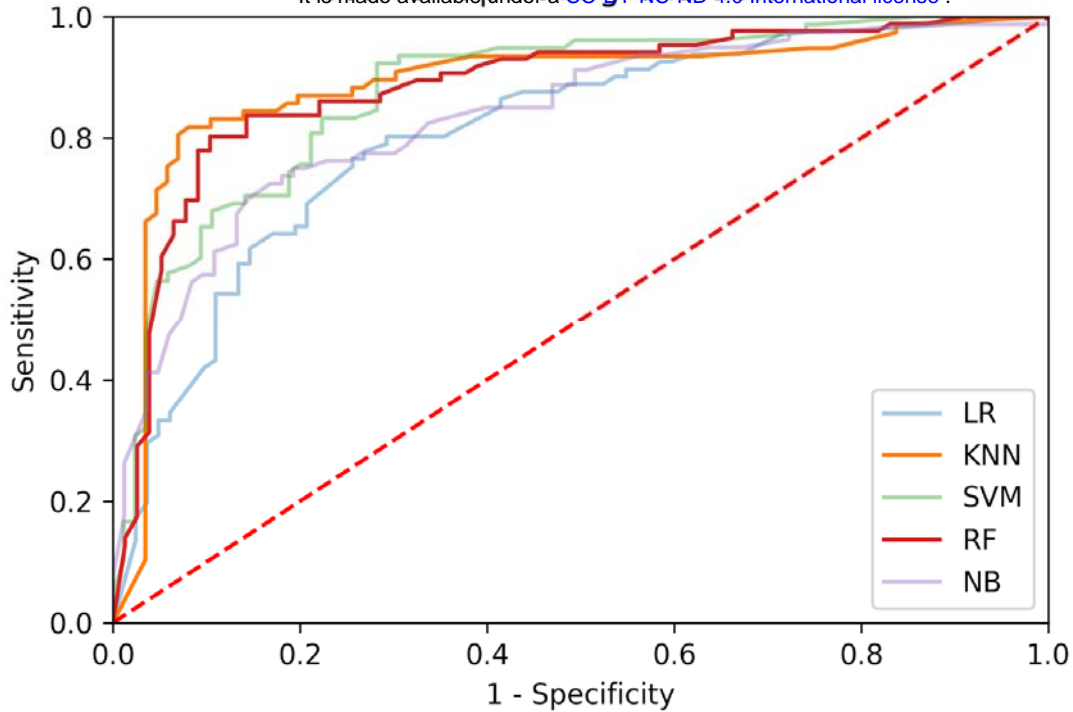
medRxiv preprint doi: <https://doi.org/10.1101/2020.10.02.20205070>; this version posted October 4, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



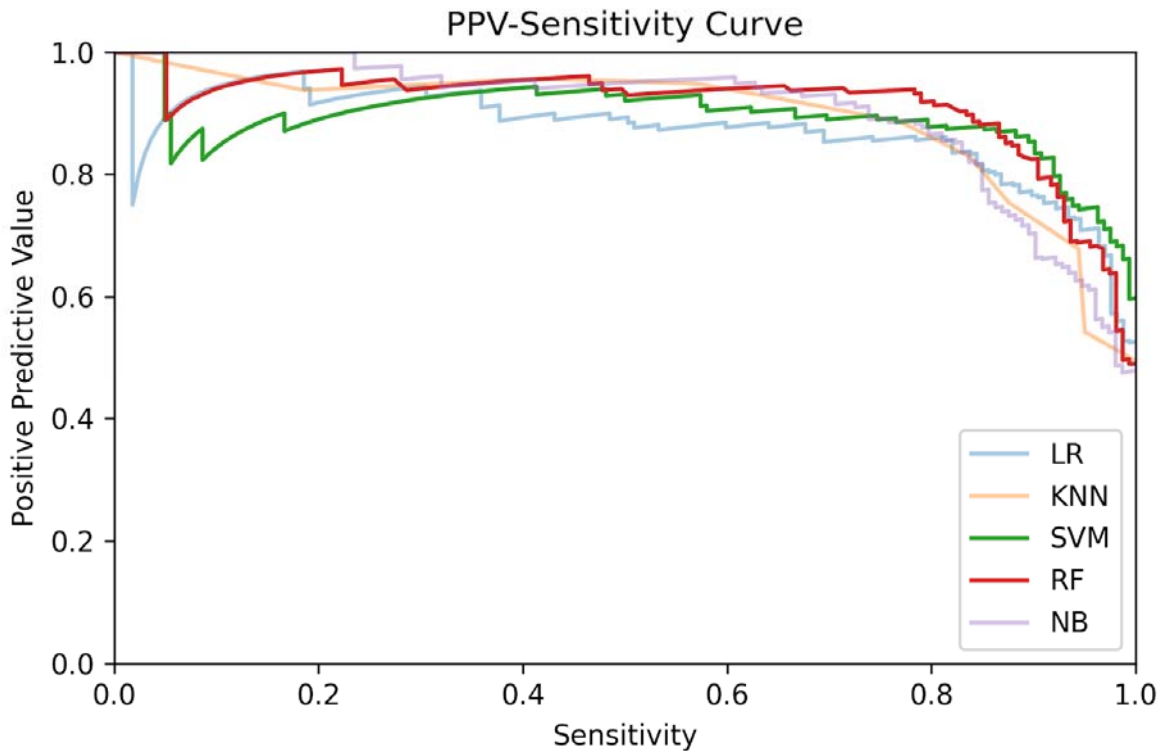
Suppl. Figure 1. Receiver operating characteristic curves for the models trained using the *OSR dataset*.



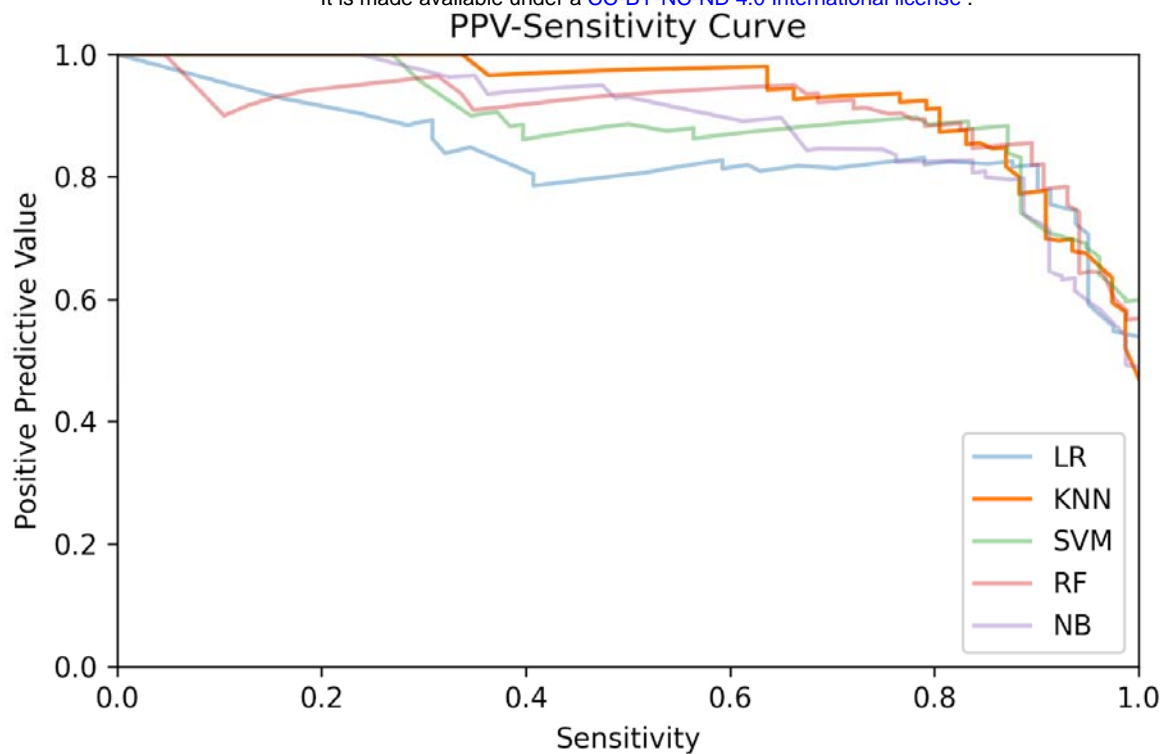
Suppl. Figure 2. Receiver operating characteristic curves for the models trained using the *COVID-specific dataset*.



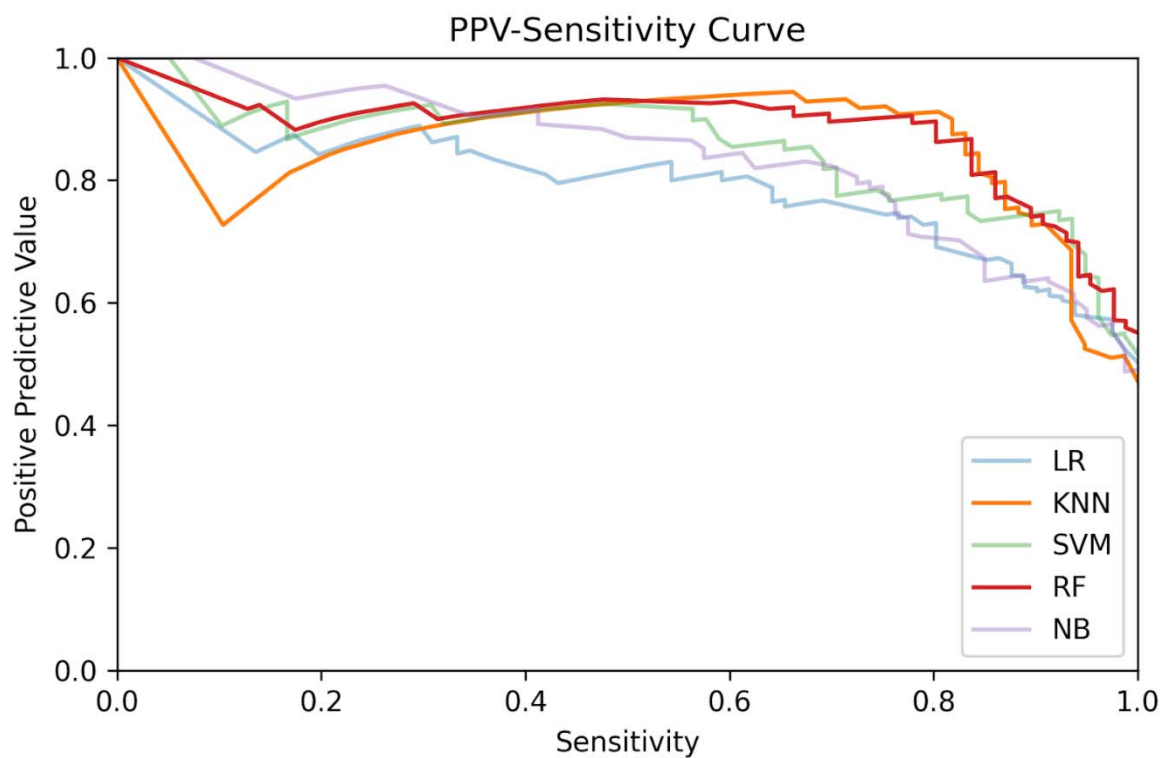
Suppl. Figure 3. Receiver operating characteristic curves for the models trained using the *CBC dataset*.



Suppl. Figure 4. Positive predictive value-sensitivity curves for the models trained using the *OSR dataset*.

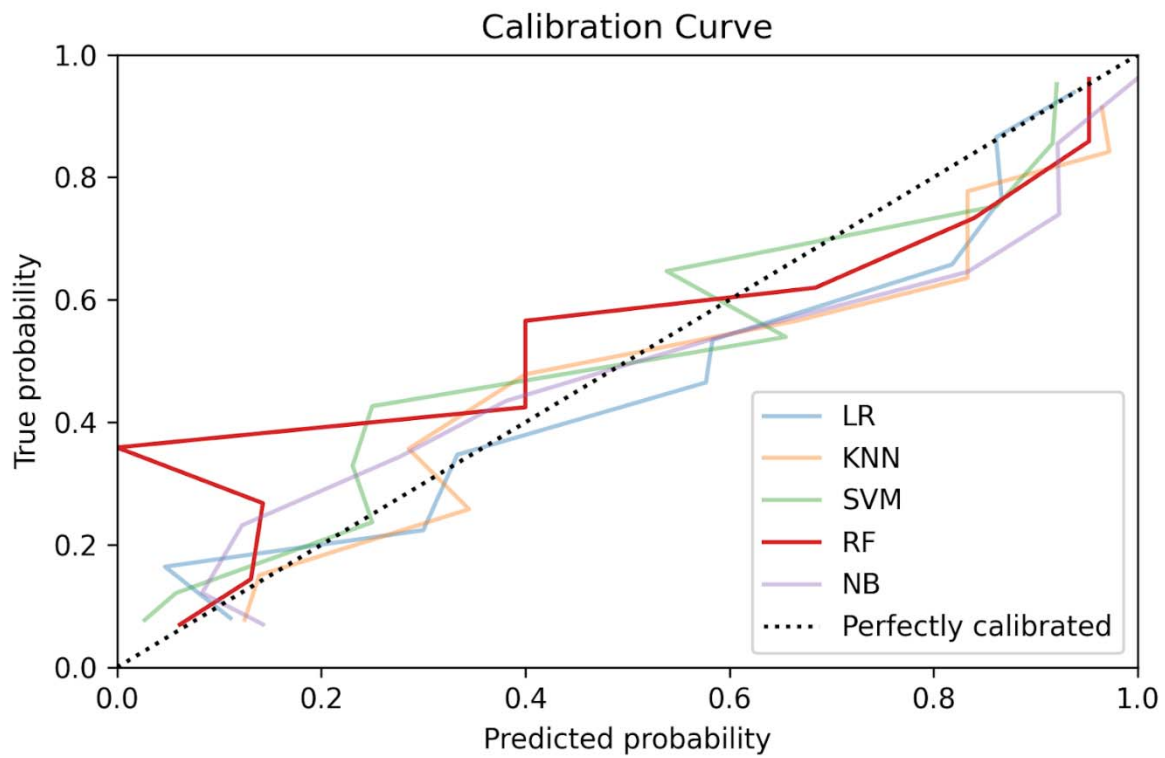


Suppl. Figure 5. Positive predictive value-sensitivity curves for the models trained using the *COVID-specific dataset*.

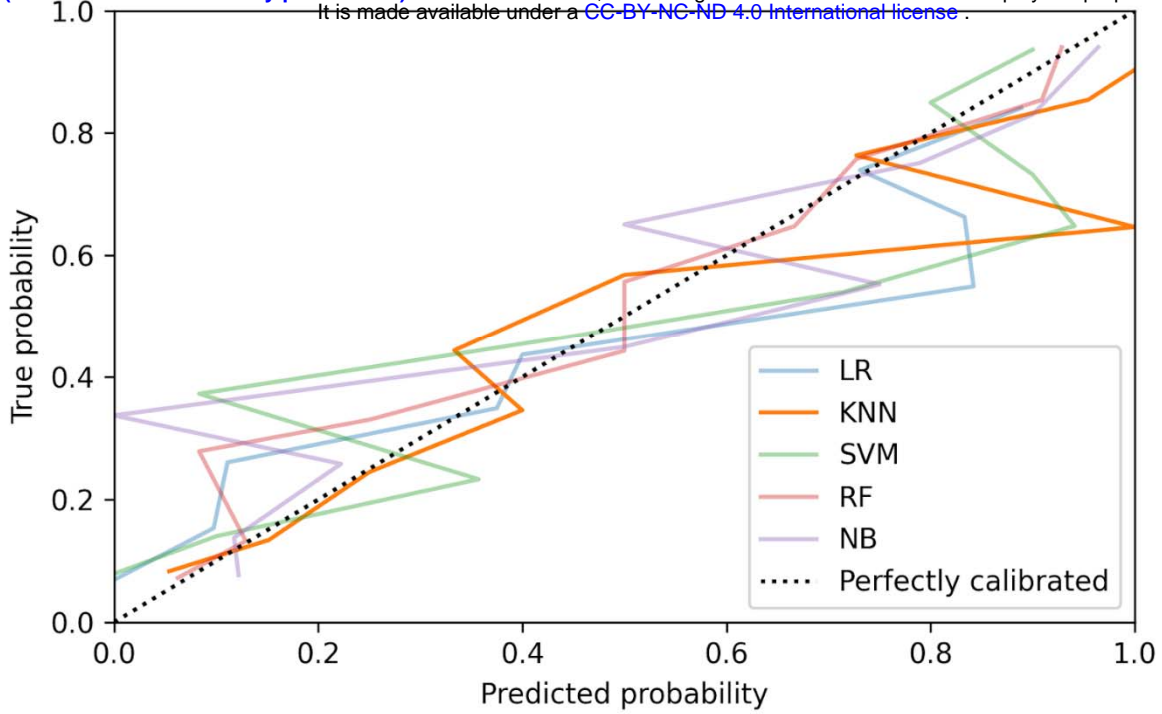


Suppl. Figure 6. Positive predictive value-sensitivity curves for the models trained using the *CBC dataset*.

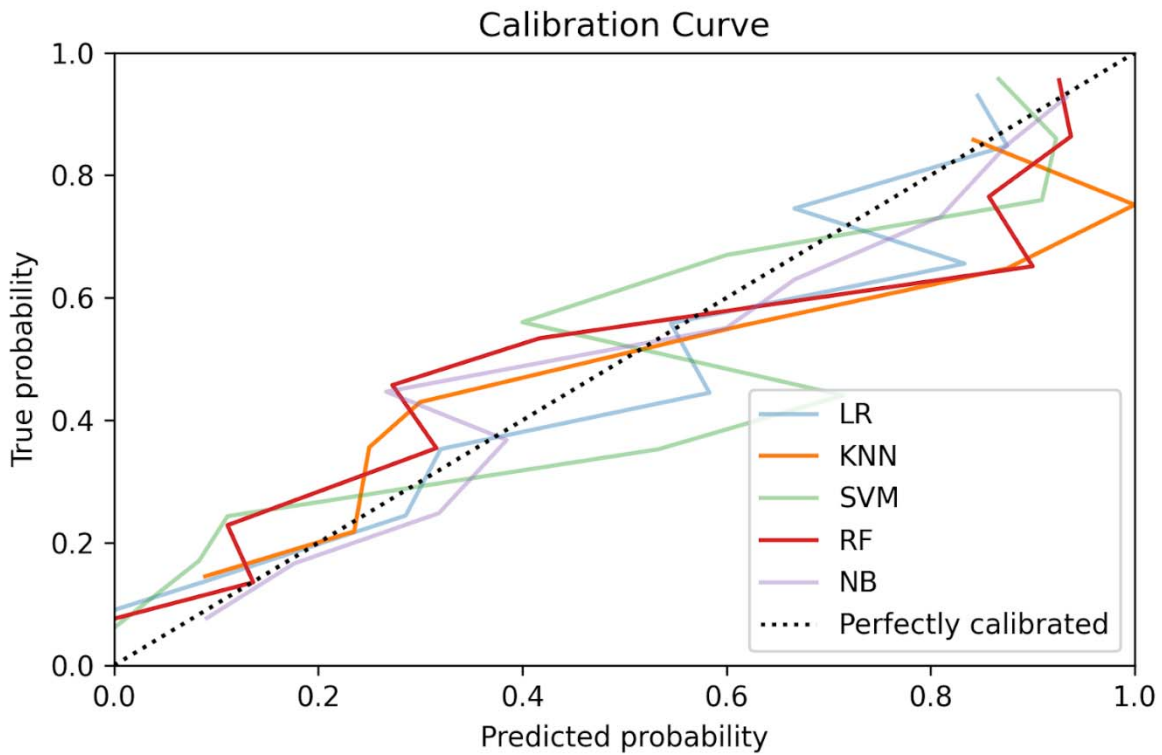




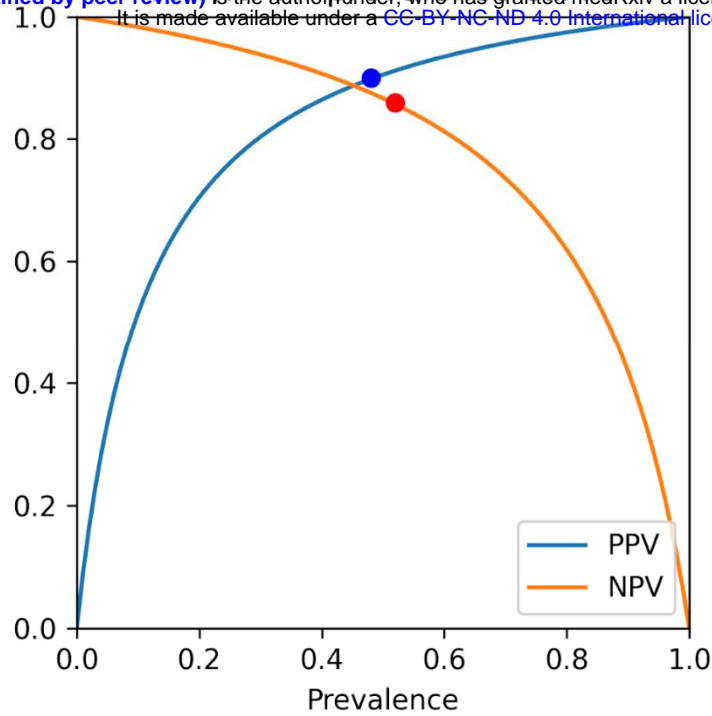
33  
34 Suppl. Figure 7. Calibration curves for the models trained using the *OSR dataset*.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



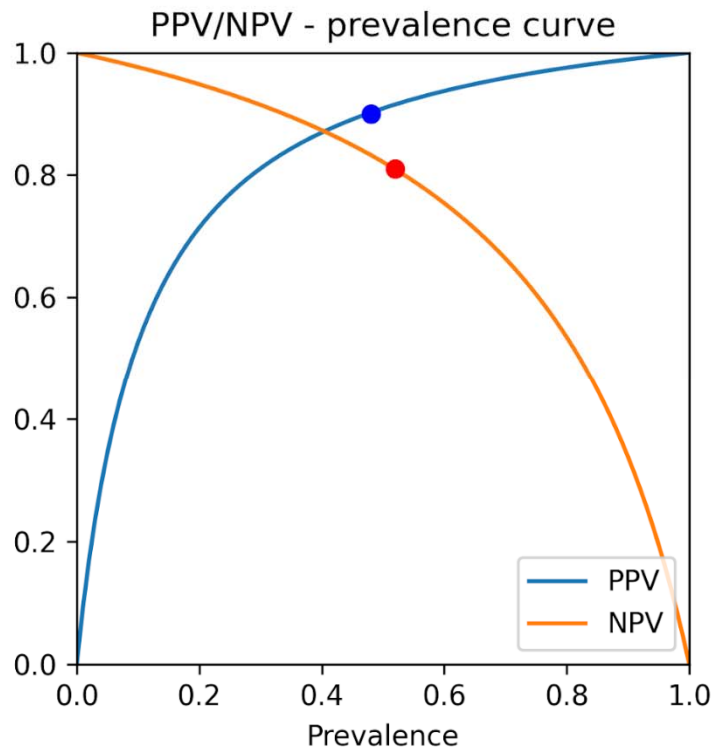
Suppl. Figure 8. Calibration curves for the models trained using the *COVID-specific dataset*.



Suppl. Figure 9. Calibration curves for the models trained using the *CBC dataset*.

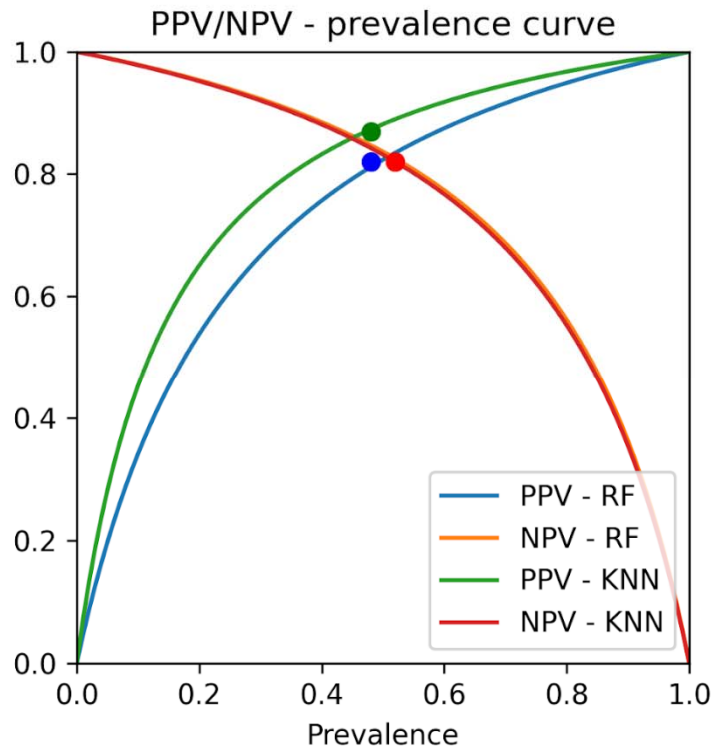


Suppl. Figure 10. Positive predictive value/negative predictive value-prevalence curve for the random forest algorithm, trained using the *OSR dataset*. The points on the curves indicate the prevalence in the dataset.

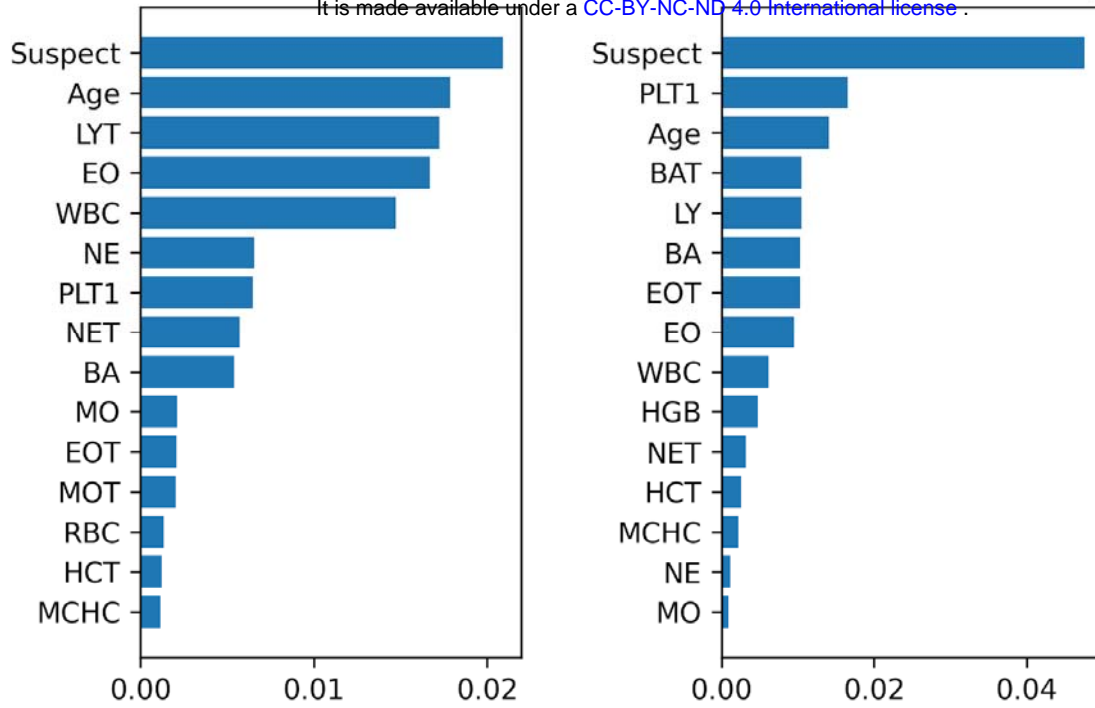


1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Suppl. Figure 11. Positive predictive value/negative predictive value-prevalence curve for the k-nearest neighbors algorithm trained using the *CBC dataset*. The points on the curves indicate the prevalence in the dataset. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).



Suppl. Figure 12. Positive predictive value/negative predictive value-prevalence curve for the random forest and k-nearest neighbors' algorithms, trained using the *CBC dataset*. The points on the curves indicate the prevalence in the dataset.



Suppl. Figure 13. Feature importances for Random Forest and k-Nearest Neighbors, trained on the *CBC* dataset.

## Identification of the Uncertain Cases

In the ground truthing process, we identified 165 uncertain cases for which we combined the results of the rRT-PCR test together with the radiologic gold standard. The uncertain cases were identified through two different methods: either patients who turned out to be positive within 72 hours after a first negative test and were admitted as inpatients despite this test result; or patients who, despite having had a negative test, had a hematochemical profile that was more similar to positive patients. For this purpose we used the k-means clustering algorithm ( $k = 2$ ) based on a set of COVID-19 characteristic biomarkers (AST, lymphocytes, calcium, LDH, CRP, WBC, XDP, fibrinogen)<sup>20,21</sup>.

## Implementation of the Internal-External Validation

The internal-external validation was performed based on the *IOG dataset*, using a bootstrap-based procedure. The goal of this procedure was to evaluate the ability of the developed models to generalize to new settings when provided with a limited amount of new data.

First, we generated 100 random, 50/50 train-test splits of the *IOG dataset*, then for each of these splits: first, the train set of the *IOG dataset* was oversampled using the SMOTE algorithm to obtain a sample of 1,624 synthetic instances; second, the oversampled train set was combined with the *COVID-specific* (respectively, *CBC*) dataset to obtain a combined training set encompassing 3,248 instances; third, the best models (obtained as described in the Methods and Results sections) were re-trained over the combined training set and evaluated on the test set. The average results over the 100 generated splits were reported.

## Hematochemical Analysis

The hematological analyses were performed on a Sysmex XE 2100 system (Sysmex, Japan) and the coagulation features were determined using the STAR Max analyzer (Stago Group, France); the biochemical parameters were measured on a Roche COBAS 6000 system (Roche Diagnostic, Basel, Switzerland) using Roche reagents, calibrators (Calibrator for automated systems [Cfas]/Cfas proteins), and control materials at two different levels (Precicontrol ClinChem Multi 1 and 2). All of the methods for the enzyme activity measurements were standardized to IFCC reference measurement procedures. The point of care (POC) measurements and the hemogas analysis were undertaken using Rapidpoint 500 (Siemens Healthcare).