



Since January 2020 Elsevier has created a COVID-19 resource centre with free information in English and Mandarin on the novel coronavirus COVID-19. The COVID-19 resource centre is hosted on Elsevier Connect, the company's public news and information website.

Elsevier hereby grants permission to make all its COVID-19-related research that is available on the COVID-19 resource centre - including this research content - immediately available in PubMed Central and other publicly funded repositories, such as the WHO COVID database with rights for unrestricted research re-use and analyses in any form or by any means with acknowledgement of the original source. These permissions are granted for free by Elsevier for as long as the COVID-19 resource centre remains active.



Diagnosis and prediction of COVID-19 severity: can biochemical tests and machine learning be used as prognostic indicators?

Alexandre de Fátima Cobre^a, Dile Pontarolo Stremel^b, Guilhermina Rodrigues Noletto^c, Mariana Millan Fachi^a, Monica Surek^a, Astrid Wiens^d, Fernanda Stumpf Tonin^a, Roberto Pontarolo^{d,*}

^a Pharmaceutical Sciences Postgraduate Programme, Universidade Federal Do Paraná, Curitiba, Brazil

^b Department of Forest Engineering and Technology, Universidade Federal Do Paraná, Curitiba, Brazil

^c Department of Biochemistry, Universidade Federal Do Paraná, Curitiba, Brazil

^d Department of Pharmacy, Universidade Federal Do Paraná, Curitiba, Brazil

ARTICLE INFO

Keywords:

COVID-19
Diagnosis
Severity
Blood test
Urine test
Machine learning model

ABSTRACT

Objective: This study aimed to implement and evaluate machine learning based-models to predict COVID-19 diagnosis and disease severity.

Methods: COVID-19 test samples (positive or negative results) from patients who attended a single hospital were evaluated. Patients diagnosed with COVID-19 were categorised according to the severity of the disease. Data were submitted to exploratory analysis (principal component analysis, PCA) to detect outlier samples, recognise patterns, and identify important variables. Based on patients' laboratory tests results, machine learning models were implemented to predict disease positivity and severity. Artificial neural networks (ANN), decision trees (DT), partial least squares discriminant analysis (PLS-DA), and K nearest neighbour algorithm (KNN) models were used. The four models were validated based on the accuracy (area under the ROC curve).

Results: The first subset of data had 5,643 patient samples (5,086 negatives and 557 positives for COVID-19). The second subset included 557 COVID-19 positive patients. The ANN, DT, PLS-DA, and KNN models allowed the classification of negative and positive samples with >84% accuracy. It was also possible to classify patients with severe and non-severe disease with an accuracy >86%. The following were associated with the prediction of COVID-19 diagnosis and severity: hyperferritinaemia, hypocalcaemia, pulmonary hypoxia, hypoxemia, metabolic and respiratory acidosis, low urinary pH, and high levels of lactate dehydrogenase.

Conclusion: Our analysis shows that all the models could assist in the diagnosis and prediction of COVID-19 severity.

1. Introduction

Coronavirus disease (COVID-19) remains an emergency of global interest; up to 21 May 2021, a total of 164.52 million confirmed cases and 3.42 million deaths had accumulated from the disease [1]. Social disparity and the scarcity of hospital resources for the treatment of patients in hospital units have been identified among the main factors associated with an increased number of deaths from this disease [2–7]. Thus, it is essential to identify potential prognostic biomarkers towards earlier and more targeted care, especially considering that some patients

with COVID-19 develop severe disease, which is associated with a higher risk of hospitalisation. Biomarkers provide a dynamic and powerful approach to understanding the spectrum of disease with applications in observational and analytic epidemiology, randomised clinical trials, screening and diagnosis, and prognosis [8]. Recently, studies investigating biomarkers to diagnose COVID-19 in early stages have been encouraged worldwide, aiming to provide a faster referral to treatment and reducing health-related problems associated with the disease [17, 18].

Machine learning (ML) is an effective and innovative tool able to

* Corresponding author.

E-mail addresses: alexandrecofre@gmail.com (A.F. Cobre), dile.stremel@gmail.com (D.P. Stremel), guilherminanoletto@ufpr.br (G.R. Noletto), marianamfachi@gmail.com (M.M. Fachi), monicasurek13@gmail.com (M. Surek), astridwiens@hotmail.com (A. Wiens), stumpf.tonin@ufpr.br (F.S. Tonin), pontarolo@ufpr.br (R. Pontarolo).

<https://doi.org/10.1016/j.combiomed.2021.104531>

Received 9 March 2021; Received in revised form 21 May 2021; Accepted 25 May 2021

Available online 29 May 2021

0010-4825/© 2021 Elsevier Ltd. All rights reserved.

assist healthcare professionals, policymakers, and other stakeholders during decision-making processes. These computer system models learn and adapt information by using algorithms and statistical networks to analyse and draw inferences from patterns in data. In the field of clinical diagnosis of COVID-19, these predictive analyses grounded on biomarkers can help optimise the screening of patients with severe disease, minimising mortality and hospitalisation, and reducing care delays. Previous machine learning studies highlight that some demographic variables, patients' comorbidities, and laboratory findings can be predictive factors for COVID-19 mortality [9–12]. However, most of these studies included a small sample size, which may impact the model's robustness and reliability of findings (e.g., low sensitivity) [14–16] and prevent its use in practice. Moreover, to date, specific biomarkers associated with the disease severity and patients' hospitalisation in intensive care units are still unknown, which may hamper the development of further targeted treatments [13].

Thus, this study aims to implement and evaluate machine learning-based algorithm models for the diagnosis and prediction of the severity of COVID-19 using data from biochemical, haematological, and urine tests from a large sample size.

2. Material and methods

2.1. Data set

Data from the public Kaggle platform [19] on individuals that had an reverse transcription polymerase chain reaction (RT-PCR) exam to detect severe acute respiratory syndrome coronavirus 2 (SARS-Cov-2) infection in the Israelita Albert Einstein hospital (São Paulo, Brazil) were collected. Regardless of the RT-PCR result (positive or negative for COVID-19), patients were included for analysis when presenting data on biochemical, haematological, and urinary parameters (Table 1). Two subgroups of samples were created according to RT-PCR results: (i) 5,643 patients' samples accounting for both negative ($n = 5,086$) and positive ($n = 557$) results; (ii) 557 positive samples of asymptomatic outpatients and patients with severe COVID-19, hospitalised in intensive care units. As the SARS-Cov-2 infection resembles other respiratory

Table 1

Biochemical, urinalysis, haematological, virological, and bacteriological tests performed on the patients included in the study.

Biochemical tests
Glucose serum, urea, C-reactive protein, creatinine, potassium, sodium, alanine transaminase, aspartate transaminase, gamma-glutamyltransferase, total bilirubin, direct bilirubin, indirect bilirubin, alkaline phosphatase, ionised pH, blood, magnesium analysis, HCO ₃ (venous blood gas analysis), lactate dehydrogenase, creatine phosphokinase, ferritin, arterial lactic acid, lipase dosage, HCO ₃ (arterial blood gas analysis), phosphorus, pCO ₂ (venous blood gas analysis), Hb saturation (venous blood gas analysis), base excess (venous blood gas analysis), pO ₂ (venous blood gas analysis), total CO ₂ (venous blood gas analysis), Hb saturation (arterial blood gases), pCO ₂ (arterial blood gas analysis), base excess (arterial blood gas analysis), pH (arterial blood gas analysis), total CO ₂ (arterial blood gas analysis), pO ₂ (arterial blood gas analysis), arterial FiO ₂ , and ctO ₂ (arterial blood gas analysis).
Haematological tests
Hematocrit, Hemoglobin, Platelets, Mean platelet volume, Red blood Cells, Lymphocytes, Mean corpuscular hemoglobin concentration, Leukocytes, Basophils, Mean corpuscular hemoglobin, Eosinophils, Mean corpuscular volume, Monocytes, Red blood cell distribution width
Urine tests
Urine pH, segmented neutrophil, promyelocytes, metamyelocytes and myeloblasts, and the international normalised ratio (INR).
Virological tests
Respiratory syncytial virus, influenza A, influenza B, parainfluenza 1, coronavirus NL63, rhinovirus/enterovirus, coronavirus HKU1, parainfluenza 3, adenovirus, parainfluenza 4, coronavirus 229E, coronavirus OC43, influenza A H1N1, influenza H1N1 test, and influenza A rapid test.
Bacteriological tests
<i>Mycoplasma pneumoniae</i> , <i>Chlamydomphila pneumoniae</i> , and <i>Streptococcus A</i> .

diseases, to minimise the chance of obtaining false-positive samples, patients who tested positive for at least one other virus or respiratory bacteria (see Table 1) were excluded from analyses.

2.2. Machine learning models

The first step for implementing any ML model is to perform exploratory analysis. The exploratory analysis intends to: (i) identify the presence of possible outliers, (ii) recognise patterns of data distribution in the multidimensional space, and (iii) identify relationships between variables [20]. In this study, both data used to build the COVID-19 diagnostic model and the data used to build the severity prediction model were previously subjected to two methods of exploratory analysis: principal component analysis (PCA) and k-means cluster analysis (KMCA). Additionally, the outliers were detected and eliminated from the dataset using the graphical method of leverage versus student residuals [21].

A total of four algorithms were used: (i) artificial neural networks (ANN), (ii) decision trees (DT), (iii) discriminant analysis by partial least squares (PLS-DA), and (iv) the method of k-nearest neighbours (KNN). For implementing these models with algorithms for the diagnosis and prediction of COVID-19 severity, 70% of the samples were used for the training set and 30% for the test set. For both the diagnostic model and the severity model, the Kennard-Stone method was employed to select samples from the training set and samples from the test set [22]. The samples used for implementing the algorithms for COVID-19 diagnosis were divided into class 1 (negative samples for COVID-19) and class 2 (positive samples for COVID-19). For the severity prediction models, samples were classified into class 1 (non-severe disease; i.e., outpatients) and class 2 (severe disease; i.e., hospitalised patients).

The number of latent variables (LVs) selected for the ML models was performed using the leave-one-out cross-validation method. The number of LVs presented the lowest square root of mean cross-validation error (RMSECV). The predictive capacity of the model was evaluated using the square root of mean error of prediction (RMSEP), where the classification models were optimised considering the lower RMSEP.

The analytical validation of the models based on the machine learning algorithms was performed using the following metrics: sensitivity, specificity, and accuracy. These figures of merit were calculated using the parameters true positive (TP), true negative (TN), false positive (FP), and false-negative (FN) [20,23–25]. In ML, a sample is called a true positive when it belongs to class one (1) and is correctly classified by the ML algorithm as belonging to class one (1). A sample is considered a false positive when it belongs to class zero (0) and is incorrectly classified by the ML algorithm as being class one (1). A true negative sample belongs to class zero (0) and is correctly classified as class zero (0). Finally, a sample is false negative when it belongs to class one (1) and is wrongly classified by the ML algorithm as class zero (0). Sensitivity and specificity are defined as the ability of the ML-based model to correctly classify negative and positive samples, respectively. Accuracy is the ability of an ML-based model to correctly classify both negative and positive samples. The values of sensitivity, specificity, and accuracy vary from zero (0) to one (1), and the closer to 1, the more sensitive, specific, and accurate the model, respectively. These parameters were calculated according to equations (1)–(3), respectively:

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (1)$$

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (2)$$

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (3)$$

where FN: false negative; FP: false positive; TN: true negative; and TP: true positive.

The accuracy of the models was calculated by integrating the area under the receiver operating characteristic (ROC) curve.

3. Results

The variation of all biomarkers (biochemical, haematological, and urinary) for the COVID-19 patients are summarised in Table 2.

The results of the exploratory analyses of the diagnostic and severity disease data using the PCA model are depicted in Fig. 1. Additional results with the KMCA model are available in Figures ESM 1 and ESM 2 in the Online Resource. These models were able to differentiate between positive and negative patients for COVID-19 (diagnostic data) and between patients with non-severe disease and vs severe disease (disease severity data).

Outlier samples were analysed for both diagnostic and disease severity data using the graph of leverage versus student residuals with 95% of confidence interval (Fig. 2). No sample was detected as an outlier.

The results of the four ML-based models for the data subsets are presented in Table 3. All models for COVID-19 diagnosis and prediction of disease severity were compared with each other using the following metrics: training time, model training error, cross-validation error, sensibility, specificity, and accuracy (area under the ROC curve). The ANN model performed better as it presented less training time and forecasting errors, and greater overall accuracy. The ROC curves of the models are depicted in Fig. 3.

According to the ML-based models, some biomarkers were judged as critical or important for predicting COVID-19 and disease severity (see Table 4). Ferritin was ranked as the most important variable in all models. Further information is available on Online Resource.

4. Discussion

In the present study, we used four ML-based models (ANN, DT, PLS-DA, and KNN) with over 5,000 RT-PCR samples and data on patients' biochemical, haematological, and urinary parameters that effectively predicted COVID-19 diagnosis and disease severity in Brazil. Considering the high likelihood of existing SARS-Cov-2 mutations (B.1.1.7, P.1, and P.2) in the sample, the complexity of these models is even higher [38–40].

Several ML-based models, including unsupervised approaches (e.g., PCA and hierarchical cluster analysis) and supervised models (e.g., artificial neural network, PLS-DA, DT, KMC, and KNN), are available in the scientific literature [26,27]. The performance of these models relies on several factors, including sample size and the type of data. ML-based models built with larger samples are usually more accurate and efficient for series forecasting; for instance, the deep neural networks that require a great amount of training data [28–30]. The larger the network architecture, the more data is needed to obtain more robust models [31–33].

Regarding COVID-19, previous models have been implemented aiming at predicting disease behaviour and severity. However, most of these studies used a small sample size, which may directly impact the performance of the model [34–36]. Banerjee (2020) developed an ML algorithm to forecast COVID-19 diagnosis using a public database with 598 patients, of which only 39 were positive for SARS-CoV-2. The authors obtained a model with good specificity (91%) but low sensitivity (43%), which can prevent the use of the model in practice for early diagnosis of the disease [14]. Similarly, Joshi (2020) implemented a logistic regression model previously trained with 390 samples, of which only 33 were positive for COVID-19, proving sensitivity and specificity values of 93% and 43%, respectively [37]. Additionally, most studies implemented ML-based models using only routine blood tests [41,42]. In our study, besides the complete blood count test, data from patients' biochemical, urinary, bacteriological, and virological tests aiming at identifying further biomarkers associated with COVID-19 were also included.

Table 2

Levels of biochemical, haematological, and urine biomarkers variation in positive patients with severe disease on a normalised scale of patients.

Biomarker	COVID-19 positive patients' samples ^a	COVID-19 severe patients' samples ^b
Hematocrit	Low	Low
Haemoglobin	Low	Low
Platelets	Low	Low
Mean platelet volume	Low	Low
Red blood Cells	Low	Low
Lymphocytes	Low	Low
Mean corpuscular haemoglobin concentration (MCHC)	Low	Low
Leukocytes	High	High
Basophils	Normal	Normal
Mean corpuscular haemoglobin (MCH)	Normal	Low
Eosinophils	Low	Low
Mean corpuscular volume (MCV)	Low	Low
Monocytes	High	Normal
Red blood cell distribution width (RDW)	Low	Normal
Serum glucose	High	High
Neutrophils	Low	Low
Urea	Low	Low
C-reactive protein	High	High
Creatinine	High	High
Potassium	Low	Low
Sodium	Low	Low
Alanine transaminase	High	High
Aspartate transaminase	High	High
Gamma-glutamyltransferase	High	High
Total bilirubin	High	High
Direct bilirubin	High	High
Indirect bilirubin	High	High
Alkaline phosphatase	High	High
Ionised calcium	Low	Low
pCO ₂ (venous blood gas analysis)	High	High
Magnesium	Low	Low
Hb saturation (venous blood gas analysis)	Low	Low
Base excess (venous blood gas analysis)	Low	Low
pO ₂ (venous blood gas analysis)	Low	Low
Total CO ₂ (venous blood gas analysis)	High	High
pH (venous blood gas analysis)	Low	Low
HCO ₃ (venous blood gas analysis)	High	High
Rods	High	High
Segmented	Low	Low
Promyelocytes	Normal	–
Metamyelocytes	Normal	–
Myelocytes	Normal	–
Urine pH	Low	Low
Urine density	Normal	Low
Urine red blood cells	Normal	Normal
International normalised ratio (INR)	High	High
Lactate dehydrogenase	High	High
Creatine phosphokinase (CPK)	Normal	Low
Ferritin	High	High
Arterial lactic acid	High	High
Hb saturation (arterial blood gases)	Low	Low
pCO ₂ (arterial blood gas analysis)	High	High
Base excess (arterial blood gas analysis)	Low	Low
pH (arterial blood gas analysis)	Low	Low
Total CO ₂ (arterial blood gas analysis)	High	High
HCO ₃ (arterial blood gas analysis)	High	High
pO ₂ (arterial blood gas analysis)	Low	Low
Arterial FiO ₂	Low	Low
Phosphorous	Low	–

^a Diagnostic data.

^b Disease severity data.

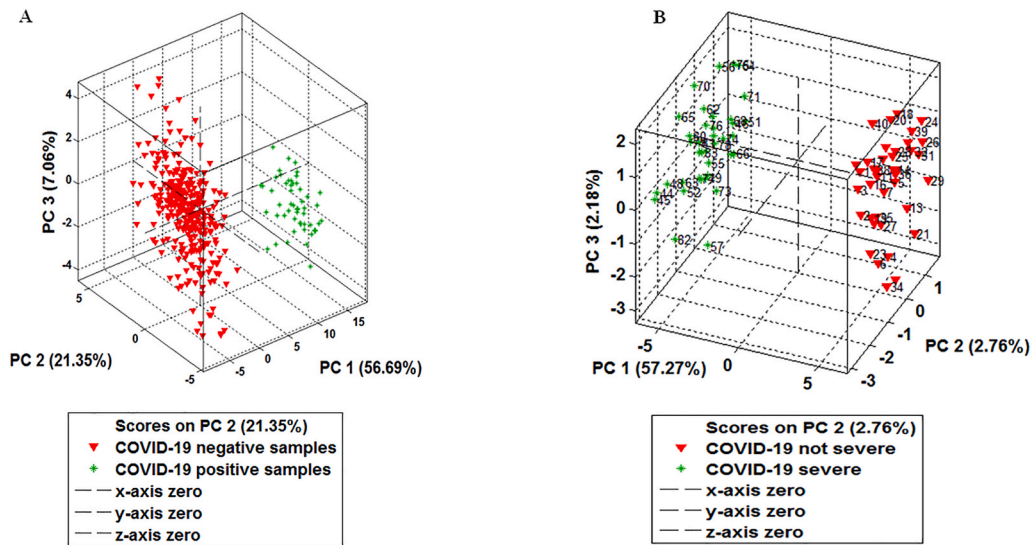


Fig. 1. Exploratory analysis. Principal component analysis (PCA) model for the discrimination of negative and positive samples (A) and samples from patients with severe and non-severe disease (B).

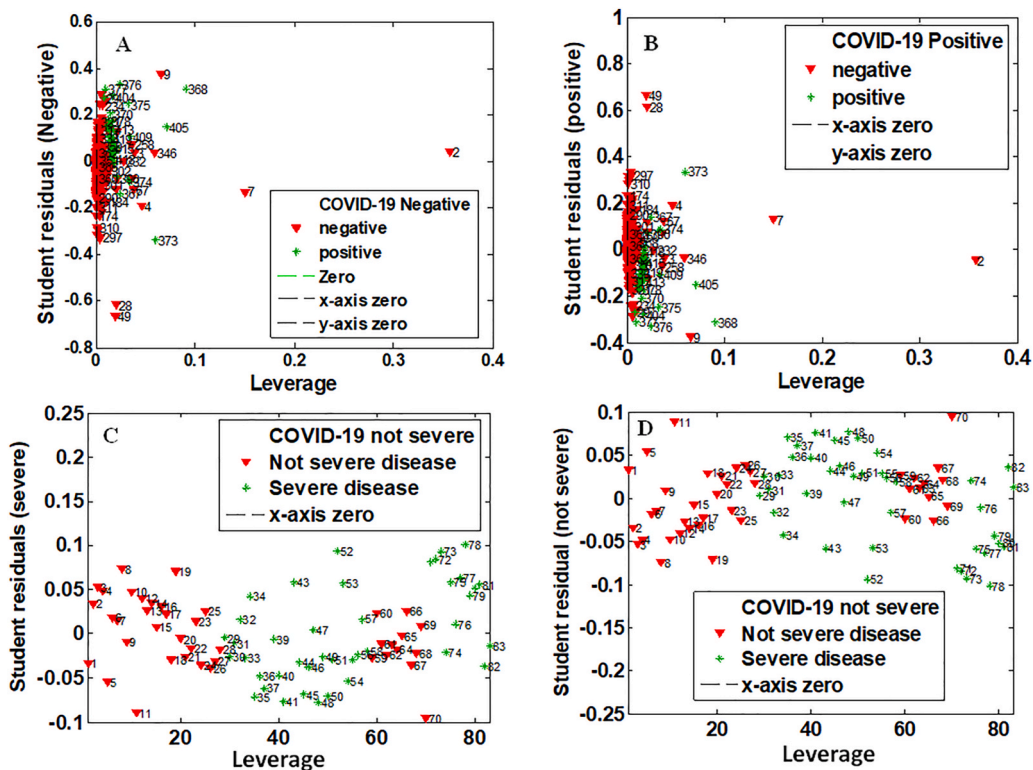


Fig. 2. Graph of leverage versus student residuals for detecting outlier samples. For diagnostic data: outlier analysis of negative samples (A) and positive samples (B). For severity data: outlier analysis for samples from patients without severity (C) and with severity (D).

The models reached 84%–98% accuracy. The biomarkers that most contributed to this result and predicting the diagnosis and severity of COVID-19 included: hyperferritinaemia, hypocalcemia (low levels of ionised calcium), hypoxemia (low arterial oxygen pressure, pO_2), pulmonary hypoxia (low inspired fraction of arterial oxygen, FiO_2), respiratory acidosis (high levels of total CO_2 and pCO_2), metabolic acidosis (high levels of lactic acid and low venous pH), low urinary pH, and high levels of lactate dehydrogenase (LDH).

Similar precision values were recently reported by Zhou (2021) (94%) and Wu (2021) (90%) after the implementation of ML-based

models for COVID-19 diagnosis and disease severity, respectively [15, 35]. However, different variables were highlighted by the authors as important for data classification. According to Zhou (2021), these were the rates of circulating lymphocytes, while Wu (2021) reported the rates of neutrophils and lymphocytes, the neutrophil/lymphocyte ratio, and platelet/lymphocyte ratio [15,35]. This variation may be associated with the different sample sizes ($n = 357$ vs $n = 51$) and type of models (decision tree model vs support vector machine model).

Recently, most patients with severe COVID-19 who require hospitalisation in intensive care units develop an atypical form of acute

Table 3
Performance comparison of the machine learning models for COVID-19.

Metric	Diagnostic model				Disease severity model			
	RNA	DT	PLS-DA	K-NN	RNA	DT	PLS-DA	K-NN
Training time	21 min. 43 s	27 min. 11 s	31 min. 19 s	22 min. 15 s	7 min. 1 s	10 min. 19 s	18 min. 3 s	09 min. 53 s
Calibration error	1.0%	0.5%	1.2%	0.5%	1.0%	8.4%	6.0%	0.4%
Cross validation error	0.8%	1.0%	0.9%	0.6%	0.5%	1.8%	4.0%	0.7%
Sensitivity	0.93	0.89	0.88	0.84	0.99	0.90	0.87	0.82
Specificity	0.94	0.89	0.90	0.83	0.97	0.94	0.88	0.88
Accuracy ^a	0.94	0.90	0.90	0.84	0.98	0.92	0.88	0.86

^a Area under the ROC curve.

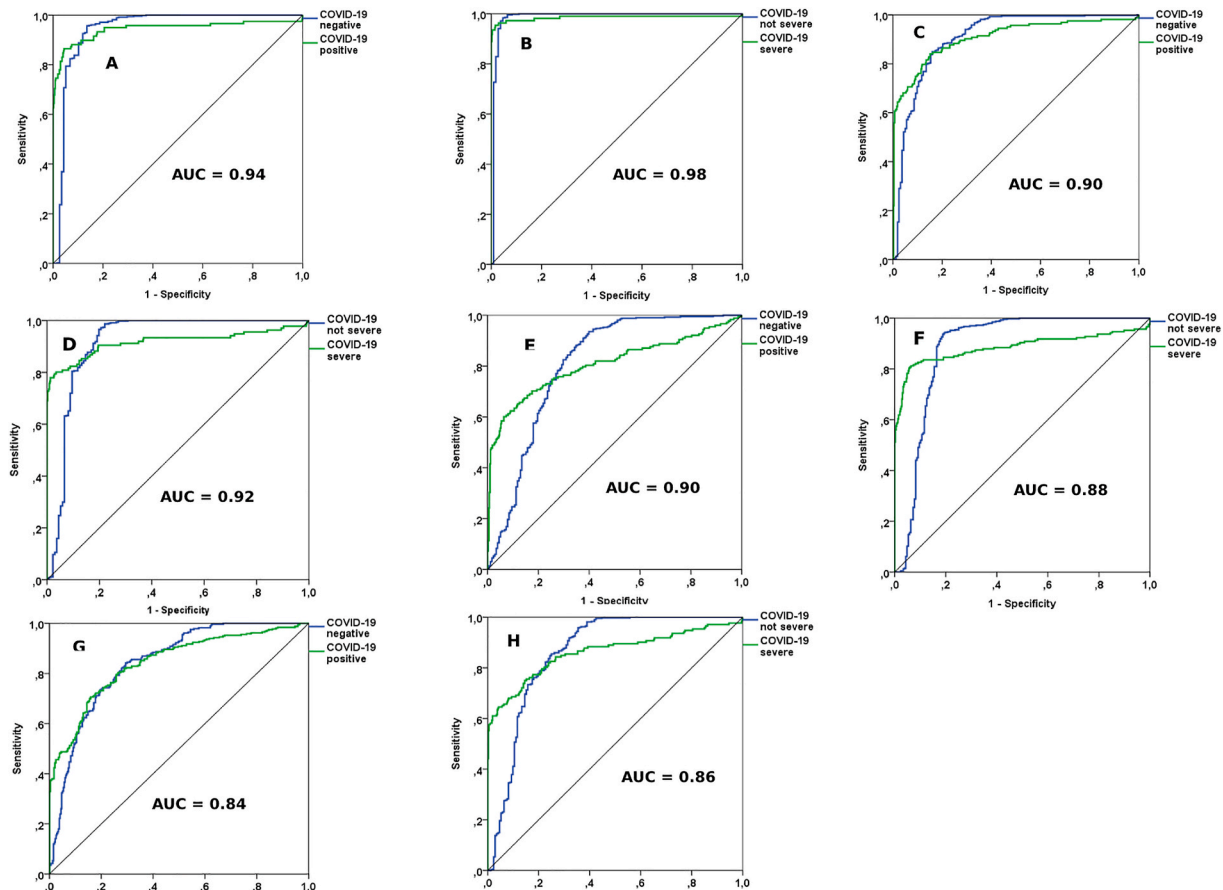


Fig. 3. ROC curves of the accuracy of the machine learning models. Artificial neural network (ANN): diagnosis (A) and severity (B). Decision tree (DT): diagnosis (C) and severity (D). Discriminant analysis by partial least squares (PLS-DA): diagnosis (E) and severity (F). K-nearest neighbours (KNN): diagnosis (G) and severity (H).

distress syndrome, which is usually accompanied by a preserved volume of pulmonary gas [43]. This suggests hypoxia, which results from the difficulty in performing gas exchanges at the level of the pulmonary alveoli. It has been observed that this pulmonary dysfunction can also compromise iron metabolism [43]. Imbalances in both haemoglobin and ferritin levels were reported in patients with severe disease or deaths caused by COVID-19. In a systematic review with a meta-analysis conducted by Taneri (2020), which included 189 studies ($n = 57,563$ patients), high-risk patients with severe disease had significantly high levels of ferritin (weighted mean difference (WMD), 473.25 ng/mL (95% CI 382.52; 563.98)) and low haemoglobin levels (WMD, 4.08 g/L (95% CI 5.12; -3.05)) when compared to patients with moderate or low-risk disease [44]. In our study, both patient groups with the disease (diagnostic model) and patients with severe disease (disease severity model) presented low levels of haemoglobin and high levels of ferritin. All developed ML-based models (ANN, PLS-DA, KNN, and DT)

highlighted that the differences in the levels of these two biomarkers were critical, both for the prediction of COVID-19 diagnosis and disease severity. Recent studies have indicated anaemia and hyperferritinaemia to be strong biomarkers for the prognosis of mortality due to SARS-CoV-2, in addition to other serious respiratory diseases [45–50].

Ferritin is a protein found mainly in the liver, bone marrow, and spleen and is the body's main source of iron storage. It is considered a key biomarker of immune dysregulation, mainly in a situation of hyperferritinaemia, through the direct route of proinflammatory and immunosuppressive effects contributing to a cytokine storm [51–53]. Some studies show that fatal outcomes caused by COVID-19 are accompanied by cytokine storm syndrome involving high levels of inflammatory markers, such as ferritin [54,55]. In the present study, ferritin was associated with predicting the severity of COVID-19, and it was the most important biomarker in predicting the diagnosis of the disease by all the ML-based models developed, corroborating the

Table 4
Important biomarkers in machine learning models for the diagnosis and classification of COVID-19 severity.

Biomarkers	Diagnostic model				Disease severity model				
	ANN	DT	PLSDA	KNN	ANN	DT	PLSDA	KNN	
Ferritin	++	++	++	++	++	++	++	++	
Gamma-glutamyltransferase	+	-	-	-	-	-	-	-	
HCO ₃ (arterial)	+	-	-	+	-	-	+	+	
Base excess (arterial)	+	-	-	-	-	-	+	+	
Base excess (venous)	-	-	-	-	-	-	+	-	
Sodium	+	-	-	-	-	-	-	-	
Total O ₂ (arterial)	-	-	-	-	+	-	-	-	
pO ₂ (arterial)	-	-	-	+	-	-	-	-	
Total CO ₂	-	+	+	+	-	-	-	+	
pCO ₂ (arterial)	+	-	-	-	-	-	+	+	
pCO ₂ (venous)	+	-	+	-	-	-	+	-	
Indirect bilirubin	+	-	-	-	-	-	-	-	
Alkaline phosphatase	+	-	-	-	+	-	-	-	
Urine pH	-	-	+	+	+	-	+	-	
pH (venous)	-	-	-	-	-	-	+	-	
pH (arterial)	-	-	-	-	-	-	-	+	
FiO ₂ (arterial)	-	-	-	-	+	-	-	-	
ctO ₂ (arterial)	-	-	-	-	-	-	-	+	
Total bilirubin	-	-	-	-	+	-	-	-	
Red blood cell distribution width	-	-	-	-	+	-	-	-	
Platelets	-	-	-	-	+	-	-	-	
C-reactive protein	-	-	-	-	+	-	+	-	
Calcium ionised	-	-	+	+	-	-	-	-	
Urine-density	-	-	+	+	-	-	-	-	
Lactate dehydrogenase	-	-	-	-	-	-	+	-	
Arterial lactic acid	-	-	-	+	-	-	-	+	
Haemoglobin saturation (arterial)	-	-	-	-	-	-	-	+	
Phosphorous	-	-	+	+	-	-	-	-	
Lipase dosage	-	-	-	-	-	-	+	-	
Rods	-	-	-	+	-	-	+	-	

Less important variable (-); Important variable (+); Critical variable (++)

literature data.

Acid and base disorders are important indicators in the pathogenesis and severity of several diseases, especially respiratory diseases of infectious origin, such as pneumonia [56–58]. Acidosis can occur as a result of a significant increase in arterial carbon dioxide pressure (respiratory acidosis) or a variety of inorganic or organic compounds (metabolic acidosis), such as bicarbonate, lactic acid, arterial, ketones, or as a result of renal failure or hyperchloremic acidosis; all of these factors act simultaneously in the increase of hydrogen protons and, consequently, the reduction of blood and respiratory pH levels [59–62]. Researchers suggest that the metabolic acidosis caused by lactic acid in COVID-19 is probably due to anaerobic glycolysis, which is favoured in consequence of hypoxemia. In this condition, pyruvate, a product of the glycolytic pathway, is not translocated to mitochondria to follow the oxidative process [63,64]; instead, it is converted into lactate in the cytosol by the LDH enzyme. As the hypoxemia impairs the tissue oxygenation and the oxidative phosphorylation, the cells obtain ATP by anaerobic glycolysis. This flow relies on the conversion of pyruvate to lactate which results in high levels of this metabolite that comes out of the cells. The excessive consumption of lactate during the process of gluconeogenesis culminates in lactic acidosis [65]. It has been reported that on the 18th day of COVID-19 disease, the levels of lactic acid begin to increase significantly, triggering metabolic acidosis, although the carbon dioxide pressure is acceptable [66].

In our study, patients positively diagnosed with COVID-19 (diagnostic model) and with severe disease (disease severity model) had high levels of carbon dioxide pressure (arterial and venous gas analysis), total carbon dioxide (arterial and venous gas analysis), arterial lactic acid and bicarbonate (arterial gas analysis), and exceptionally low venous pH (pH = 0.3), which also suggests respiratory and metabolic acidosis. These results are similar to other countries, indicating that these metabolic imbalances are prevalent in COVID-19 patients [64,67,68]. In general, all the available data show that most patients with severe disease have accompanying comorbidities, such as diabetes. Recent studies

indicate that metabolic acidosis is influenced by the use of metformin for the treatment of diabetes mellitus [69,70].

The present analysis detected levels of five function biomarkers, including bilirubin, direct bilirubin, indirect bilirubin, alanine transaminase, and aspartate transaminase, increased in patients with positive (diagnostic model) and severe (severity model) disease compared to patients with non-severe disease. A systematic review and meta-analysis performed by Parohan (2020) revealed similar findings, where all 1,455 patients with severe disease had extremely high levels of total bilirubin (WMD 2.30 mmol/l; 95% CI, 1.24; 3.36; $p < 0.001$), alanine aminotransferase (WMD 7.35 U/L; 95% CI, 4.77; 9.93; $p < 0.001$) and aspartate aminotransferase (WMD 8.84 U/L; 95% CI 5.97; 11.71; $p < 0.001$), compared to 1,973 patients with non-severe disease [71]. Liver damage has also been reported in other viral pneumonia (e.g., MERS and SARS) and is directly associated with disease severity and mortality [72–75]. However, the biochemical/pathophysiological mechanisms that explain the liver dysfunction caused by COVID-19 are still unknown. It is unclear if the liver dysfunction is due to SARS-Cov-2 or is a consequence of multiple organ failure caused by the virus [71].

Additionally, high levels of the C-reactive protein (CRP) in the samples were found. An increase in this important biomarker in COVID-19 patients' was previously reported by Chen (2020) (up to 86%) [50]. Recently, a systematic review with meta-analysis concluded that extremely high levels of CRP were statistically associated with COVID-19 severity [77,78]. CRP is an inflammatory protein in the acute phase of inflammatory and infectious processes that is synthesised mainly in liver cells but also smooth muscle cells, macrophages, lymphocytes, and adipocytes. High levels of CRP (increasing up to 100 times) are commonly found during infections (plasma CRP levels increase around 1–500 µg/mL within 24–72 h) [76]. However, the role of CRP isoforms and their involvement in the progression of infectious diseases is still unknown [77,78].

Finally, in addition to the biochemical parameters already mentioned, low calcium levels were also detected in our analysis as a

predictor for COVID-19 diagnosis and disease severity. Calcium is essential for a wide variety of processes in the body, ranging from normal muscle contraction to enzymatic activities [79]. It is known that alteration in Ca^{2+} homeostasis can contribute to cell death by necrosis and apoptosis. Evidence shows that calcium metabolism disorders are associated with cardiovascular disease and early cell death [80,81]. Low calcium levels have already been associated with increased in-hospital mortality in patients with severe coronary artery disease [82], septic patients [79], bacterial pneumonia [83], and patients with dengue [84]. More recent studies have also associated hypocalcemia as an important predictor of hospitalisation and mortality risk by COVID-19 [67,85–87].

Although the present study has shown consistent results, it has some limitations. Cross-sectional studies, with no follow-up analysis of patients' data, are prone to selection bias, information bias, and confounding bias. In addition, biomarkers' levels may change during the disease.

5. Conclusion

All the ML-based models (ANN, DT, PLS-DA, and KNN) were able to effectively predict COVID-19 diagnosis and disease severity with an accuracy above 84%, which is similar to the results obtained by RT-PCR and the minimum recommended threshold for diagnostic tests. The ANN was the model with the best performance (94% and 98%) and, thus, could be used as a supporting decision tool for healthcare professionals in practice. Hyperferritinaemia, hypocalcemia, hypoxemia, pulmonary hypoxia, respiratory acidosis, metabolic acidosis, low urinary pH, and high levels of lactate dehydrogenase were associated with COVID-19 diagnosis and disease severity. These biomarkers are potential therapeutic targets that should be more effectively investigated in further clinical trials.

Funding

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

Authors' contributions

Study concepts: AFC, RP, FST, AW, GRN.
 Study design: AFC, FST, RP, GRN.
 Data acquisition: AFC, MMF, DPS.
 Statistical analysis: AFC, DPS, RP, MS, GRN.
 Manuscript preparation: AFC, AW, MMF, MS, GRN.
 Manuscript editing: AFC, AW, DPS, MMF, MS.
 Manuscript review: AFC, DPS, MMF, MS, FST, RP, AW, GRN.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

The authors express their gratitude to the CAPES (Brazilian Federal Agency for Support and Evaluation of Graduate Education within the Ministry of Education of Brazil) for research funding - Finance Code 001.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.compbiomed.2021.104531>.

References

- [1] World Health Organization, WHO, WHO coronavirus disease (COVID-19) Dashboard 2021, Available from: <https://covid19.who.int/>, 2021.
- [2] A.F. Cobre, B. Böger, M.M. Fachi, R. de O. Vilhena, E.L. Domingos, F.S. Tonin, R. Pontarolo, Risk factors associated with delay in diagnosis and mortality in patients with Covid-19 in the city of Rio de Janeiro, Brazil, *Cienc. Saúde Coletiva* 25 (2020) 4131–4140, <https://doi.org/10.1590/1413-812320202510.2.26882020>.
- [3] A.F. Cobre, M. Surek, R. de O. Vilhena, B. Böger, M.M. Fachi, D.R.O. Momade, F. S. Tonin, F.M. Sarti, R. Pontarolo, Influence of foods and nutrients on COVID-19 recovery: a multivariate analysis of data from 170 countries using a generalized linear model, *Clin. Nutr.* (2021), <https://doi.org/10.1016/j.clnu.2021.03.018>.
- [4] A. Supady, J.R. Curtis, D. Abrams, R. Lorusso, T. Bein, J. Boldt, C.E. Brown, D. Duerschmied, V. Metaxa, D. Brodie, Allocating scarce intensive care resources during the COVID-19 pandemic: practical challenges to theoretical frameworks, *Lancet Respir. Med.* 9 (2021) 430–434, [https://doi.org/10.1016/S2213-2600\(20\)30580-4](https://doi.org/10.1016/S2213-2600(20)30580-4).
- [5] A.L. Ribeiro, N.W. Alves-Sousa, P.R. Martins-Filho, V.O. Carvalho, Social disparity in magnifying glass: the inequality among the vulnerable people during COVID-19 pandemic, *Int. J. Clin. Pract.* 75 (2021) 2–3, <https://doi.org/10.1111/ijcp.13839>.
- [6] A.F. Cobre, B. Böger, R. de O. Vilhena, M.M. Fachi, J.M. dos Santos, F.S. Tonin, A multivariate analysis of risk factors associated with death by Covid-19 in the USA, Italy, Spain, and Germany, *J. Public Health* (2020), <https://doi.org/10.1007/s10389-020-01397-7>.
- [7] E.J. Emanuel, G. Persad, R. Upshur, B. Thome, M. Parker, A. Glickman, C. Zhang, C. Boyle, M. Smith, J.P. Phillips, Fair allocation of scarce medical resources in the time of Covid-19, *N. Engl. J. Med.* 382 (2020) 2049–2055, <https://doi.org/10.1056/NEJMs2005114>.
- [8] P. Weiss, D.R. Murdoch, Clinical course and mortality risk of severe COVID-19, *Lancet* 395 (2020) 1014–1015, [https://doi.org/10.1016/S0140-6736\(20\)30633-4](https://doi.org/10.1016/S0140-6736(20)30633-4).
- [9] S. Bolourani, M. Brenner, P. Wang, T. McGinn, J.S. Hirsch, D. Barnaby, T.P. Zanos, M. Barish, S.L. Cohen, K. Coppa, K.W. Davidson, S. Debnath, L. Lau, T.J. Levy, A. Makhnevich, M.D. Paradis, V. Tóth, A machine learning prediction model of respiratory failure within 48 hours of patient admission for COVID-19: model development and validation, *J. Med. Internet Res.* 23 (2021) 1–15, <https://doi.org/10.2196/24246>.
- [10] W. Liang, H. Liang, L. Ou, B. Chen, A. Chen, C. Li, Y. Li, W. Guan, L. Sang, J. Lu, Y. Xu, G. Chen, H. Guo, J. Guo, Z. Chen, Y. Zhao, S. Li, N. Zhang, N. Zhong, J. He, Development and validation of a clinical risk score to predict the occurrence of critical illness in hospitalized patients with COVID-19, *JAMA Intern. Med.* 180 (2020) 1081–1089, <https://doi.org/10.1001/jamainternmed.2020.2033>.
- [11] X. Guan, B. Zhang, M. Fu, M. Li, X. Yuan, Y. Zhu, J. Peng, H. Guo, Y. Lu, Clinical and inflammatory features based machine learning model for fatal risk prediction of hospitalized COVID-19 patients: results from a retrospective cohort study, *Ann. Med.* 53 (2021) 257–266, <https://doi.org/10.1080/07853890.2020.1868564>.
- [12] M. Pourhomayoun, M. Shakibi, Predicting mortality risk in patients with COVID-19 using machine learning to help medical decision-making, *Smart Health (Amsterdam, Netherlands)* 20 (2021), 100178, <https://doi.org/10.1016/j.smhl.2020.100178>.
- [13] A.J. Heffernan, Host diagnostic biomarkers of infection in the ICU : where are we and where are we going ? *Curr. Infect. Dis. Rep.* 23 (2021) 4, <https://doi.org/10.1007/s11908-021-00747-0>.
- [14] A. Banerjee, S. Ray, B. Vorselaars, J. Kitson, M. Mamalakis, S. Weeks, M. Baker, L. S. Mackenzie, Use of machine learning and artificial intelligence to predict SARS-CoV-2 infection from full blood counts in a population, *Int. Immunopharmacol.* 86 (2020) 106705, <https://doi.org/10.1016/j.intimp.2020.106705>.
- [15] X. Zhou, Z. Wang, S. Li, T. Liu, X. Wang, J. Xia, Y. Zhao, Machine learning-based decision model to distinguish between covid-19 and influenza: a retrospective, two-centered, diagnostic study, *Risk Manag. Healthc. Policy* 14 (2021) 595–604, <https://doi.org/10.2147/RMHP.S291498>.
- [16] V. Schöning, E. Liakoni, C. Baumgartner, A.K. Exadaktylos, W.E. Hautz, A. Atkinson, F. Hammann, Development and validation of a prognostic COVID-19 severity assessment (COSA) score and machine learning models for patient triage at a tertiary hospital, *J. Transl. Med.* 19 (2021) 1–11, <https://doi.org/10.1186/s12967-021-02720-w>.
- [17] A. Banerjee, S. Ray, B. Vorselaars, J. Kitson, M. Mamalakis, S. Weeks, M. Baker, L. S. Mackenzie, Use of machine learning and artificial intelligence to predict SARS-CoV-2 infection from full blood counts in a population, *Int. Immunopharmacol.* 86 (2020) 106705, <https://doi.org/10.1016/j.intimp.2020.106705>.
- [18] J. Wu, P. Zhang, L. Zhang, W. Meng, J. Li, C. Tong, Y. Li, J. Cai, Z. Yang, J. Zhu, M. Zhao, H. Huang, X. Xie, S. Li, Rapid and accurate identification of COVID-19 infection through machine learning based on clinical available blood test results, *MedRxiv* (2020), <https://doi.org/10.1101/2020.04.02.20051136>.
- [19] Kaagle. <https://www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge/discussion/139347>, 2020.
- [20] D. Ballabio, V. Consonni, Classification tools in chemistry. Part 1: linear models. PLS-DA, *Anal. Methods* 5 (2013) 3790–3798, <https://doi.org/10.1039/c3ay40582f>.
- [21] B. Walczak, D.L. Massart, Multiple outlier detection revisited, *Chemometr. Intell. Lab. Syst.* 41 (1998) 1–15, [https://doi.org/10.1016/S0169-7439\(98\)00034-3](https://doi.org/10.1016/S0169-7439(98)00034-3).
- [22] P. Taylor, R.W. Kennard, L.A. Stone, Technometrics computer aided design of experiments, *Technometric* 11 (1969) 137–148.
- [23] D. Ruiz-Perez, H. Guan, P. Madhivanan, K. Mathee, G. Narasimhan, So you think you can PLS-DA? *BMC Bioinf.* 21 (2020) 1–10, <https://doi.org/10.1186/s12859-019-3310-7>.

- [71] M. Parohan, S. Yaghoubi, A. Seraji, Liver injury is associated with severe coronavirus disease 2019 (COVID-19) infection: a systematic review and meta-analysis of retrospective studies, *Hepatol. Res.* 50 (2020) 924–935, <https://doi.org/10.1111/hepr.13510>.
- [72] A. Assiri, J.A. Al-Tawfiq, A.A. Al-Rabeeh, F.A. Al-Rabiah, S. Al-Hajjar, A. Al-Barrak, H. Flemban, W.N. Al-Nassir, H.H. Balkhy, R.F. Al-Hakeem, H. Q. Makhdoom, A.I. Zumla, Z.A. Memish, Epidemiological, demographic, and clinical characteristics of 47 cases of Middle East respiratory syndrome coronavirus disease from Saudi Arabia: a descriptive study, *Lancet Infect. Dis.* 13 (2013) 752–761, [https://doi.org/10.1016/S1473-3099\(13\)70204-4](https://doi.org/10.1016/S1473-3099(13)70204-4).
- [73] F. Al-Hameed, A.S. Wahla, S. Siddiqui, A. Ghabashi, M. Al-Shomrani, A. Al-Thaqafi, Y. Tashkandi, Characteristics and outcomes of middle east respiratory syndrome coronavirus patients admitted to an intensive care unit in Jeddah, Saudi Arabia, *J. Intensive Care Med.* 31 (2016) 344–348, <https://doi.org/10.1177/0885066615579858>.
- [74] M. Saad, A.S. Omrani, K. Baig, A. Bahloul, F. Elzein, M.A. Matin, M.A.A. Selim, M. Al Mutairi, D. Al Nakhli, A.Y.A. Aidaroos, N. Al Sherbeeni, H.I. Al-Khashan, Z. A. Memish, A.M. Albarrak, Clinical aspects and outcomes of 70 patients with Middle East respiratory syndrome coronavirus infection: a single-center experience in Saudi Arabia, *Int. J. Infect. Dis.* 29 (2014) 301–306, <https://doi.org/10.1016/j.ijid.2014.09.003>.
- [75] H.L. Chang, K.T. Chen, S.K. Lai, H.W. Kuo, L.J. Su, R.S. Lin, F.C. Sung, Hematological and biochemical factors predicting SARS fatality in Taiwan, *J. Formos. Med. Assoc.* 105 (2006) 439–450, [https://doi.org/10.1016/S0929-6646\(09\)60183-2](https://doi.org/10.1016/S0929-6646(09)60183-2).
- [76] N.R. Sproston, J.J. Ashworth, Role of C-reactive protein at sites of inflammation and infection, *Front. Immunol.* 9 (2018) 1–11, <https://doi.org/10.3389/fimmu.2018.00754>.
- [77] J.R. Thiele, J. Habersberger, D. Braig, Y. Schmidt, K. Goerendt, V. Maurer, H. Bannasch, A. Scheichl, K.J. Woollard, E. Von Dobschütz, F. Kolodgie, R. Virmani, G.B. Stark, K. Peter, S.U. Eisenhardt, Dissociation of pentameric to monomeric C-reactive protein localizes and aggravates inflammation: in vivo proof of a powerful proinflammatory mechanism and a new anti-inflammatory strategy, *Circulation* 130 (2014) 35–50, <https://doi.org/10.1161/CIRCULATIONAHA.113.007124>.
- [78] I. Huang, R. Pranata, M.A. Lim, A. Oehadian, B. Alisjahbana, C-reactive protein, procalcitonin, D-dimer, and ferritin in severe coronavirus disease-2019: a meta-analysis, *Ther. Adv. Respir. Dis.* 14 (2020) 1–14, <https://doi.org/10.1177/1753466620937175>.
- [79] M.K. Holowaychuk, L.G. Martin, Review of hypocalcemia in septic patients: state-of-the-art review, *J. Vet. Emerg. Crit. Care* 17 (2007) 348–358, <https://doi.org/10.1111/j.1476-4431.2007.00246.x>.
- [80] S.A. Appel, N. Molshatzki, Y. Schwammenthal, O. Merzeliak, M. Toashi, B.A. Sela, D. Tanne, Serum calcium levels and long-term mortality in patients with acute stroke, *Cerebrovasc. Dis.* 31 (2010) 93–99, <https://doi.org/10.1159/000321335>.
- [81] J.W. Chung, W.S. Ryu, B.J. Kim, B.W. Yoon, Elevated calcium after acute ischemic stroke: association with a poor short-term outcome and long-term mortality, *J. Stroke* 17 (2015) 54–59, <https://doi.org/10.5853/jos.2015.17.1.54>.
- [82] S. Di Yan, X.J. Liu, Y. Peng, T.L. Xia, W. Liu, J.Y. Tsauo, Y.N. Xu, H. Chai, F. Y. Huang, M. Chen, D.J. Huang, Admission serum calcium levels improve the GRACE risk score prediction of hospital mortality in patients with acute coronary syndrome, *Clin. Cardiol.* 39 (2016) 516–523, <https://doi.org/10.1002/clc.22557>.
- [83] R.T. Sankaran, J. Mattana, S. Pollack, P. Bhat, T. Ahuja, A. Patel, P.C. Singhal, Laboratory abnormalities in patients with bacterial pneumonia, *Chest* 111 (1997) 595–600, <https://doi.org/10.1378/chest.111.3.595>.
- [84] G.R. Constantine, S. Rajapakse, P. Ranasinghe, B. Parthithpan, P. Jayawardana, A. Wijewickrama, Hypocalcemia is associated with disease severity in patients with dengue, *J. Infect. Dev. Ctries* 8 (2014) 1205–1209, <https://doi.org/10.3855/jidc.4974>.
- [85] J. Liu, P. Han, J. Wu, J. Gong, D. Tian, Prevalence and predictive value of hypocalcemia in severe COVID-19 patients, *J. Infect. Publ. Health* 13 (2020) 1224–1228, <https://doi.org/10.1016/j.jiph.2020.05.029>.
- [86] A. Akirov, A. Gorshtein, I. Shraga-Slutzky, I. Shimon, Calcium levels on admission and before discharge are associated with mortality risk in hospitalized patients, *Endocrine* 57 (2017) 344–351, <https://doi.org/10.1007/s12020-017-1353-y>.
- [87] L. di Filippo, A.M. Formenti, M. Doga, S. Frara, P. Rovere-Querini, E. Bosi, M. Carlucci, A. Giustina, Hypocalcemia is a distinctive biochemical feature of hospitalized COVID-19 patients, *Endocrine* 71 (2021) 9–13, <https://doi.org/10.1007/s12020-020-02541-9>.