

Evaluation of Supervised machine learning classifiers for detecting the degrees of possibility of Coronavirus disease infection

Abbas Rammal (✉ rammal_abbass@hotmail.com)

Lebanese University

Elie Yammine

Lebanese University

Walid Fahs

Islamic University of Lebanon

Rida Khatoun

Telecom ParisTech

Research Article

Keywords: Machine learning, Supervised learning algorithms, Classification, Coronavirus disease, Model prediction, Support Vector Machine, Naïve Bayesian Classifier, Decision Tree Classifier, Linear Discriminant Analysis, Artificial Neural Network, Extended Gamma, Associative Classifier Naive Associative Classifier

Posted Date: March 7th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-958969/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Evaluation of Supervised machine learning classifiers for detecting the**
2 **degrees of possibility of Coronavirus disease infection**

3
4 Abbas Rammal*^{1,2,3}, Elie Yammine², Walid Fahs¹, Rida Khatoun⁴

5
6 ¹Islamic University of Lebanon, Faculty of Engineering, 30014, Wardenieh, Lebanon.

7 ²Lebanese University, Faculty of Science, Statistics and Informatics department, Beirut, Lebanon.

8 ³Lebanese University, Faculty of Information, Data Science department, Beirut, Lebanon.

9 ⁴ Telecom ParisTech, Paris, France

10
11 rammal_abbass@hotmail.com, e.yammine@st.ul.edu.lb, walid.fahs@iul.edu.lb,

12 rida.khatoun@telecom-paris.fr

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 **Corresponding author**

29 Correspondence to Abbas Rammal

30 **First author:** Abbas Rammal, Ph.D., Assistant professor, Lebanese University, Faculty of
31 Science, Statistics and computer sciences department, Beirut, Lebanon,
32 abbas.rammal@ul.edu.lb

33 **Second author:** Elie Yammine, Master's degree in applied statistics, Lebanese University,
34 Faculty of Science, computer science and statistics department, Beirut, Lebanon,
35 e.yammine@st.ul.edu.lb

36 **Third author:** Walid Fahs, Phd, Assistant professor, Islamic University of Lebanon, Faculty
37 of Engineering, transportation, and civil engineering department, Wardenieh, Lebanon,
38 walid.fahs@iul.edu.lb

39 **Fourth author:** Rida Khatoun, Ph.D., Assistant professor, Telecom ParisTech, Paris, France,
40 rida.khatoun@telecom-paris.fr

41

42

43

44

45

46

47

48

49

50

51

52 **1. Abstract**

53 The COVID-19 virus spread from China throughout the world, causing out of control
54 challenges to the public health community. During the epidemic, it is difficult to act against
55 infectious disease due to its unknown trends, so the spread prediction becomes difficult in light
56 of the scarce data. With the absence of treatment for COVID-19 infection, countries have taken
57 some mitigation steps and policies, such as general lockdown and social distance measures, but
58 there has been variation in the extent of viral spread due to several additional factors. Prediction
59 methods based on Artificial intelligence (AI) and Machine Learning (ML) can help in
60 suggesting new policies or even assessing the effectiveness of the existing ones. Such methods
61 have attracted wide attention from researchers implementing statistical modeling and machine
62 learning methods.

63 The objective of this study is to examine different supervised classification approaches to detect
64 the degrees of possibility of Coronavirus disease infection in different countries. Naïve
65 Bayesian Classifier (NBC), Decision Tree Classifier (DTC), Linear Discriminant Analysis
66 (LDA), Support Vector Machine (SVM) and Artificial Neural Network classification (ANN)
67 are machine learning algorithms used for the prediction of coronavirus disease cases according
68 to the time of their evolution while considering data collected from official reports and scientific
69 journals. Since we collected mixed data, we suggested to also apply the supervised classifiers
70 available for mixed data, such as Extended Gamma (EG) and Naive Associative Classifier
71 (NAC).

72 The results showed that ANN and DTC supervised algorithms allow better discrimination
73 between the degrees of possibility of Coronavirus disease infection among advanced methods
74 like NBC, SVM, NAC, EG and LDA.

75

76

77 **2. Keywords**

78 Machine learning, Supervised learning algorithms, Classification, Coronavirus disease,
79 Model prediction, Support Vector Machine, Naïve Bayesian Classifier, Decision Tree
80 Classifier, Linear Discriminant Analysis, Artificial Neural Network, Extended Gamma,
81 Associative Classifier Naive Associative Classifier.

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109 **3. Introduction**

110

111 In December 2019, the novel coronavirus appeared in Wuhan city in China [1] as the first case
112 was reported to the World Health Organization (W.H.O) on 31 December 2019. The virus
113 created a global threat and was announced as coronavirus disease (COVID-19) by the W.H.O
114 on 11 February 2020. The COVID-19 epidemic caused serious global socio-economic turmoil
115 [2]. As of 29 July 2020, more than 17 million cases were recorded. The virus was detected in
116 215 countries causing more than 665,000 deaths [3].

117

118 The early detection of any disease, whether it is infectious or not, is a critical task for rapid
119 treatment to save more lives [4]. Fast diagnosis and screening processes help prevent the spread
120 of pandemic diseases like SARS-CoV-2 [5]. Moreover, predicting the probability degree of
121 getting infected with covid-19 may help national authorities in taking proactive mitigation
122 measures to contain possible epidemic waves of infections. The biggest challenge is to build a
123 technique to improve the detection of covid-19 disease. Several Data Science and Artificial
124 Intelligence techniques have been used to aid the Coronavirus Response [6, 7]. In the literature
125 review, deep learning methods to predict the structure of proteins and their interactions with
126 chemical compounds to facilitate new antiviral drugs/vaccines or recommend current drugs is
127 widely used [8]. A machine learning-based prognostic model to predict if a patient infected
128 with Covid-19 would survive the infection based on age and other risk factors [9]. Data-Based
129 Analysis, Modelling and Forecasting infection rates and spread/patient prognosis to enable
130 hospitals/health officials to better plan resourcing and response [10].

131

132 Machine learning technologies are used to improve the accuracy of prediction for screening
133 both infectious and non-infectious diseases [11]. The machine learning methods were used in
134 modeling former pandemics (e.g., Ebola, Cholera, swine fever, H1N1 influenza [12], dengue

135 fever [13], Zika, oyster norovirus [14]). Table 1 represents notable Machine Learnings methods
 136 used for outbreak prediction.

137

138 Table 1. Notable Machine Learning techniques for outbreak prediction

References	Outbreak infection	Machine learning methods	Resources used to collect data
Geospatial Health [15]	Dengue fever	Neural Network	The computerized database for 2016 has been obtained from Taiwan Centers for Disease Control (Taiwan CDC, 2016), which contains all national dengue fever records daily from 1998, including age, gender, township of residence and time of disease onset.
Malaysian Journal of Public Health Medicine [16]	Dengue/Aedes	Bayesian Network	The Aedes outbreak database sourced from the Ministry of Health at pre-determined localities around the Klang Valley, Malaysia.
Global Ecology and Biogeography [12]	H1N1 flu	Neural Network	Data on the prefectural distribution of the H1N1 influenza pandemic in 2009 were provided by the Ministry of the Environment of Japan.
Current Science [17] Infectious Disease Modelling [18]	Dengue	Naïve Bayes Classification and regression tree (CART)	The data has been collected on the monthly dengue cases from the Government of the National Capital Territory (NCT) of Delhi (India).
Environment International [19]	Oyster norovirus	Neural Network	The epidemiological data were obtained from various online data sources. Specifically, historical norovirus outbreaks in Louisiana oyster harvest areas were collected from Louisiana Morbidity Reports released by Louisiana Department of Health and Hospitals

139

140 Another recent study presents a comparative analysis of machine learning and soft computing
141 models to predict the COVID-19 outbreak as an alternative to Susceptible- Infectious-Removed
142 (SIR) and Susceptible-Exposed-Infectious-Removed (SEIR) models. Among a wide range of
143 machine learning models investigated, two models showed promising results (i.e., multi-
144 layered perceptron, MLP, and adaptive network-based fuzzy inference system, ANFIS) [20].
145 Also, Machine learning and deep learning can replace humans by giving an accurate diagnosis
146 [21]. The perfect diagnosis can save radiologists' time and can be cost-effective than standard
147 tests for COVID-19. X-rays and computed tomography (CT) scans can be used for training the
148 machine learning model. Several initiatives are under-way in this regard. Wang and Wong
149 developed COVID-Net, which is a deep convolutional neural network, which can diagnose
150 COVID-19 from chest radiography images [22]. Moreover, Recent studies [28] show the
151 potential of Artificial Intelligence and Machine Learning tools by suggesting a new model that
152 comes with rapid and valid method SARS-CoV-2 diagnosis using Deep Convolutional Network
153 [23, 24]. In this regard, the remarkable performance suggests the use of the convolutional neural
154 network (Resnet-101) as an adjuvant tool for increase the accuracy of Covid-19 diagnosis.
155 Recent studies used the supervised machine learning techniques for classifying the text into
156 four different categories COVID, SARS, ARDS and Both (COVID, ARDS). Logistic
157 regression and Multinomial Naïve Bayes showed better results than other ML algorithms for
158 detecting COVID-19 using clinical text data set [25].

159

160 In a comparative analysis framework between various supervised machine learning algorithms
161 in diagnosing COVID-19 infections, Pijush Dutta implemented bagging algorithm, k-nearest
162 neighbor, and random forest for classifying the datasets of COVID-19 [26]. Random Forest
163 gave better results with employing accuracy of 85.71%. Haochen Yao and Nan Zhang built a
164 COVID-19 severeness detection model based on supervised machine learning algorithms [27],

165 and obtained the best model using the Support Vector Machine algorithm (SVM) with 28
166 features and overall accuracy of 81.48%. Mahbubunnabi Tamal trained a supervised machine
167 learning algorithms to distinguish between COVID-19 and other diseases, where he found that
168 SVM and Ensemble Bagging Model Trees (EBM) when trained on 71 radiomics features can
169 distinguish between COVID-19 and other diseases with an overall sensitivity of 99.6% and
170 87.8% and specificity of 85% and 97% respectively [28]. Davide Brinati developed two
171 machine learning classification models using hematochemical values from routine blood exams
172 to discriminate between patients who are either positive or negative to the SARS-CoV-2 [29]:
173 their accuracy ranges between 82% and 86%, and sensitivity between 92% and 95%.

174

175 Due to the highly complex nature of the COVID-19 outbreak and variation in its behavior from
176 nation-to-nation, this study suggests applying the machine learning as an effective tool to model
177 the detection and improve the accuracy of prediction for screening coronavirus disease
178 (COVID-19) according to the time of their evolution (in days) in different country. Various
179 supervised machine learning algorithms exist in the literature. The best-known are Linear
180 Discriminant Analysis (LDA) [30], Naïve Bayes Classification (NBC) [31], Decision Tree
181 Classification (DTC) [32, 33], and Support Vector Machine Classification (SVM) [34],
182 Artificial Neural Network (ANN) [35]. Since we collected mixed data, I strongly suggested
183 comparing the classifiers available for mixed data, such as Extended Gamma Associative
184 Classifier (EG) [36], Naive Associative Classifier (NAC) [37]. Irrespective of the method
185 chosen, the use of Cross-validation is a model assessment technique used to evaluate a
186 supervised classification algorithm's performance in making predictions on new COVID-19
187 datasets that it has not been trained on. This paper aims to investigate the generalization ability
188 of the proposed Machine Learning models and the accuracy of the proposed models for

189 detecting the degrees of possibility of Coronavirus disease infection from official reports of
190 COVID-19 data.

191

192 **4. Mathematical Methods**

193 The goal of supervised learning is to build a concise model of the distribution of class labels in
194 terms of predictor features. The resulting classifier is then used to assign class labels to the
195 testing samples where the values of the predictor feature are known, but the value of the class
196 labels are unknown. The classifier's evaluation is most often based on prediction accuracy (the
197 percentage of correct prediction divided by the total number of predictions) [38]. There are at
198 least three techniques which are used to calculate a classifier's accuracy. One technique is to
199 split the available data by using part of them for training and the other for estimating
200 performance, for example two-thirds for training and the other third for testing. In the second
201 technique, known as cross-validation, the training set is divided into mutually exclusive and
202 equal sized subsets and for each subset the classifier is trained on the union of all the other
203 subsets. The average of the error rate of each subset is therefore an estimate of the error rate of
204 the classifier. The third technique is the leave one out validation, which is a special case of cross
205 validation. All test subsets consist of a single sample. This type of a validation is more
206 expensive computationally, but useful when the most accurate estimate of classifier's error rate
207 is required.

208 Supervised classification is one of the tasks most frequently carried out by so-called artificial
209 intelligence. Many techniques have been developed based such as Logic-based techniques,
210 Perceptron-based techniques, Bayesian Networks, etc. In this application, we have chosen
211 Decision Trees classifiers as logic learning methods, and Linear Discriminant Analysis,
212 Support Vector Machines and Naïve Bayes Classification as statistical learning algorithms.
213 These techniques are widely used in several research areas, such as biology, medicine fields

214 and advanced technology [35]. Conversely to artificial neural network, statistical approaches
 215 are characterized by having an explicit underlying probability model, which provides a
 216 probability that a sample belongs in each class, rather than simply a classification. In addition,
 217 the perceptron-based techniques require several parameters which must be studied such that the
 218 size of hidden layers and number of neurons. For this, it is difficult to adapt the neural network
 219 method as an objective function in the methods of selection of variables such as Genetic
 220 Algorithms.

221 Let $\mathbf{X} = \{\mathbf{x}_j\}_{j=1}^n$ be a training set of n samples of observed variables, where each sample is
 222 represented by an S -dimensional vector and let c_k denoting the class membership of \mathbf{x}_j , where
 223 $c_k \in \{c_1, \dots, c_K\}$ with K denoting the number of classes of training set.

224

225 **4.1. Naïve Bayesian classifier**

226 The naïve Bayesian classifier (NBC) is based on Bayes' theorem and predicting independence
 227 requirements. A naïve Bayesian model is simple to construct since it does not need costly
 228 iterative parameter estimates, resulting in a quick computing time for training [39].

229 Given $\mathbf{x}_j = [x_{j1}, \dots, x_{jS}, \dots, x_{jS}]^T \in \mathbb{R}^S$, selection of observed data retrieved from the training
 230 data set \mathbf{X} and belonging to the class c_k . The Naïve Bayes classifier can predict that \mathbf{x}_j belongs
 231 to the predicted class \hat{C} with the following Maximum a Posteriori Probability (MAP) based on
 232 \mathbf{x}_j :

$$233 \quad \text{MAP}(\mathbf{x} = \mathbf{x}_j) = \arg \max_{k \in \{1, \dots, K\}} P(C = c_k | \mathbf{x} = \mathbf{x}_j). \quad (\text{eq. 1})$$

234 According to Bayes' theorem, the probability $P(C = c_k | \mathbf{x} = \mathbf{x}_j)$ that we wish to calculate
 235 could be defined in terms of $P(C = c_k)$, $P(\mathbf{x} = \mathbf{x}_j | C = c_k)$ and $P(\mathbf{x} = \mathbf{x}_j)$ as

$$236 \quad P(C = c_k | \mathbf{x} = \mathbf{x}_j) = \frac{P(\mathbf{x} = \mathbf{x}_j | C = c_k) P(C = c_k)}{P(\mathbf{x} = \mathbf{x}_j)}, \quad (\text{eq. 2})$$

237 where $P(\mathbf{x} = \mathbf{x}_j)$ in the denominator may be discarded since it does not rely on C , and the

238 value of C , and the value of $P(\mathbf{x} = \mathbf{x}_j)$ is a known constant. $P(\mathbf{x} = \mathbf{x}_j|C = c_k)$ is called the
 239 Class-conditional Probability Distribution (CPD). Thus, [40, 41] gives the computation of the
 240 Nave Bayes classifier:

$$241 \quad \text{MAP}(\mathbf{x} = \mathbf{x}_j) = \arg \max_{k \in \{1, \dots, K\}} P(\mathbf{x} = \mathbf{x}_j|C = c_k)P(C = c_k). \quad (\text{eq. 3})$$

242 The likelihood of a classification error, often known as the classifier's risk, is defined as [31]:

$$243 \quad R^{NBC}(\mathbf{x}_j) = \sum_{\mathbf{x}_j \in \mathbf{X}} \mathbf{1}_{C(\mathbf{x}=\mathbf{x}_j) \neq \hat{C}(\mathbf{x}=\mathbf{x}_j)} P(\mathbf{X} = \mathbf{x}_j) = E_{\mathbf{X}} \left\{ \mathbf{1}_{C(\mathbf{x}=\mathbf{x}_j) \neq \hat{C}(\mathbf{x}=\mathbf{x}_j)} \right\}, \quad (\text{eq. 4})$$

244 where $\mathbf{1}_{C(\mathbf{x}=\mathbf{x}_j) \neq \hat{C}(\mathbf{x}=\mathbf{x}_j)}$ is the indicator function and $E_{\mathbf{X}}$ is the expectancy over \mathbf{X} and
 245 $C(\mathbf{x} = \mathbf{x}_j)$ is the real class which the vector \mathbf{x}_j corresponds and $\hat{C}(\mathbf{x} = \mathbf{x}_j)$ is the expected
 246 class of \mathbf{x}_j provided by the NBC method. The Bayes error is the rate of misclassification (Bayes
 247 risk). When addressing classification issues, the Bayes risk is frequently specified as the
 248 reference. It is the least error rate when the distribution is given.

249 **4.2.Linear Discriminant Analysis**

251 Linear Discriminant Analysis (LDA) is a dimensional reducing and classification approach
 252 that has been widely used to the analysis of spectral data. Consider the challenge of allocating
 253 a sample of observed variables retrieved from the training set \mathbf{x}_j to the class c_k in a classification
 254 task. The categorization score is calculated as follows:

$$255 \quad cf(\mathbf{x}_j) = (\mathbf{x}_j - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_j - \mu_k) + \ln|\Sigma_k| - 2 \ln[P(C = c_k)] \quad (\text{eq. 5})$$

256 $P(C = c_k)$ is the prior probability of class c_k , Σ_k is the class covariance matrix of class c_k ,
 257 and μ_k is the mean vector of class c_k , and they are calculated by

$$258 \quad \widehat{\mu}_k = \frac{1}{n_k} \sum_{j=1}^{n_k} \mathbf{x}_j \quad (\text{eq. 6})$$

$$259 \quad \widehat{\Sigma}_k = \frac{1}{n_k} \sum_{j=1}^{n_k} (\mathbf{x}_j - \mu_k)(\mathbf{x}_j - \mu_k)^T \quad (\text{eq. 7})$$

$$260 \quad P(\widehat{C} = c_k) = \frac{n_k}{n} \quad (\text{eq. 8})$$

261 where n_k represents the number of samples in class c_k , and n represents the total number of

262 samples in the training dataset. The sample \mathbf{x}_j is classified as belonging to the class with the
263 lowest classification score [32].

264 When the class covariance matrices are believed to be equal, a pooled covariance matrix is
265 formed.

$$266 \quad \Sigma_{\text{pooled}} = \frac{1}{n} \sum_{k=1}^K n_k \Sigma_k \quad (\text{eq. 9})$$

267 and is substituted for the class covariance matrix in eq. (5). Without taking into account
268 constants, the following classification rule for LDA is produced.

$$269 \quad cf(\mathbf{x}_j) = (\mathbf{x}_j - \mu_k)^T \Sigma_{\text{pooled}}^{-1} (\mathbf{x}_j - \mu_k) - 2 \ln[P(C = c_k)] \quad (\text{eq. 10})$$

270 Equation (10) relates to the Mahalanobis distance when the prior probability n_i is constant. The
271 LDA classifier's error rate R^{LDA} , which is defined as:

$$272 \quad R^{LDA}(\mathbf{x}_j) = \frac{1}{n} \sum_{j=1}^n |sign(c_k(\mathbf{x}_j) - \hat{c}_k(\mathbf{x}_j))|, \quad (\text{eq. 11})$$

273 with \mathbf{x}_j belonging to the class c_k as determined by the LDA method classifier in the class \hat{c}_k

274 **4.3. Decision Tree Classifier**

275 A decision tree classifier is a non-parametric classifier that does not require any a priori
276 statistical assumptions to be made regarding the distribution of data [43]. The basic structure of
277 the decision tree, however, consists of one root node, a number of internal nodes and finally a
278 set of terminal nodes [44]. A node is a subset of the predictors that is used to determine a split.
279 A non-terminal node or parent node is a node that is further split into two child nodes. Growing
280 a tree consists of selecting the optimal splits to determine a non-terminal node, and the
281 assignment of each terminal node to a class. The data is recursively divided down the decision
282 tree according to the defined classification framework.

283 Classes are simply assigned to a terminal node by observing which class is mostly commonly
284 observed in that region of the tree. Thus, the challenge is to optimally choose the best variable
285 and split that variable to maximize the purity or similarity among the responses. The impurity

286 of a parent node τ , denoted $i(\tau)$, is zero when all observations are in the same class. A split s
 287 is determined by selecting the best predictor and split value that optimizes the highest reduction
 288 in purity [45, 46]

$$289 \quad \Delta(s, \tau) = i(\tau) - \sum_{b=1}^B p(\tau_b/\tau) i(\tau_b) \quad (\text{eq. 12})$$

290
 291 where τ_b denotes child node b , $p(\tau_b/\tau)$ is the proportion of observations in τ that are assigned
 292 to τ_b , and B is the number of branches after splitting. Two common impurity functions are the
 293 entropy criterion [33]

$$294 \quad i(\tau) = - \sum_{k=1}^K p_k \log_2(p_k) \quad (\text{eq. 13})$$

295 and the Gini index criterion

$$296 \quad i(\tau) = - \sum_{k=1}^K p_k^2 \quad (\text{eq. 14})$$

297 where p_k is the proportion of observations in class c_k with $k \in \{1, \dots, K\}$. Pruning is based
 298 upon successive steps of removing lower branches that lead to improved classification rates.

299 Once the final tree is determined by $\Delta(s, \tau)$, it is natural to evaluate its predictive performance
 300 by comparing the observed class to the predicted class from the CT for observation \mathbf{x}_j . In a
 301 terminal node m , representing a region R_m with n_m observations, let

$$302 \quad \hat{p}_{mk} = \frac{1}{n_m} \sum_{j=1}^{n_m} \mathbf{1}_{c_k}(\mathbf{x}_j) \quad (\text{eq. 15})$$

303 denote the proportion of class c_k observations in terminal node m . We classify the observations
 304 in node m to class

$$305 \quad \widehat{c}_k(\mathbf{x}_j) = \underset{k}{\operatorname{argmax}} \hat{p}_{mk} \quad (\text{eq. 16})$$

306 The misclassification error rate is simply the proportion of observations in the node that are not
 307 members of the majority class in that node.

$$308 \quad R^{DTC}(\mathbf{x}_j) = \frac{1}{n} \sum_{j=1}^n \left(1 - \max_k(\hat{p}_{mk}(\mathbf{x}_j)) \right) \quad (\text{eq. 17})$$

309

310 4.4. Multi-Class Support Vector Machine (SVM)

311 Support vector machines (SVMs) are a type of learning algorithms that are used for
312 classification and regression. SVM classifiers, like Decision Trees, are non-parametric [47].
313 The one-against-all technique [48] is now the most fundamental technique for implementing
314 SVM multi-class classification. K binary SVM models are built in this simplest application of
315 the SVM to a K-class issue. Class c_k is isolated from the other classes in the k th class SVM
316 problem. To create a final multi-class classifier, all k binary SVM classifiers are concatenated.
317 The term "remaining" refers to the fact that all data points from classes other than c_k are merged
318 to form a single class c_l . Using the usual SVM technique [42], the ideal hyperplane that
319 separates data points from the class c_i and the combined class c_l is obtained. The best separating
320 hyperplane distinguishing the classes c_i and c_k is denoted as:

$$321 \quad g_k(\mathbf{x}_j) = \mathbf{w}_k \cdot \varphi(\mathbf{x}_j) + \mathbf{b}_k \quad \mathbf{k} \in \{1, \dots, K\} \quad (\text{eq. 18})$$

322 where $\mathbf{w}_k \in \mathbb{R}^S$ is the weight vector, \mathbf{b} is the bias, and the transfer function φ maps the training
323 data into an appropriate feature space \mathbb{R}^S to allow for complex nonlinear surfaces. The
324 following minimization is used to estimate the parameters of the decision function $g_k(\mathbf{x}_j)$:

$$325 \quad \min J(\mathbf{w}_k, \xi) = \frac{1}{2} \|\mathbf{w}_k\|^2 + C \sum_{j=1}^n \xi_j \quad (\text{eq. 19})$$

326 subject to

$$327 \quad y_j(\mathbf{w}_k^T \varphi(\mathbf{x}_j) + \mathbf{b}_k) \geq 1 - \xi_j \quad \xi_j \geq 0 ; j = 1 \dots n, \quad (\text{eq. 20})$$

328 To loosen the separability restrictions in with scalar $y_j \in \{-1, +1\}$ signifying its class label,
329 $C \in \mathbb{R}^+$ is a regularization constant, and ξ_j is a slack variable (eq. 19).

330 The $f_k(\mathbf{x}_j)$ decision rule, which allocates the vector \mathbf{x}_j to the class c_k provided by:

$$331 \quad f_k(\mathbf{x}_j) = \text{sign}(g_i(\mathbf{x}_j)) \quad (\text{eq. 21})$$

332 The fundamental challenge with this strategy is that the classifier outputs $f_k(\mathbf{x}_j)$ are binary
333 numbers. The standard approach to this problem is to disregard the sign-operator in eq. 21. We

334 say \mathbf{x}_j is in the class with the biggest value of the decision function and is provided by after
 335 finding all the optimum hyperplanes given by $g_k(\mathbf{x}_j)$ for $\mathbf{k} \in \{1, \dots, K\}$.

$$336 \quad \widehat{c}_k(\mathbf{x}_j) = \underset{k}{\operatorname{argmax}} g_k(\mathbf{x}_j) \quad (\text{eq. 22})$$

337 The index of the biggest component of the different classifiers $g_i(\mathbf{x}_j)$ for $\mathbf{k} \in \{1, \dots, K\}$ is
 338 allocated to the data point \mathbf{x}_j in this manner. The SVM classifier's error rate R^{SVM} , which is
 339 defined as:

$$340 \quad R^{SVM}(\mathbf{x}_j) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}_{c_k \neq \widehat{c}_k}(\mathbf{x}_j), \quad (\text{eq. 23})$$

341 with \mathbf{x}_j belonging to the class c_k as assessed by the method classifier in the class \widehat{c}_k and
 342 $\mathbf{1}_{c_k \neq \widehat{c}_k}(\mathbf{x}_j)$ being the indicator function defined as:

$$343 \quad \mathbf{1}_{c_k \neq \widehat{c}_k}(\mathbf{x}_j) = \begin{cases} 1 & \text{if class } \widehat{c}_k \text{ of } \mathbf{x}_j \neq \text{class } c_k, \\ 0 & \text{if class } \widehat{c}_k \text{ of } \mathbf{x}_j = \text{class } c_k. \end{cases} \quad (\text{eq. 24})$$

344

345

346 **4.5.Artificial Neural Network**

347 Artificial Neural network (ANN) method is an algorithm in machine learning and artificial
 348 intelligence, which is inspired by the human nervous system to analyze and model complex,
 349 nonlinear systems, and parallel computations [49]. Haykin introduces an artificial neural
 350 network as a massively parallel learning machine, made up of simple processing units called
 351 neurons [50].

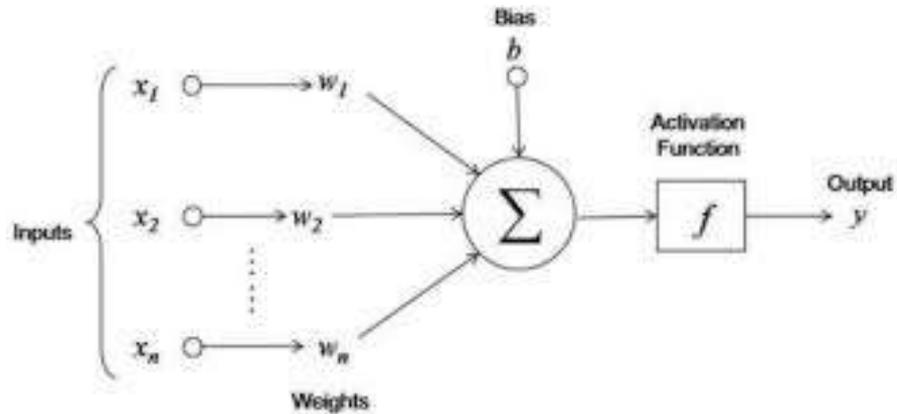
352 A neuron has asset of inputs x_1, x_2, \dots, x_R and each connection from the inputs to the processing
 353 element is affected by different connection strengths known as synaptic weights. A neuron is
 354 showed in mathematical and schematic terms by the following equations (25 ,26) and Figure 1.

$$355 \quad \Sigma = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b \quad (\text{eq. 25})$$

$$356 \quad y = f(\Sigma), \quad (\text{eq. 26})$$

357 Where x_1, x_2, \dots, x_n are the inputs and w_1, w_2, \dots, w_n are the synaptic weights of the neuron
358 L. b is the bias, f is the activation function, and y is the output of the neuron.

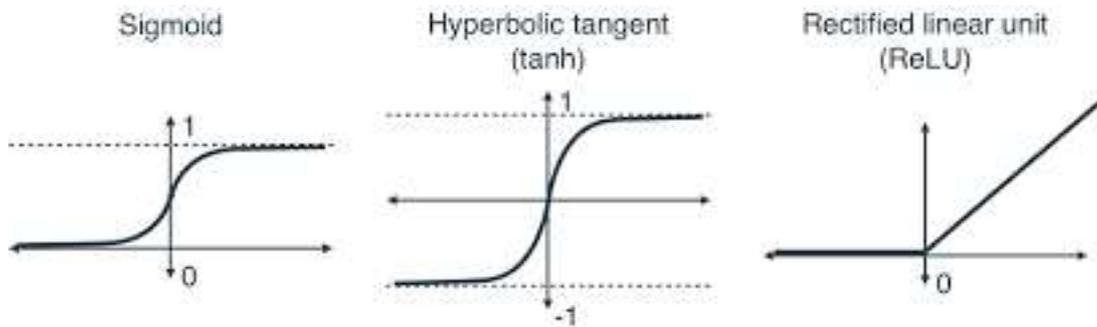
359



360
361
362
363

Figure 1. Structure of a neuron

364 The output of the neuron is defined by applying an activation function to the linear combination
365 of inputs and weights. There are many different activation functions as shown in Figure 2.



366
367

Figure 2. Different activation functions

369

370 The procedure in which neurons are arranged in a neural network is known as the network
371 topology or architecture. There are many different types of neural network such as MLP, RBF
372 and PNN. One of the famous ANN architecture which is most widely used in remote sensing
373 application is multilayer perceptron, or MLP that has been used in this study. An MLP is a feed-

374 forward neural network with one or more layers of neurons between the input and output layer
 375 called hidden layers [51].

376

377 **4.6. Extended Gamma Associative Classifier**

378 The Gamma Associative Classifier (EG) is a supervised learning used as prediction model for
 379 classification. This algorithm has able to handle mixed data, or absence of information in the
 380 data. The extended gamma association classifier (EG) includes two stages: training and
 381 classification. Let X and Y be the training and testing datasets, respectively, where each
 382 instance $x_j \in X$, $y_j \in P$ is described by a set of attributes or factors. [36]

383 The training stage of the classifier starts with the storage of the training set and incorporates the
 384 calculations of various parameters such that the pause and stop parameters. Within the
 385 classification stage, EG employments an iterative handle, based on the calculation of the
 386 average similarity with each class of the samples to be classified [52]. To analyze the similitude
 387 between test sample and training sample, the extended generalized similitude γ_{ext} is utilized.
 388 After getting the similitudes, the average of the generalized similarity of said test pattern for
 389 each class c_k is calculated. The number of samples belonging to the class k in the training data
 390 is given by n , and x_{ij} represents the value of the j -th attribute of the i -th sample of the class c_k ,
 391 and w_j represents the weight of the j th feature with $w_j \in [0 1]$.

$$392 \quad c_k = \frac{\sum_{i=1}^n \sum_{j=1}^m w_j \gamma_{ext}(x_j, y_j, \theta)}{n} \quad (\text{eq. 27})$$

$$393 \quad \gamma_{ext}(x_j, y_j, \theta) = \begin{cases} \gamma_{num}(x_j, y_j, \theta) & \text{if } j\text{th feature is numeric} \\ \gamma_{cat}(x_j, y_j) & \text{if } j\text{th feature is categoric} \\ \gamma_{miss}(x_j, y_j) & \text{if } x_j \text{ or } y_j \text{ are missing} \end{cases} \quad (\text{eq. 28})$$

394 where

$$395 \quad \gamma_{num}(x_j, y_j, \theta) = \begin{cases} 1 & \text{if } |x_j - y_j| \leq \theta \\ 0 & \text{otherwise} \end{cases} \quad (\text{eq. 29})$$

396
$$\gamma_{cat}(x_j, y_j) = \begin{cases} 1 & \text{if } x_j = y_j \\ 0 & \text{otherwise} \end{cases} \quad (\text{eq. 30})$$

397
$$\gamma_{miss}(x_j, y_j) = \begin{cases} 1 & \text{if } x_j = y_j = "?" \\ 0 & \text{otherwise} \end{cases} \quad (\text{eq. 31})$$

398 In case a single greatest is found among all the values of c_k , the process ends. If not, the values
 399 of the stop and pause parameters, as well as the value of the θ parameter, are taken under
 400 consideration in an iterative process [53].

401

402 **4.7. Naïve Associative Classifier**

403 The Naïve Associative Classifier (NAC) a novel supervised learning model, directly handles
 404 mixed and incomplete data [37]. In its training stage, the NAC stores the training set and
 405 calculates, for each numerical attribute, the standard deviation. However, the NAC classifier
 406 needs to store the entire training set of data. It also uses similarity comparisons with respect to
 407 the entire training set to classify a new sample [54].

408 In the classification stage, the Mixed and Incomplete Data Similarity Operator (MIDSO) and
 409 the total similarity operator are used to analyze the similarity between the test samples and the
 410 training samples $s^t(x, y)$.

411

412
$$s^t(x, y) = \sum_{i=1}^m w_i * MIDSO(x, y, A_i) \quad (\text{eq. 32})$$

413
$$MIDSO(x, y, A_i) = \begin{cases} s_n(x, y, A_i) & \text{if } A_i \text{ is numeric} \\ s_c(x, y, A_i) & \text{if } A_i \text{ is categoric} \end{cases} \quad (\text{eq. 33})$$

414
$$s_c(x, y, A_i) = \begin{cases} 0 & \text{if } ((x_j \neq y_j) \vee (x_j = '?' \vee y_j = '?')) \\ 1 & \text{otherwise} \end{cases} \quad (\text{eq. 34})$$

415

416
$$s_n(x, y, A_i) = \begin{cases} 0 & \text{if } (|x_j - y_j| > \sigma_i) \vee (x_j = '?' \vee y_j = '?') \\ 1 & \text{otherwise} \end{cases} \quad (\text{eq. 35})$$

417 With w_i represents the weight of the i th attribute. After having obtained the similarities, we
418 calculated the average of the generalized similarity of the test instance for each class k_l , noted
419 $s_l(\cdot)$.

$$420 \quad s_l(\cdot) = \frac{1}{n} \sum_{y \in k_l} s^t(\cdot, y) \quad (\text{eq. 36})$$

421 If a single maximum is found among all the values of $s_l(\cdot)$, the process ends. If not, any of the
422 classes with maximum similarity is assigned [55].

423

424 **5. Application and results**

425 In this section, we first present data collection that contains the several indicators that directly
426 or indirectly affects the spread of the coronavirus pandemic. The discrimination analysis of
427 samples of infected patients using the proposed approach of supervised classification
428 algorithms is presented and discussed regarding the size of training and validation dataset. All
429 data treatments were performed using the MATLAB® R2019b environment, and scripts are
430 available upon request.

431

432 **5.1.Data set collection**

433 The essential pavement procedure to examine the test hypothesis and to evaluate the predicted
434 outcome is in collecting and processing the data. Thus, to predict the behavior of coronavirus
435 spread, four countries have been chosen. These countries have adopted different methodologies
436 to deal with the coronavirus pandemic and achieved different results related to the methodology
437 used: China, Lebanon, Italy, and Iran. The date range for the data for each of the four countries
438 is as the following:

439 China: 9 January 2020 – 28 March 2020 (80 days)

440 Lebanon: 21 February 2020 – 31 March 2020 (40 days)

441 Italy: 31 January 2020 – 31 March 2020 (61 days)

442 Iran: 19 February 2020 – 31 March 2020 (42 days)

443

444 After studying the situation of the virus in these countries, we noticed that there are several
445 indicators that affect directly or indirectly the spread of the coronavirus. These indicators are
446 presented as follows :

- 447 • The governance reactions: the different responses and reactions from governments of
448 the four countries during the outbreak. These indices are used to explore whether the
449 government response affects the rate of infection and identify correlates of intense
450 responses.
- 451 • The medical resources: This factor refers to the health system policies such as the
452 COVID-19 testing regime or emergency investments into healthcare (ICU beds...) and
453 the health services quality in these four countries. The sensitivity effects of this factor
454 on the results are proposed to be investigated in this study.
- 455 • The quarantine commitment of the people in the country with the government
456 Guidelines. This factor must have theoretically a straight depending on relation with the
457 variation of the spread intensity of the virus between the countries.
- 458 • The special events: this factor takes into consideration the existence of simultaneous
459 events that affect the spread of the coronavirus: other disasters, economic problems,
460 war, political problems or disturbance, official holidays, etc.
- 461 • The economical level and governmental aids: This factor refer to the economic policies,
462 such as income support to citizens or the provision of foreign aid. Depending on the
463 direct relationship between this factor and the quarantine commitment of the people in
464 the country, we have proposed it to be present in this study.
- 465 • Previous experience history of the four governments that determines the existence of
466 experience in such kind of critical disaster management or not.

- 467 • The use of technology devoted to control the virus spread in these countries that help in
 468 the health and hospitalization services, the lockdown control, and the restrictions of the
 469 infected zones.
- 470 • The population age average/ The population density which is considered as the number
 471 of the people in one km²/ in the four countries which causes a meaningful variation in
 472 the situation of the infected records.

473 These direct and indirect factors are considered as parameters used by our model to predict the
 474 future behavior of the spread of COVID-19. The data is combined into a series of novel indices
 475 that aggregate various measures of each factor. Parameters are measured depending on a
 476 specific criterion presented in the following table 2:

477 Table 2. The direct and indirect factors that affecting the spreading level of the coronavirus in
 478 different countries.

479

Effective factors	Indications	Values
Governance reactions (Lockdown)	No lockdown measures at all (value =0).	0, 0.25, 0.5, 0.75, and 1
	Disabling major facilities: Stop flights at airports, public transport, ports, tourist places, universities and school closures and restrictions in movement in the public places (value =0.25)	
	Partially lockdown measures (value=0.5)	
	Fully lockdown measures (value=0.75)	
	Providing virus checks abundantly (value=1)	
Medical resources	0 = no medical resources, 0.25 = low, 0.5 = medium, 0.75 = good, 1 = high	0, 0.25, 0.5, 0.75, and 1
	• Number of nurses for 1000 people	
	• Number of doctors for 1000 people	
	• Number of icu beds for 1000	
Previous experience history	yes / no	1/0

The used technology	0 = no technology, 0.25 = low, 0.5 = medium, 0.75 = good, 1 = high	0, 0.25, 0.5, 0.75, and 1
Special events	yes / no	1/0
Economical resources and governmental aids	0 = no economic resources, 0.25 = low, 0.5 = medium, 0.75 = good, 1 = high	0, 0.25, 0.5, 0.75, and 1
Population density / km2	*	Number
Family number	*	Number
Quarantine commitment (Procedure)	1 = no commitment, 0.75 = low, 0.5 = medium, 0.25 = good, 0 = high	0, 0.25, 0.5, 0.75, and 1
The infected patients by the Covid-19	- new confirmed daily cases	Number

480

481 Thus, we have collected real data for the four countries, from different official sources with all
482 parameters and daily records (recovered, death, confirmed, PCR tests).

483

484 Many factors have been known to be associated with the initial levels of the Coronavirus
485 outbreak at the country level including geographical factors [56], demographical parameters
486 [57], healthcare services [58], and economic status [59].

487 As for the technique of data collection, we relied on studies conducted in these four countries
488 to obtain them. In Iran, the demographical data were extracted from the Statistical Centre of
489 Iran [60], healthcare services and economic levels were extracted from [61]. The number of
490 daily confirmed cases was obtained from the Iranian Ministry of Health and Medical Education
491 website.

492 In China, these factors have been selected using data from the Figshare site, nationwide survey
493 in China and World Health Organization (WHO) [62]. In Italy, the collected data were extracted
494 from the most recently updated databases on the website of the Italian Institute of Statistics

495 (ISTAT) [63], and the ISTAT's Health for All database [64, 65]. In Lebanon, the data has been
496 collected from the Ministry of Health [66], Interior and Finance to obtain these factors and the
497 number of daily confirmed cases.

498 In another study, we conducted a Path analysis which is a statistical technique that uses both
499 bivariate and multiple linear regression techniques to test the causal relations among the
500 variables specialized in the model [67]. We evaluated the parameter estimates of the previously
501 stated regressors on covid-19 daily cases to assess the most influential factors that contribute to
502 the prediction of day-by-day covid-19 infections. The results showed that quarantine
503 commitment ($\beta = -0.823$) and full lockdown measures ($\beta = -0.775$), has the largest direct effect
504 on covid-19 daily infections. The results also show that more experience ($\beta = -0.35$), density
505 in society ($\beta = -0.288$), medical resources ($\beta = 0.136$), and economy resources ($\beta = 0.142$) have
506 indirect effects on covid-19 daily infections.

507 These direct and indirect factors allow to measure the score of each sample which represents
508 the times in days in four countries: China, Iran, Italy, and Lebanon. The daily records are
509 basically the cumulative records for a dependent day-by-day scale. In other words, the record
510 of the next day is the sum of the records of the current day and the new records obtained in the
511 same day. The samples were classified into 3 categories which are the degrees of possibility of
512 Coronavirus disease infection according to their evolution (in days): Low (if the standardized
513 cumulative cases are between 0 and 0.09), medium (if the standardized cumulative cases are
514 between 0.09 and 0.3) and High (if the standardized cumulative cases are above 0.3). The data
515 collected is an imbalanced data meaning that a certain degree of possibility of infection is
516 dominant over other degrees, where in our case the low-risk degree dominates the other two
517 degrees. This is due to the data collected at the beginning of the covid-19 pandemic where most
518 countries are still at low risk of infection and the virus is not fully spread. Regarding the
519 distribution of the degree of possibility of getting infected by covid-19 among the countries

520 studied, China (80 days) and Lebanon (40 days) remained at low risk for the whole durations
521 while Iran maintained a low risk of infection in the first 20 days, medium risk of infection for
522 the next 15 days and high risk of infection for the last 7 days of the 42 total days studied in this
523 country. As for Italy, it remained at a low risk for the first 36 days, medium risk for the next 7
524 days, and high risk for the rest 18 days of the 61 total days studied in Italy.
525 In addition, the data collected do not include missing values. The application of supervised
526 learning methods may result in better accuracy, unless a missing value is expected to have a
527 very high variance. All Covid-19 data were preprocessed with a Standard Normal Variate
528 (SNV) preprocessing method [68].

529

530 **5.2. Analysis of discrimination of samples of infected patients by the Covid-19**

531 To deepen the obtained results, we analyzed the confusion matrices and the ROC curves. A
532 confusion matrix is used to describe the performance of a classification model on a set of test
533 data for which the true values are known. It is a summary of prediction results on a classification
534 problem and allows the visualization of the performance of an algorithm. Hence, the elements
535 in the diagonal are the elements correctly classified, while the elements out of the diagonal are
536 misclassified. It can easily be seen that the overall accuracy can be computed by dividing the
537 sum of the diagonal elements (number of correctly classified samples) by the total sum of the
538 elements of the matrix (total number of samples in the dataset).

A Confusion Matrix

Output Class	Low	59 79.7%	2 2.7%	0 0.0%	96.7% 3.3%
	Medium	0 0.0%	5 6.8%	6 8.1%	45.5% 54.5%
	High	0 0.0%	0 0.0%	2 2.7%	100% 0.0%
		100% 0.0%	71.4% 28.6%	25.0% 75.0%	89.2% 10.8%
	Low	Medium	High	Target Class	

B Confusion Matrix

Output Class	Low	59 79.7%	1 1.4%	0 0.0%	98.3% 1.7%
	Medium	0 0.0%	6 8.1%	3 4.1%	66.7% 33.3%
	High	0 0.0%	0 0.0%	5 6.8%	100% 0.0%
		100% 0.0%	85.7% 14.3%	62.5% 37.5%	94.6% 5.4%
	Low	Medium	High	Target Class	

C Confusion Matrix

Output Class	Low	59 79.7%	7 9.5%	0 0.0%	89.4% 10.6%
	Medium	0 0.0%	0 0.0%	4 5.4%	0.0% 100%
	High	0 0.0%	0 0.0%	4 5.4%	100% 0.0%
		100% 0.0%	0.0% 100%	50.0% 50.0%	85.1% 14.9%
	Low	Medium	High	Target Class	

D Confusion Matrix

Output Class	Low	59 79.7%	0 0.0%	0 0.0%	100% 0.0%
	Medium	0 0.0%	7 9.5%	1 1.4%	87.5% 12.5%
	High	0 0.0%	0 0.0%	7 9.5%	100% 0.0%
		100% 0.0%	100% 0.0%	87.5% 12.5%	98.6% 1.4%
	Low	Medium	High	Target Class	

E Confusion Matrix

Output Class	Low	59 79.7%	1 1.4%	0 0.0%	98.3% 1.7%
	Medium	0 0.0%	5 6.8%	0 0.0%	100% 0.0%
	High	0 0.0%	1 1.4%	8 10.8%	88.9% 11.1%
		100% 0.0%	71.4% 28.6%	100% 0.0%	97.3% 2.7%
	Low	Medium	High		Target Class

F Confusion Matrix

Output Class	Low	54 73.0%	0 0.0%	0 0.0%	100% 0.0%
	Medium	5 6.8%	7 9.5%	3 4.1%	46.7% 53.3%
	High	0 0.0%	0 0.0%	5 6.8%	100% 0.0%
		91.5% 8.5%	100% 0.0%	62.5% 37.5%	89.2% 10.8%
	Low	Medium	High		Target Class

544

G Confusion Matrix

Output Class	Low	57 77.0%	0 0.0%	0 0.0%	100% 0.0%
	Medium	2 2.7%	7 9.5%	3 4.1%	58.3% 41.7%
	High	0 0.0%	0 0.0%	5 6.8%	100% 0.0%
		96.6% 3.4%	100% 0.0%	62.5% 37.5%	93.2% 6.8%
	Low	Medium	High		Target Class

545

546 Table 2. Results of average confusion matrix (of 100 runs) for the training (2/3 of samples) and

547 testing (1/3 of samples) dataset using the supervised algorithms (A) LDA (B) NBC (C) SVM

548 (D) DTC (E) ANN (F) EG (G) NAC methods for the detection of infected patients by the Covid-
549 19 according to the time of evolution (in days) of coronavirus disease in different country

550

551 Cross-validation is a model assessment technique used to evaluate a supervised classification
552 algorithm's performance in making predictions on new datasets that it has not been trained on.
553 This is done by partitioning a dataset and using a subset to train the algorithm and the remaining
554 data for testing. Several methods of cross validation exist such as k-fold, Holdout, Leave-out
555 and Re-substitution. We use k-fold method that partitions samples of infected patients by the
556 Covid-19 into $k = 2$ until 10 randomly chosen subsets of roughly equal size. With this approach
557 there is not a possibility of high bias if we have limited data, because we would not miss some
558 information about the data which we have not used for training. the discussion of choosing the
559 value of K in K-fold cross validation, is shown in section 5.3. but here we will display the
560 results of K = 3-fold. One subset is used to validate the model (1/3 of data = 74 samples) trained
561 using the remaining subsets (2/3 of data = 149 samples). This process is repeated T=100 times
562 to get the best stable model and to ensure that most of our samples of infected patients by the
563 Covid-19 have been used for validation.

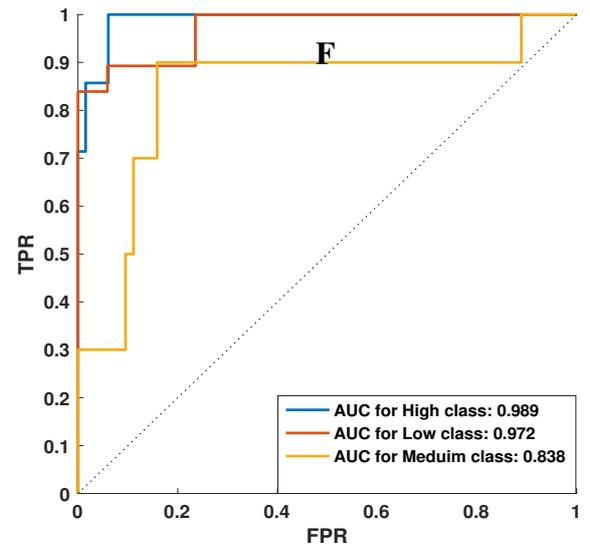
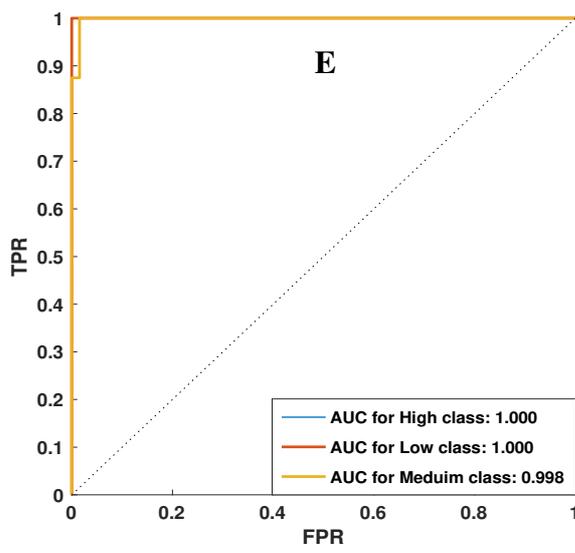
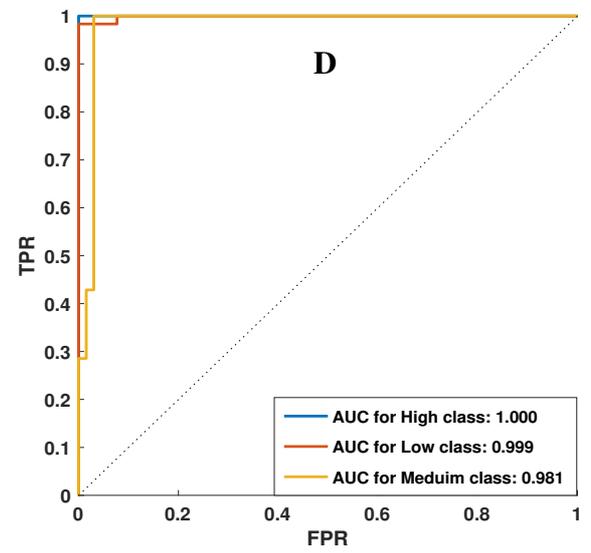
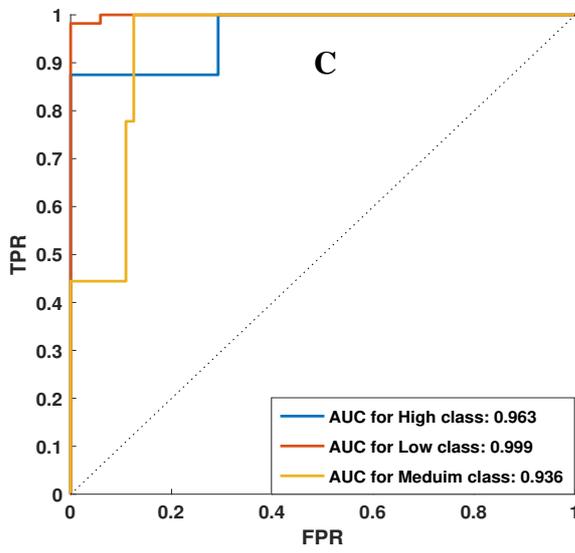
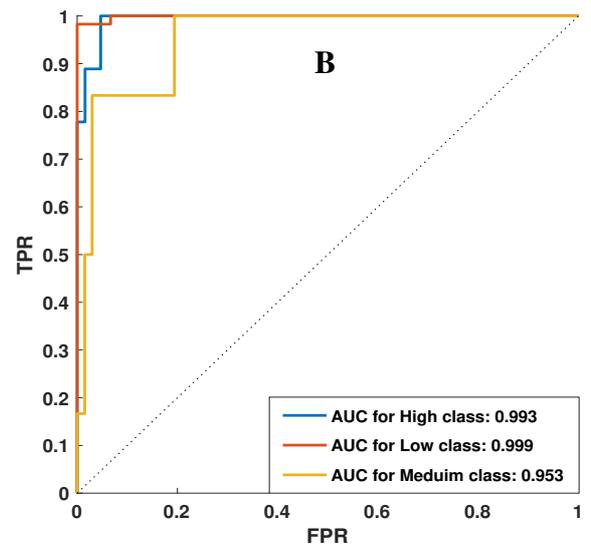
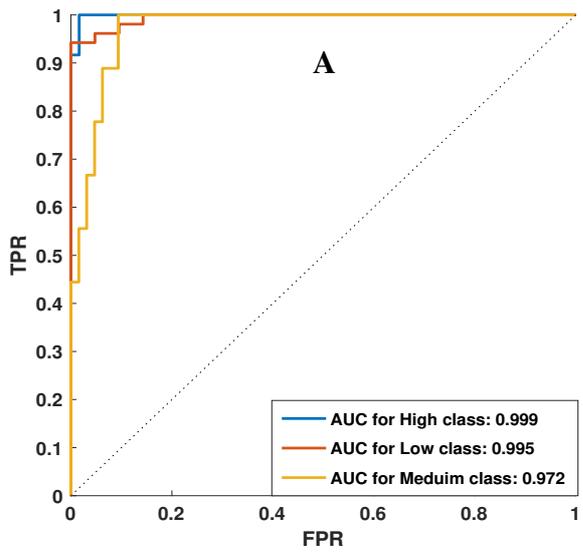
564 The Linear Discriminant Analysis (LDA), Naïve Bayes Classifier (NBC), Support Vector
565 Machine (SVM), Decision Tree Classifier (DTC), Artificial Neural Network (ANN), Extended
566 Gamma Associative Classifier (EG) and Naive Associative Classifier (NAC) proposed
567 supervised algorithms has been applied on the coronavirus data. Tables 2 (A)(B)(C)(D)(E)
568 presents the performance of the LDA, NBC, SVM, DTC and ANN classifications as the mean
569 over the T=100 samplings of the confusion matrix information for the samples of infected
570 patients by the Covid-19 according to the time of evolution (in days) of coronavirus disease in
571 different country. Each table show the confusion matrix for the true label's targets and predicted
572 labels outputs. The rows correspond to the predicted class (Output Class) extract by
573 classification methods and the columns correspond to the infected samples degree true class

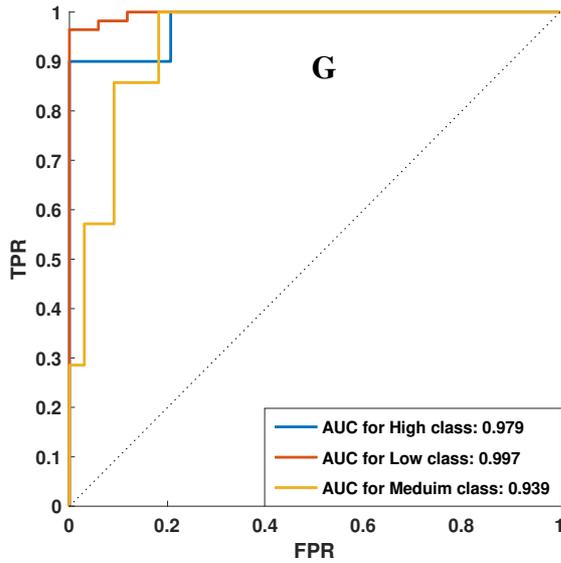
574 (Target Class). We have three degrees of possibility according to their evolution of coronavirus
575 that represents the classes of our COIVID-19 data: Low, Medium and Hight. The diagonal cells
576 correspond to spectra that are correctly classified. The off-diagonal cells correspond to
577 incorrectly classified of degree of infected samples (time of evolution) by the Covid-19.
578 The results of classification accuracy for Detection True Classifier (DTC) and Artificial Neural
579 Networks (ANN) algorithms registered 98.6% and 97.4% of validation datasets describing the
580 degree of possibility of Coronavirus disease infection. However, the classification accuracy is
581 worse for the Naïve Bayes Classifier (NBC) algorithm with 94.6 % of validation datasets
582 correctly classified, the Support Vector Machine (SVM) algorithm with 85.1 % of validation
583 datasets correctly classified and Linear Discriminant Analysis (LDA) algorithm with 89.8 % of
584 validation datasets correctly classified. These values, presented in Table 3, confirm that the
585 DTC and ANN algorithm allows a better discriminate within the degrees of possibility of
586 Coronavirus disease infection according to days of their evolution than state of art methods of
587 supervised classifications such as NBC, SVM, EG, NAC and LDA algorithms.

588
589 A ROC (Receiver Operating Characteristic) curve summarizes the performance of a classifier
590 over all possible thresholds. It is generated by plotting the True Positive Rate (y-axis) against
591 the False Positive Rate (x-axis) by varying the threshold for assigning observations to a given
592 class. A ROC curve is a performance measurement for classification problem at various
593 thresholds settings. The plot of True Positive Rate (TPR; sensitivity) versus False Positive Rate
594 (FPR; 1-specificity) across varying cut-offs generates a curve in the unit square called a ROC
595 curve. ROC curves are used to compare different supervised classifiers. In our case we have
596 used the multi-class ROC curves, which is a kind of multi-objective optimization to compare
597 the supervised classification algorithm; LDA, NBC, SVM, ANN, DTC, EG and NAC. As in
598 several multi-class problem, the idea is generally to carry out pairwise comparison: one class

599 vs. all other classes, one class vs. another class [69]. If the curve has high values, it leads to the
600 greater of area under the curve obtained, and the less error the classifier makes. ROC curve
601 corresponding to progressively greater discriminant capacity are located progressively closer
602 to the upper left-hand corner in "ROC space". To numerically assess the discrimination of the
603 classes, the Area Under the Curve (AUC) was used to measure of the ability of a classifier to
604 distinguish between classes and is used as a summary of the ROC curve. Higher the AUC,
605 better the model is at distinguishing between classes of degrees of possibility of Coronavirus
606 disease.

607 The figure below represents the comparisons between ROC curves by LDA, NBC, SVM, ANN
608 and DTC for each class of degrees of possibility of Coronavirus disease infection according to
609 days of their evolution for validation datasets (1/3 samples). The ROCs of each model are
610 plotted using the average of 100 runs. We can see that the classification threshold in ANN and
611 DTC methods (figures 3D and 3E) are very close to 1 and very far from diagonal for all the
612 classes of degrees of possibility of Coronavirus disease compared to other methods. But the
613 classification threshold is relatively reduced for the LDA and SVM methods (figures 3A and
614 3B). The area under the curves of each class of degrees of possibility of Coronavirus disease
615 for the ANN and DTC methods are higher than that for NBC, SVM, LDA, EG and NAC
616 methods. This also confirms that the measures of sensitivity and specificity of ROC curves for
617 ANN and DTC algorithms are generally stronger than the other ROC curves for NBC, SVM,
618 LDA EG and NAC algorithms for each class of degrees of possibility of Coronavirus disease
619 in all figures 3(A-G), which indicates better in-sample accuracy than the other classifier
620 methods. The AUC values, presented in Figures 3 (A-G), confirm that the ANN and DTC
621 supervised methods allows a better classifier of degrees of possibility of Coronavirus disease.





625

626 Figure 3. The average of 100 iterations of Receiver operating characteristic (ROC) curve of A)
 627 LDA, (B) NBC, (C) SVM, (D) ANN and (E) DTC (F) EG (G) NAC supervised classifications
 628 algorithms for classify the degrees of possibility of Coronavirus disease infection according to
 629 days of their evolution using the validation datasets.

630

631 **5.3.Effect of the size of testing and training datasets**

632 As we already mentioned, cross validation is the most important part of machine learning
 633 models, and the results might depend on it. K-fold cross-validation is extremely useful with
 634 the appropriate K which is the number of folds chosen.

635 Unfortunately, there is no theoretically perfect procedure of determining the appropriate value
 636 of K in K-fold cross validation. Many researchers used $K = 10$, for example, see Bengio and
 637 Grandvalet (2004), Vehtari and Lampinen (2004). Davison and Hinkley (1997) recommend K
 638 $= \min(n^{1/2}, 10)$ in practice. Clark (2003) compares different values of K for his data and
 639 suggested that the choice of $K = 4$ will probably be good in general, though not necessarily
 640 optimal. Duan et al. (2003) used $K = 5$ for the number of folds. Anderson et al (2006)
 641 suggested, leaving 20% of the n-samples at a time for final model validation. For this, we have

642 tested the method of supervised learning classification LDA, NBC, SVM, ANN and DTC on
643 different partitions of testing and training dataset depending on different values of K between
644 2 and 10. To examine the influence of the cross-validation method, we have computed the ratio
645 of classification (%) obtained for each method of the supervised learning classification LDA,
646 NBC, SVM, ANN and DTC on COVID-19 data. The values in Table 3, indicate that the
647 detection tree classifier (DTC) algorithm gives the best results of classification for different
648 sizes of testing and training datasets for the samples of infected patients by Covid-19 according
649 to the time of evolution (in days) of coronavirus disease in different countries. This sampling
650 analysis is done to assess the importance of K-fold cross-validation method in determining the
651 best classifier among the classifiers studied.

652

653 Table 3. The Means and Standard Deviation of classification performance (expressed as %) for
654 100 runs of LDA, NBC, SVM, ANN, DTC, EG and NAC supervised methods on different sizes
655 of testing and training dataset of coronavirus diseases or different values of K in K-fold cross
656 validation.

	K-fold	LDA	NBC	SVM	ANN	DTC	EG	NAC
1/2 training & 1/2 testing	2	87.4 ± 2.7	91.4 ± 2.2	82.3 ± 4.5	91.2 ± 2.8	96.9 ± 1.7	83.1 ± 5.2	91.8 ± 3.3
2/3 training & 1/3 testing	3	89.5 ± 3.4	94.5 ± 2.5	85.6 ± 5.2	97.3 ± 1.8	98.6 ± 0.7	89.2 ± 4.2	93.2 ± 2.6
3/4 training & 1/4 testing	4	89.5 ± 3.2	94.5 ± 1.9	87.1 ± 4.4	98.1 ± 0.8	98.9 ± 0.5	90.4 ± 3.8	93.3 ± 2.6
4/5 training & 1/5 testing	5	90.2 ± 3.4	95.4 ± 2.6	87.4 ± 3.8	99.1 ± 0.5	99.2 ± 0.3	91.2 ± 3.8	93.5 ± 2.4
9/10 training & 1/10 testing	10	91.1 ± 2.8	95.8 ± 2.2	88.9 ± 3.8	99.5 ± 0.2	99.5 ± 0.3	91.4 ± 3.2	94.5 ± 2.2

657

658

659 To quantify the results of the AUROC of each model performed 100 runs, we calculated in The
 660 Table 4 the mean and standard deviation of the AUROC for three classes of degrees of
 661 possibility of infection with Coronavirus disease infection and the proposed supervised learning
 662 methods.

663

664 Table 4. The mean and SD of the AUROC for three classes of degrees of possibility of
 665 Coronavirus disease infection.

	Low	Medium	High
LDA	0.995 ± 0.004	0.972 ± 0.024	0.999 ± 0.001
NBC	0.999 ± 0.001	0.953 ± 0.042	0.993 ± 0.006
SVM	0.999 ± 0.001	0.936 ± 0.064	0.963 ± 0.032
DTC	0.999 ± 0.001	0.981 ± 0.018	1.000 ± 0.000
ANN	1.000 ± 0.000	0.998 ± 0.001	1.000 ± 0.000
EG	0.972 ± 0.027	0.838 ± 0.158	0.989 ± 0.001
NAC	0.979 ± 0.026	0.939 ± 0.058	0.979 ± 0.028

666

667 In application of machine learning algorithms for classification tasks, we need to apply
 668 approximate statistical test for determining whether one learning algorithm out-performs
 669 another on a particular learning task. For this we propose to apply McNemar's test that is based
 670 on a X2 test for goodness of fit. So, we may reject the null hypothesis in favor of the hypothesis
 671 that the two algorithms have different performance if P-value is greater than 0.05. We found
 672 that the application of McNemar test between seven used supervised algorithms, gives the
 673 existence of performance difference between DTC and ANN algorithms and other algorithms,
 674 and there is no performance difference between DTC and ANN algorithms.

675

676 **6. Conclusions and Future Research**

677 COVID-19 has shocked the world due to the speed of its spread and its non-availability of
678 vaccine or drug. Various researchers are working for conquering this deadly virus. We used
679 224 reports recorded on specified dates in four countries which are labelled in three classes of
680 degrees of infections for each day: Low, medium and High possibility of Coronavirus disease
681 infection. Various direct and indirect factors like the economical level and governmental aids,
682 medical resources, previous experience, Governance reactions and the used technology are
683 being extracted from these COVID-19 samples according to the time of evolution (in days) of
684 coronavirus disease in different countries.

685
686 The machine learning algorithms are used for classifying daily samples in four countries into
687 three different classes degrees of possibility of Coronavirus disease infection. After performing
688 classification, it was revealed that decision trees classifier (DTC) and artificial neural networks
689 (ANN) gives excellent results by having 98.6% and 97.3% accuracy well classified. Other
690 machine learning algorithm that showed better results were Naive Bayesian classifier (NBC)
691 by having 94.5% accuracy well classified. The efficiency of models can be improved by
692 increasing the amount of data.

693
694 Further studies should be performed to investigate some issues, for example, a study that takes
695 into account a priori information such as the degrees of possibility of Coronavirus disease into
696 the objective function through Partial Least Square (PLS) models. Furthermore, it would be
697 very interesting to apply of new technology in the search for the degrees of possibility of
698 Coronavirus disease infection such as the deep learning technique which allows to extract the
699 spectral features to classify the state of COVID-19.

700

701 **7. Declarations**

702 **Acknowledgements**

703 This study was carried out by Faculty of Engineering-Islamic University of Lebanon in
704 partnership with Lebanese University, Statistics and Informatics Department - Faculty of
705 Science in Beirut.

706 **Funding**

707 This project has been funded with the support of the National Council for Scientific Research
708 in Lebanon CNRS-L.

709 **Ethics approval and consent to participate**

710 The study does not involve human participants. Not applicable

711 **Consent for publication**

712 Not applicable

713 **Availability of data and material**

714 All data, models, and code generated or used during the study appear in the submitted article
715 and are provided upon request by contacting Abbas Rammal via email:
716 rammal_abbass@hotmail.com

717 **Competing interests**

718 The authors declare that there are no competing interests.

719 **Authors' contributions**

720 AR: Conceptualization, Methodology, Software, Data curation, Visualization, Validation,
721 Investigation, Original and revised draft preparation, Writing-Reviewing, and Editing. EY:
722 Data curation, Investigation, Original and revised draft preparation, Writing-Reviewing, and
723 Editing. WF and RK: Conceptualization, Methodology, Data Collection, and Writing. All
724 authors read and approved the final manuscript.

725 **8. References**

- 726 [1] Wu, F., Zhao, S., Yu, B., Chen, Y., Wang, W., Song, Z.G., Hu, Y., Tao, Z.W., Tian, J.H.,
727 Pei, Y.Y., Yuan, M.L., Zhang, Y.L., Dai, F.H., Liu, Y., Wang, Q.M., Zheng, J.J., Xu, L.,
728 Holmes, E.C., Zhang, Y.Z. 2020. A new coronavirus associated with human respiratory disease
729 in china. *Nature*. 44(59), 265–269.
- 730 [2] Ardabili, S.F., Mosavi, A., Ghamisi, P., Ferdinand, F., Varkonyi-Koczy, A.R., Reuter, U.,
731 Rabczuk, T., Atkinson, P.M. 2020. COVID-19 Outbreak Prediction with Machine Learning
732 <http://dx.doi.org/10.2139/ssrn.3580188>
- 733 [3] Ivanov, D. 2020. Predicting the impacts of epidemic outbreaks on global supply chains: A
734 simulation-based analysis on the coronavirus outbreak (COVID-19/SARS-CoV-2) case.
735 *Transp. Res. Part E Logist.* 136, <https://doi:10.1016/j.tre.2020.101922>
- 736 [4] Vaishya, R., Javaid, M., Khan, I.H., Haleem, A. 2020. Artificial Intelligence (AI)
737 applications for COVID-19 pandemic. *Diabetes Metab Syndr.* 14(4), 337–339.
- 738 [5] Ai, T., Yang, Z., Hou, H., Zhan, C., Chen, C., Lv, W., Tao, Q., Sun, Z., Xia, L. 2020
739 Correlation of chest CT and RT-PCR testing in coronavirus disease 2019 (COVID-19) in China:
740 a report of 1014 cases. *Radiology*. 296(2), 32-40.
- 741 [6] Pham, Q., Nguyen, D.C., Huynh-The, T., Hwang, W., Pathirana, P.N. 2020. Artificial
742 Intelligence (AI) and Big Data for Coronavirus (COVID-19) Pandemic: A Survey on the State-
743 of-the-Arts. *IEEE Access*. 8, 130820-130839.
- 744 [7] Luengo-Oroz, M., Pham, K.H., Bullock, J., Kirkpatrick, R., Luccioni, A., Rubel, S.,
745 Wachholz, C., Chakchouk, M., Biggs, P., Nguyen, T., Purnat, T., Mariano, B. 2020. Artificial
746 intelligence cooperation to support the global response to COVID-19. *Nature Machine*
747 *Intelligence*. 2, 295–297.
- 748 [8] Robson, B. 2020. COVID-19 Coronavirus spike protein analysis for synthetic vaccines, a
749 peptidomimetic antagonist, and therapeutic drugs, and analysis of a proposed achilles' heel
750 conserved region to minimize probability of escape mutations and drug resistance. *Comput.*

751 Biol. Med. 121, 103749.

752 [9] Wynants, L., Van Calster, B., Collins, G.S., Riley, R.D., Heinze, G., Schuit, E., Bonten,
753 M.M.J., Damen, J.A.A., Debray, T.P.A., De Vos, M., Dhiman, P., Haller, M.C., Harhay, M.O.,
754 Henckaerts, L., Kreuzberger, N., Lohman, A., Luijken, K., Ma, J., Andaur, C.L., Reitsma, J.B.,
755 Sergeant, J.C., Shi, C., Skoetz, N., Smits, L.J.M., Snell, K.I.E., Sperrin, M., Spijker, R.,
756 Steyerberg, E.W., Takada, T., van Kuijk, S.M.J., van Royen, F.S., Wallisch, C., Hooft, L.,
757 Moons, K.G.M., van Smeden, M. 2020. Prediction models for diagnosis and prognosis of covid-
758 19 infection: systematic review and critical appraisal. *BMJ*. 369.

759 [10] Anastassopoulou, C., Russo, L., Tsakris, A., Siettos, C. 2020 Data-based analysis,
760 modelling and forecasting of the COVID-19 outbreak. *PLoS ONE*. 15(3).

761 [11] Agrebi, S., Larbi, A. 2020. Use of artificial intelligence in infectious diseases. *Artificial*
762 *Intelligence in Precision Health*. 415-438.

763 [12] Koike, F., Morimoto, N. 2018. Supervised forecasting of the range expansion of novel
764 non-indigenous organisms: Alien pest organisms and the 2009 H1N1 flu pandemic. *Global*
765 *Ecol. Biogeogr.* 27, 991-1000.

766 [13] Agarwal, N., Koti, S.R., Saran, S., Kumar, S. 2018. Data mining techniques for predicting
767 dengue outbreak in geospatial domain using weather parameters for New Delhi, India. *Curr.*
768 *Sci.* 114, 2281-2291.

769 [14] Raja, D.B., Mallol, R., Ting, C.Y., Kamaludin, F., Ahmad, R., Ismail, S., Jayaraj, V.J.,
770 Sundram, B.M. 2019. Artificial intelligence model as predictor for dengue outbreaks. *Malays.*
771 *J. Public Health Med.* 19, 103-108.

772 [15] Anno, S., Hara, T., Kai, H., Lee, M.A., Chang, Y., Oyoshi, K., Mizukami, Y., Tadono, T.
773 2019. Spatiotemporal dengue fever hotspots associated with climatic factors in taiwan including
774 outbreak predictions based on machine-learning. *Geospatial Health*. 14, 183-194.

775 [16] Raja, D.R., Mallol, R., Ting, C. Y., Kamaludin, F., Rohani Ahmad, Suzilah Ismail, Vivek

776 Jason Jayaraj, & Bala Murali Sundram. (2019). ARTIFICIAL INTELLIGENCE MODEL AS
777 PREDICTOR FOR DENGUE OUTBREAKS. *Malaysian Journal of Public Health Medicine*,
778 19(2), 103-108.

779 [17] Agarwal N, Koti SR, Saran S, Kumar AS. Data mining techniques for predicting dengue
780 outbreak in geospatial domain using weather parameters for New Delhi, India. *Current Science*
781 journal. 2018; 114: 2281-2291.

782 [18] Titus Muurlink, O., Stephenson, P., Islam, M.Z., Taylor-Robinson, A.W. 2018. Long-term
783 predictors of dengue outbreaks in Bangladesh: A data mining approach. *Infectious Disease*
784 *Modelling*. 3, 322-330.

785 [19] Chenar, S.S., Deng, Z. 2018. Development of artificial intelligence approach to forecasting
786 oyster norovirus outbreaks along Gulf of Mexico coast. *Environment International*. 111, 212-
787 223.

788 [20] Yang, Z., Zeng, Z., Wang, K., Wong, S.S., Liang, W., Zanin, M., Liu, P., Cao, X., Gao,
789 Z., Mai, Z. 2020. Modified SEIR and AI prediction of the epidemics trend of COVID-19 in
790 China under public health interventions. *Journal of Thoracic Disease*.

791 [21] Ozturk, T., Talo, M., Yildirim, E.A., Baloglu, U.B., Yildirim, O., Rajendra Acharya, U.
792 2020. Automated detection of COVID-19 cases using deep neural networks with X-ray images.
793 *Computer Biol Med*. 121, 103792.

794 [22] Wang, L., Wong, A. 2020. COVID-Net: a tailored deep convolutional neural network
795 design for detection of COVID-19 Cases from chest radiography images.

796 [23] Hemdan, E.E.D., Shouman, M.A., Karar, M.E. 2020. COVIDX-Net: A Framework of
797 Deep Learning Classifiers to Diagnose COVID-19 in X-Ray Images.
798 <https://arxiv.org/abs/2003.11055>

799 [24] Wang, S., Kang, B., Ma, J., Zeng, X., Xiao, M., Guo, J., Cai, M., Yang, J., Li, Y., Meng,
800 X., Xu, B. 2020. A deep learning algorithm using CT images to screen for Corona Virus Disease

801 (COVID-19).

802 [25] Khanday, A.D., Rabani, S.T., Khan, Q.R. 2020. Machine learning based approaches for
803 detecting COVID-19 using clinical text data. *International Journal of Information Technology*.
804 12, 731–739.

805 [26] Dutta, P., Paul, S., Kumar, A. 2021. Comparative analysis of various supervised machine
806 learning techniques for diagnosis of COVID-19. *Electronic Devices, Circuits, and Systems for*
807 *Biomedical Applications*, 521–540. <https://doi.org/10.1016/B978-0-323-85172-5.00020-4>

808 [27] Yao, H., Zhang, N., Zhang, R., Duan, M., Xie, T., Pan, J., Peng, E., Huang, J., Zhang, Y.,
809 Xu, X., Xu, H., Zhou, F., Wang, G. 2020. Severity Detection for the Coronavirus Disease 2019
810 (COVID-19) Patients Using a Machine Learning Model Based on the Blood and Urine
811 Tests. *Front. Cell Dev. Biol.* 8, 683. doi: 10.3389/fcell.2020.00683

812 [28] Tamal, M., Alshammari, M., Alabdullah, M., Hourani, R., Abu Alola, H., Hegazi, T.
813 2020. An Integrated Framework with Machine Learning and Radiomics for Accurate and Rapid
814 Early Diagnosis of COVID-19 from Chest X-ray. *Expert Systems with Applications*,
815 <https://doi.org/10.1101/2020.10.01.20205146>

816 [29] Brinati, D., Campagner, A., Ferrari, D., Locatelli, M., Banfi, G., Cabitza, F. 2020.
817 Detection of COVID-19 Infection from Routine Blood Exams with Machine Learning: a
818 Feasibility Study. *Journal of Medical Systems*, <https://doi.org/10.1101/2020.04.22.20075143>

819 [30] Huerta, E.B., Duval, B., Hao, J. K., 2010. A hybrid LDA and genetic algorithm for gene
820 selection and classification of microarray data. *Neurocomputing*. 73 (13-15), 2375–2383.

821 [31] Rammal, A., Perrin, E., Vrabie, V., Assaf, R., Fenniri, H. 2017. Selection of discriminant
822 mid-infrared wavenumbers by combining a naïve Bayesian classifier and a genetic algorithm:
823 Application to the evaluation of lignocellulosic biomass biodegradation. *Mathematical*
824 *Biosciences*. 289, 153-161.

825 [32] Punia, M., Joshi, P.K., Porwal, M.C. 2010. Decision tree classification of land use land

826 cover for Delhi, India using IRS-P6 AWiFS data. *Expert Systems with Applications*. 38 (5),
827 5577–5583.

828 [33] Safavian, S.R., Landgrebe, D. 1991. A survey of decision tree classifier methodology.
829 *IEEE Transactions on Systems, Man, and Cybernetics*. 21(3), 660-674.

830 [34] Liu, M., Wang, M., Wang, D. 2013. Comparison of random forest, support vector machine
831 and back propagation neural network for electronic tongue data classification: Application to
832 the recognition of orange beverage and Chinese vinegar. *Sensors and Actuators B: Chemical*.
833 177, 970-980

834 [35] Kotsiantis, S.B. 2007. Supervised Machine Learning: A Review of Classification
835 *Informatica*. 31, 249-268.

836 [36] Villuendas-Rey, Y., Yáñez-Márquez, C., Antón-Vargas, J., López-Yáñez, I. 2019. An
837 Extension of the Gamma Associative Classifier for Dealing with Hybrid Data. *Access IEEE*, 7,
838 64198-64205.

839 [37] Villuendas-Rey, Y., Rey-Benguría, C., Ferreira-Santiago, Á., Camacho-Nieto, O., Yáñez-
840 Márquez, C. 2017. The naïve associative classifier (NAC): A novel, simple, transparent, and
841 accurate classification model evaluated on financial data. *Neurocomputing*. 265, 105–115

842 [38] Kotsiantis, S.B., 2007. Supervised Machine Learning: A Review of Classification
843 *Informatica*. 31, 249-268.

844 [39] Rish I. 2001. An Empirical Study of the naïve Bayes Classifier, *Artificial Intelligence*
845 *Journal IJCAI*. 3: 41-46.

846 [40] Fukunaga, K. 1990. *Introduction to Statistical Pattern Recognition*. Academic Press,
847 second edition.

848 [41] Garber, F.D., Djouadi, A. 1988. Bounds on the Bayes classification error based on
849 pairwise risk functions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 10,
850 281–288.

- 851 [42] Park, C.H., Lee, M. 2008. On applying linear discriminant analysis for multi-labeled
852 problems. *Pattern Recognition Letters*. 29 (7), 878–887.
- 853 [43] Sharma, R., Ghosh, A., Joshi, P.K. 2013. Decision tree approach for classification of
854 remotely sensed satellite data using open-source support. *Journal of Earth System Science*, 122,
855 1237–1247.
- 856 [44] Bressan, G.M., De Azevedo, B.C., Ap, E., Lizzi, S. 2017. A Decision Tree Approach for
857 the Musical Genres Classification. *Applied Mathematics & Information Sciences*. 11 (6), 1703-
858 1713.
- 859 [45] Rokach, L., Maimon, O. 2008. *Data mining with decision trees: Theory and applications*,
860 World Scientific. 69, doi.org/10.1142/6604
- 861 [46] Tsang, S., Kao, B., Yip, K.Y., Ho, W.S., Lee, S.D. 2009. Decision trees for uncertain data.
862 *IEEE 25th International Conference on Data Engineering*, 441-444.
- 863 [47] Luts, J., Ojeda, F., Plas, R., Van De Moor, B., De Huffel, S. 2010. A tutorial on support
864 vector machine-based methods for classification problems in chemometrics, *Analytica Chimica*
865 *Acta*. 665, 129–145.
- 866 [48] Cortes, C., Vapnik, V. 1995. Support-vector network. *Machine Learning*. 20, 273–297,
867 doi.org/10.1007/BF00994018
- 868 [49] Jain, A.K., Mao, J., Mohuiddin, K. 1996. Artificial Neural Networks: A Tutorial. *IEEE*
869 *Computer*. 29(3), 31-44.
- 870 [50] Simon, H. 1999. *Neural Networks and Learning Machines*. Third Edition. McMaster
871 University Hamilton, Ontario, Canada. Rev. ed of: *Neural networks*. 2nd ed.
- 872 [51] Sildir, H., Aydin, E., Kavzoglu, T. 2020. Design of Feedforward Neural Networks in the
873 Classification of Hyperspectral Imagery Using Superstructural Optimization. *Remote Sensing*.
874 12, 956, doi:10.3390/rs12060956
- 875 [52] ItzamA, L.Y., Amadeo, J.A.C., Oscar, C.N., Cornelio Y.M. 2011. Pollutants Time-Series

876 Prediction using the Gamma Classifier, *International Journal of Computational Intelligence*
877 *Systems*. 4(4), 680-711, DOI: 10.1080/18756891.2011.9727822

878 [53] Rangel-Diaz-de-la-Vega, A., Villuendas-Rey, Y., Yanez-Marquez, C., Camacho-Nieto,
879 O., Lopez-Yanez, I. 2020. Impact of Imbalanced Datasets Preprocessing in the Performance of
880 Associative Classifiers. *Applied Sciences*. 10(8), 2779; doi:10.3390/app10082779

881 [54] Rangel-Diaz De La Vega, A., Villuendas-Rey, Y., Yanez-Marquez, C., Camacho-Nieto,
882 O. 2020. The Naïve Associative Classifier With Epsilon Disambiguation. *IEEE Access*. 8,
883 51862-51870.

884 [55] Villuendas-Rey, Y., Hernandez-Castano, J., Camacho-Nieto, O., Yanez-Marquez, C.,
885 Lopez-Yanez, I. 2019. NACOD: A Naïve Associative Classifier for Online Data. *IEEE Access*,
886 7, 117761-117767.

887 [56] Sun, Z., Zhang, H., Yang, Y., Wan, H., Wang, Y. 2020. Impacts of Geographic Factors
888 and Population Density on the COVID-19 Spreading under the Lockdown Policies of China.
889 *Science of the Total Environment*. 746(666).

890 [57] Dowd, J. B., Andriano, L., Brazel, D., Rotondi, V., Block, P., Ding, X., Liu, Y., Mills, M.
891 2020. Demographic Science Aids in Understanding the Spread and Fatality Rates of COVID-
892 19. *Proceedings of the National Academy of Sciences of the United States of America* 117(18),
893 9696–98

894 [58] Emanuel, E.J., Persad, G., Upshur, R., Thome, B., Parker, M., Glickman, A., Zhang, C.,
895 Boyle, C., Smith, M., Phillips, J.P. 2020. Fair Allocation of Scarce Medical Resources in the
896 Time of Covid-19. *N Engl J Med*. 382(21), 2049-2055. doi: 10.1056/NEJMsb2005114.

897 [59] Lai, C.C., Wang, C.Y., Wang, Y.H., Hsueh, S.C., Ko, W.C., Hsueh, P.R. 2020. Global
898 epidemiology of coronavirus disease 2019 (COVID-19): disease incidence, daily cumulative
899 index, mortality, and their association with country healthcare resources and economic status.
900 *Int J Antimicrob Agents*. 55(4), 105946. doi: 10.1016/j.ijantimicag.2020.105946.

- 901 [60] Safavi, P.E., Rahimian, K., Doustmohammadi, A., Safari dastjerdei, M., Rasouli, A.,
902 Zahiri, J. 2021. A prediction model for COVID-19 prevalence based on demographic and
903 healthcare parameters in Iran. <https://doi.org/10.1101/2021.01.27.21250551>
- 904 [61] Payam, S., Bordbar, N., Ghanbarzadegan, A., Najibi, M., Bastani, P. 2020. Ranking of
905 Iranian Provinces Based on Healthcare Infrastructures: Before and after Implementation of
906 Health Transformation Plan. *Cost Effectiveness and Resource Allocation* 18(1).
- 907 [63] Nie, P., Ding, L., Chen, Z. 2021. Income-related health inequality among Chinese adults
908 during the COVID-19 pandemic: evidence based on an online survey. *Int J Equity Health* 20,
909 106. <https://doi.org/10.1186/s12939-021-01448-9>
- 910 [63] Buja, A., Paganini, M., Cocchio, S., Scioni, M., Rebba, V., Baldo, V. 2020. Demographic
911 and socio-economic factors, and healthcare resource indicators associated with the rapid spread
912 of COVID-19 in Northern Italy: An ecological study. *PLoS ONE*. 15(12),
913 <https://doi.org/10.1371/journal.pone.0244535>
- 914 [64] National Institute of Statistics (ISTAT). Main database. [cited 2020 April 13]. Available
915 from: <http://dati.istat.it/Index.aspx>
- 916 [65] National Institute of Statistics (ISTAT). Health for All database. [cited 2020 April 13].
917 Available from: <https://www.istat.it/it/archivio/14562>.
- 918 [66] Bizri, A.R., Khachfe, H.H., Fares, M.Y., Musharrafieh, U. 2021. COVID-19 Pandemic:
919 An Insult Over Injury for Lebanon. *J Community Health*. 46(3), 487-493.
- 920 [67] Sydow, J., Windeler, A., Muller-Seitz, G. 2012. Path Constitution Analysis: A
921 Methodology for Understanding Path Dependence and Path Creation. *Bus Res*, 5, 155–176.
- 922 [68] Rammal, A., Perrin, E., Vrabie, V., Bertrand, I., Chabbert, B. 2017. Classification of
923 lignocellulosic biomass by weighted-covariance factor fuzzy C-means clustering of mid-
924 infrared and near-infrared spectra, *Journal of Chemometrics*. 31(2), 1–10.
- 925 [69] Everson, R., Fieldsend, J. 2006. Multi-class ROC analysis from a multi-objective

926 optimization perspective. Pattern Recognition Letters 27(8), 918-927.

927

928

929

930

931

932

933

934

935 **Supplementary Information**

936 Table 3: Correlation between the degrees of possibility and the predictors

Degrees of possibility		
	Pearson Correlation	Sig (2-tailed)
Technology used	-0.052	0.442
Special event	0.061	0.362
Procedure	0.027	0.69
Experience	-0.363	000
Density	-0.22	0.001
Family Nb	0.152	0.022
Temperature	-0.009	0.893
Lockdown	-0.232	000
Medical Resources	-0.224	0.001
Economy Resources	-0.081	0.23

937

938