# Initial Stage COVID-19 Detection System Based on Patients' Symptoms and Chest X-Ray Images

Muhammad Attaullah, Mushtaq Ali, Maram Fahhad Almufareh, Muneer Ahmad, Lal Hussain, Nz Jhanjhi & Mamoona Humayun

Published online: 18 Apr 2022.

Submit your article to this journal ↗

Article views: 1581

View related articles ↗

View Crossmark data ↗

Citing articles: 19 View citing articles ↗

# Initial Stage COVID-19 Detection System Based on Patients' Symptoms and Chest X-Ray Images

Muhammad Attaullah[a], Mushtaq Ali[a], Maram Fahhad Almufareh[b], Muneer Ahmad[c], Lal Hussain [iD][d], Nz Jhanjhi[e], and Mamoona Humayun [iD][b]

[a]Department of Information Technology, Hazara University, Mansehra, Pakistan; [b]Department of Information Systems, College of Computer and Information Sciences, Jouf University, Ksa; [c]School of Electrical Engineering and Computer Science (SEECS), National University of Sciences and Technology (NUST), Islamabad, Pakistan; [d]Department of Computer Science IT, University of Azad Jammu and Kashmir, Muzaffarabade, Pakistan; [e]School of Computer Science, SCS, Taylor's University, Subang Jaya, Selangor, Malaysia

**ABSTRACT**
The accurate diagnosis of the initial stage COVID-19 is necessary for minimizing its spreading rate. The physicians most often recommend RT-PCR tests; this is invasive, time-consuming, and ineffective in reducing the spread rate of COVID-19. However, this can be minimized by using noninvasive and fast machine learning methods trained either on labeled patients' symptoms or medical images. The machine learning methods trained on labeled patients' symptoms cannot differentiate between different types of pneumonias like COVID-19, viral pneumonia, and bacterial pneumonia because of similar symptoms, i.e., cough, fever, headache, sore throat, and shortness of breath. The machine learning methods trained on labeled patients' medical images have the potential to overcome the limitation of the symptom-based method; however, these methods are incapable of detecting COVID-19 in the initial stage because the infection of COVID-19 takes 3 to 12 days to appear. This research proposes a COVID-19 detection system with the potential to detect COVID-19 in the initial stage by employing deep learning models over patients' symptoms and chest X-Ray images. The proposed system obtained average accuracy 78.88%, specificity 94%, and sensitivity 77% on a testing dataset containing 800 patients' X-Ray images and 800 patients' symptoms, better than existing COVID-19 detection methods.

## Introduction

The virus that causes respiratory diseases is called corona virus (institute, N. c 2020). The corona virus is a large family of viruses like SARS-COV-1, MERS-COV, and SARS-COV-2 (Chorba 2020; F. Li 2016). The first case of SARS-COV-1 appeared on 16 November 2002 in China. This virus infected more than 8,000 people belonging to 29 countries, and around 774 were died (Gov 2021). The MER-COV appeared in Saudi Arabia in 2012. This virus

**CONTACT** NZ Jhanjhi ✉ noorzaman.jhanjhi@taylors.edu.my School of Computer Science and Engineering, Sce Taylor's University, Subang Jaya 47500, Selangor, Malaysia

infected both human and animals. The SARSCOV-2 virus which is also known as COVID-19 (Liu, Kuo, and Shih 2020) infected patients in China. The public health emergency regarding COVID-19 was declared on 30 January 2020 (Organization, W. H 2020). According to WHO statistics taken on May, 2021, the average number of COVID-19 cases and deaths happening per day are 372,522, and 7795, respectively (Gov 2021). These statistics frightens the whole world and the researchers started proposing methods for controlling the infection and death rate of COVID-19. Two methods, namely medical lab based and machine learning based, are usually used for diagnosing COVID-19.

The medical lab methods are further divided into three types namely, antibody blood based, antigen swab based and Reverse Transcriptase-Polymerase Chain Reaction (RT-PCR) based.

Most of the experts do not recommend antibody blood test for diagnosis of COVID-19 because it is misleading due to false negative results. The antigen swab test checks the specific type of protein in the liquid taken from nose or throat using antigen kits but due to low sensitivity, most often it produces false negative results. The accuracy rate of RT-PCR test is comparatively better, therefore it is considered as a gold standard for COVID-19 testing. However, this test is costly, invasive, and time consuming (Gov 2021). The researchers introduced noninvasive, cheaper, and fast machine learning-based COVID-19 detection methods trained either on labeled patients' symptoms, or labeled patients' medical images. The symptoms-based machine learning techniques (Ahamad et al. 2020; Zoabi, Deri-Rozov, and Shomron 2021) have the potential to classify the case of a patient as COVID-19 or normal using patient's symptoms. However, they are incapable to differentiate between different kinds of pneumonia, and often generate false results. This is because of the reason that these three diseases hold similar symptoms, i.e., cough, fever, headache, sore throat, and shortness of breath (Gov 2021; Medlineplus 2021). The machine learning methods trained on labeled medical images have been introduced in (Hussain et al. 2020; Zheng et al. 2020) for classifying the different kinds of pneumonia. The classification results of these methods are satisfactory, but they are incapable to detect COVID-19 in initial stage because COVID-19 patients have no chest infection in initial stage (Wong et al. 2020). Detecting COVID-19 in initial stage is necessary because it transfers from person to person very quickly (Chan et al. 2020). The initial stage COVID-19 detection is one of the ways through which the COVID-19 spreading rate can be minimized. In this research work, an initial stage COVID-19 detection system is proposed. In this system, we trained the regression model on publicly available COVID-19 symptoms dataset (Gov 2021) and CNN model on the chest X-Ray images dataset formed from the combination of two datasets (Bganglia 2020; Mooney 2018), respectively. After this, the

decision tree model is trained on the results obtained from the previously two trained models using domain knowledge. The contributions of our research work are described as follows:

(1) We identified that none of the existing COVID-19 detection methods are based both on patients' symptoms and medical images.
(2) It has also been identified that the available COVID-19 detection methods lack the potential to predict COVID-19 in initial stage.
(3) We proposed COVID-19 detection system based on both patients' symptoms and medical images capable of detecting COVID-19 in initial stage.

The other sections of the paper are organized as follows. The detail overview of literature is provided in section 2. The section 3 and 4 present the training dataset and working of our proposed system. The details regarding the testing dataset, evaluation metrics, and results' comparison are provided in section 5. The conclusion and future work are specified in section 6, followed by references.

## Related Work

In this section, the available machine learning approaches used for COVID-19 detection are presented. First, we describe the working details and limitation(s) of the symptoms-based machine learning techniques used for COVID-19 detection, and then present the same regarding the image-based machine learning techniques.

In (Ahamad et al. 2020) the authors proposed a COVID-19 detection method based on five symptoms (i.e., cough, fever, headache, sore throat, and shortness of breath) using machine learning algorithms. They trained the gradient-boosting machine model on publicly available dataset (Gov 2021) containing 99,232 tested patients' symptoms along with their respective labels. This dataset contains 8,393 confirmed COVID-19 cases, and the rest are normal. During testing, the patient needs to provide information regarding five symptoms to the model. Additionally, the patient also needs to provide information about his/her sex, age, and contact with confirmed COVID-19 patient. The proposed model predicts the result based on provided information. The prediction accuracy of the proposed model in terms of sensitivity, and specificity were reported as 87.30%, and 71.98%, respectively. The prediction accuracy of this model is comparatively better. However, it is incapable to differentiate between different forms of pneumonia.

In (Zoabi, Deri-Rozov, and Shomron 2021) the authors proposed a method to rank the symptoms on the basis of age. They trained different models, like XGBoost, GBM, SVM, Random Forest, and decision tree on the

publicly available dataset (Chunxiaozheng 2020) to find out the most dominant symptoms of COVID-19 patient. This work assists physicians in examining the specific set of symptoms instead of looking to all symptoms. However, this research work possesses no potential to classify the pneumonias.

In (Hussain et al. 2020) the authors used deep learning models to classify different kinds of pneumonias based on X-Ray image of lung. They trained five models, i.e., XGB-L, XGB-Tree, CART, KNN, and Naïve Bayes on publicly available dataset (Hussain et al. 2020). The prediction results of XGB-L are comparatively better than the other four models. However, it fails to detect COVID-19 in initial stage.

In (Horry et al. 2020) the authors proposed a method to detect the COVID-19 from medical chest images. They trained the deep learning model named as VGG-19 model on the publicly available dataset (Cohen et al. 2020) using transfer learning technique. The VGG-19 model got 84%, 86%, and 100% accuracy on chest CT-Scan, chest X-Ray, and chest Ultrasound images, respectively. However, this model is unable to classify different kinds of pneumonia. Furthermore, its COVID-19 detection rate is not satisfactory in initial stage COVID-19.

In (Ohata et al. 2020) the authors proposed the COVID-19 detection method to classify three different classes (COVID-19, Common pneumonia, Normal) from X-Ray images using transfer learning. They trained the CNN model on publicly available dataset (Kaggle 2018) which contains 257 patient's images. This proposed method got 98.5% accuracy on testing dataset. However, the proposed system has no potential to detect the COVID-19 in initial stage and classify different type of pneumonia diseases.

Deep neural network for detection and diagnosis method of COVID-19 from chest x-ray images have been proposed in (Khan, Shah, and Bhat 2020). The CNN, CoroNet model have trained on publicly available dataset (Bganglia 2020) containing 1251 images to classify into three different classes: COVID-19, non-COVID-19 Pneumonia, and Normal; they got 89% accuracy on testing dataset. However, the proposed system cannot detect the COVID-19 in initial stage.

In (Brunese et al. 2020) the authors proposed CNN and VGG16 model for diagnosis of COVID-19, non-COVID-19 Pneumonia, Normal, and Highlight the infected area in the X-Ray image. They trained models on publicly dataset (Dataset and I. O. H. C. X.-r 2017; Talo 2020) which contains 6523 X-Ray images of different patients. They got 97% accuracy on testing dataset this system is comparatively better to detect the COVID-19 from X-Ray image and highlighting the infected area in X-Ray image, However, this system cannot detect the COVID-19 in initial stage and also cannot differential different type of non-COVID-19 pneumonias.

In (Pham 2021) the authors proposed Classification of COVID-19 from chest X-rays images with deep learning method to improve the accuracy of existing system and fine tuning of pretrained convolutional neural networks. In the proposed system CNN, AlexNet, GoogleNet, and SqueezeNet have been trained on publicly available datasets (Bganglia 2020; Rahman 2021) which contains 3666 X-Ray images of different patients. They got 99% accuracy on testing dataset this is better to investigate and fine tun the pretrained convolutional neural networks. However, this system cannot detect the COVID-19 in initial stage and further more they have no potential to differentiate different king of pneumonia.

In (Ozturk et al. 2020) the authors proposed an automated COVID-19 detection system using deep neural networks like CNN and DarkCovidNet. Its accuracy was tested on publicly available dataset (Bganglia 2020) containing 1750 X-Ray images and got an accuracy score of 98.08% for binary classification. However, it fails to detect initial stage COVID-19. In addition, it has no potential to differentiate between different types of pneumonia. The CT-Scan base COVID-19 detection systems have been proposed in (Li et al. 2020; Zheng et al. 2020). However, these all-classification methods are incapable to detect the COVID-19 in initial stage. Furthermore, they are incorrect due to the misclassification of COVID-19 images with other types of lunge's disease (viral pneumonia, bacterial pneumonia) (Rehman et al. 2021).

In (Jain et al. 2021) the authors proposed deep learning-based technique to detect COVID-19 using chest X-Rays images. They trained different deep learning CNN models, such as, inception V3, Xception and ResNeXt models on publicly available Kaggle dataset (Patel) containing 6432 X-Ray images in which 5467 were used for training and 965 for validation. In the result comparison analysis, the Xception model got high accuracy 97.97% as compared to other models. However, this work only focuses on comparing the performance of different classifiers in detecting COVID-19 infection.

In (Marateb et al. 2021) the authors proposed reliable automatic computer aided system to classify the COVID-19 from non-COVID-19 pneumonia using demographics, symptoms and blood tests features. They train ensemble classifier like XGBOOST, SVM, and Gradiant Boosting classifiers on three publicly available dataset (Einstein 2020; Gov 2021; Sami et al. 2020). However, this proposed system is not able to further classify the non-COVID-19 pneumonia.

In (Canas et al. 2021) the authors proposed the hierarchical gaussian process model to detect the early signs of COVID-19 infection in UK. They trained the different models like hierarchical gaussian process, logistic regression and NHS algorithm on the patient self-reported dataset containing 182,991 training and 15,049 validation patients' symptoms. The hierarchical gaussian process model got high accuracy AUC 0.80% as compare to other two models. The aim of this proposed system is to estimate the probability of an

individual being infected with COVID-19 on the bases of early self-reported symptoms to enable timely self-isolation and urgent testing. This work is useful for early test and isolation of patients using only symptom. However, it does not provide information about different kinds of pneumonias.

In (Koushik, Bhattacharjee, and Hemalatha 2021) authors proposed the hybrid modeling to improve the accuracy of symptom based COVID-19 detection system. They trained the MaxVoting ensemble model, comprising gradient boosting and random forest algorithms on publicly available dataset (Gov 2021) containing 112345 different patents' symptoms along with their PCR result in which 154k and 51k were used for training and testing, respectively. They got 90% accuracy which is promising. However, it is incapable to differentiate between different forms of pneumonia.

Our proposed system differs from the aforementioned COVID-19 detection methods because it detects initial stage COVID-19 by utilizing the symptoms and medical images of the patients.
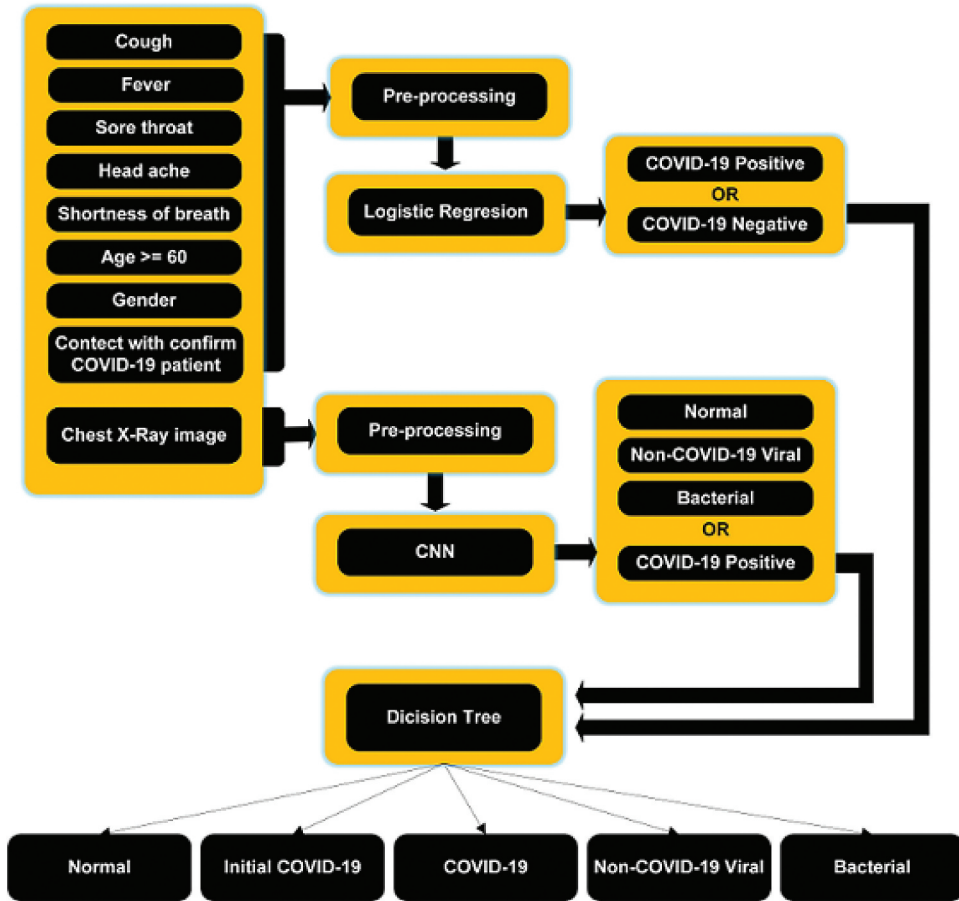
## Training Datasets

We have collected chest X-Ray images of COVID-19 patients from COVID-19 Radiography dataset (Rehman et al. 2021) and non-COVID Viral, Bacterial, and Normal patients from Chest X-Ray images (Pneumonia) (Bganglia 2020) on April 22, 2021. The COVID-19 Radiography dataset (Mooney 2018) and Chest X-Ray images (Pneumonia) (Bganglia 2020) contains 3616, and 5863 Chest X-ray images, respectively. We obtained 1000 chest X-Ray images from dataset (Mooney 2018) and 1000 chest X-Ray images from each class of dataset (Bganglia 2020) randomly and formed a new dataset containing 4000 chest X-Ray images belonging to four classes of different pneumonia.

In addition to chest X-Ray images, the symptoms of COVID-19, and non-COVID-19 patients are collected from the Israeli Ministry of Health (Gov 2021) containing 45 hundred thousand records of patients in Hebrew language. Among these 45 hundred thousand patients, 4 hundred thousand are COVID-19 positive, while the rest are COVID-19 negative cases. The English version of this dataset containing 27 hundred thousand patients' records is publicly available (Nshomron 2020). In order to improve the prediction accuracy, we converted the whole dataset into English form.

The dataset is decomposed into training and testing datasets. The training dataset contains 80% images of the new dataset and 600, 000 patients' symptoms (300,000 COVID- 19 + v, and 300,000 COVID-19 -v) while the testing dataset contains 20% images of the new dataset and 800 patients' symptoms taken randomly from the dataset (Gov 2021). According to the aforementioned statistics, every class in the training dataset contains 800 training images and every class in the testing dataset contains 200 testing images.

**Figure 1.** Architecture of our proposed system.

## The Proposed System

Our proposed system performs pre-processing operations on the X-Ray images, and symptoms datasets first. Next, it trains the logistic regression model and CNN model on the pre- processed datasets. The decision tree model is then trained on the labeled results of the previous two trained models. The architecture of our proposed system is shown in Figure 1.

### *Pre-Processing*

The images in the publicly available dataset are converted into numpy and their dimensions are changed into 150*150 pixels. The symptoms dataset contains nine symptoms, i.e., cough, fever, sore throat, head ache, shortness of breath, age above 60, gender, indication, and result. Four out of these nine symptoms, i.e., age above 60, gender, indication and result contain English alphabets, i.e., Yes/No, male/female, contact with confirm or not, COVID-19

positive or COVID- 19 negative, respectively. The aforementioned four symptoms' data are converted into numeric values, i.e., binary values because the logistic regression model works on numeric values. Next, the duplicate rows and null values in the symptom dataset are dropped and last to balance the minority and majority imbalance classes resampling technique is applied to balance it.

### Logistic Regression Model

We selected the logistic regression model for predicting COVID-19 based on symptoms because of its outstanding results generated in various application domains like weather forecasting, rainfall prediction, and e-mail span detection (Imon, Roy, and Bhattacharjee 2012; Sánchez, Ortega, and Marcos 2001). We trained the logistic regression model on the publicly available dataset. The entire dataset is decomposed into training and testing dataset contains 80% and 20% data of the whole dataset, respectively. The values of the first eight attributes in the training dataset are stored in the X_train and the value of the result attribute is stored in the Y_train. Similarly, the values of the first eight attributes in the testing dataset are stored in the X_test and the value of the result attribute is store in the Y _test. The logistic regression linear model is then trained by passing X_train and Y _train data to it. The trained model is then stored in the pickle file.

### The CNN Model

The Yann Le Cun introduced CNNs in 1988 to simulate the visual perception of human brain (Maladkar 2018). In recent past, the researchers used CNN for medical images classification, and its results were found promising as compare to other algorithms (Krizhevsky, Sutskever, and Hinton 2012; Sun et al., 2016). In this work, we used CNN model to classify infection type in the chest X-Ray image, i.e., COVID-19, non-COVID viral, bacterial pneumonia, or normal.

The data augmentation is performed during CNN training which takes 30 training images in a batch and generates 30 alternative images for each image in the batch by employing

random transformations like zooming, rotation, translation on them. The CNN is then trained on the transformed images instead of original ones. Using the transformed images instead of original ones improves the generalization of the network. The 32 filters with dimension 3 by 3 each are learnt by convolving them on each image in the dataset in the convolutional layer. The activation function, i.e., ReLu is applied on each of the 32 feature maps to transform them into abstract representation. The dimensions of the feature maps are then reduced by applying 2 by 2 max pooling layer. The 64 filters each with 3 by 3 size are then learnt by convolving them on the reduced feature maps obtained

from the previous max pooling layer in the convolutional layer. The activation function, i.e., ReLu is applied on each of the 64 feature maps to transform them into abstract representation. The size of the feature maps obtained from the previous layer are then reduced further by applying 2 by 2 max pooling layer The 128 filters with 3 by 3 dimension each are learnt by convolving them on the reduced feature maps obtained from the previous max pooling layer in the convolutional layer. The activation function, i.e., ReLu is applied on each of the 128 feature maps to transform them into abstract representation. The flatten layer converts each of the feature map into one- dimensional vector because the following dense layer works on only one-dimensional vector data. The dense layer takes the one-dimensional vector obtained from the previous flatten layer and generates 128 outputs. The Relu activation function is applied on the resultant output to transform them into a form acceptable to the next layer, that is Dropout layer. The dropout layer changes 50% of the values obtained from the previous layer to 0 to minimize the problem of overfitting. Finally, the dense layer is applied with softmax activation function which return the probabilities of four classes.

### *Decision Tree*

The logistic regression model generates binary output 0 (COVID-19 negative) or 1 (COVID-19 Positive) based on provided symptoms. While the CNN model produces one of the four outputs i.e., 0001 (Normal), 0100 (COVID-19), 0010 (Non COVID-19 viral), or 1000 (bacterial pneumonia) based on provided X-Ray image. A dataset based on the results generated by logistic regression model and CNN model using domain knowledge of radiologists serving in Ayub Medical Hospital, Abbottabad, Pakistan is created for the training purpose of decision tree model. This dataset, hereafter called knowledge dataset contains three attributes, i.e., Logistic Regression Predictions (LR_P), CNN Pre-dictions (CNN_P), and Disease (D). The LG_P attribute contains all possible predictions of logistic regression model, the CNN_P attribute contains all possible predictions of CNN model and the Disease attribute contains the labels.

The labels of Disease attribute are obtained based on domain knowledge. The Diseases attribute contains five types of labels, i.e., 0,1,2,3 and 4 referring to Normal, COVID-19, Non COVID-19 viral, bacterial pneumonia and initial stage COVID-19, respectively. The process used for labeling the Disease attribute of the knowledge dataset is described as follows.

(1) The radiologists often declare the report of a patient as COVID-19 negative if both COVID-19 infection in the chest X-Ray and COVID-19 symptoms are not found. Based on this observation the value of attribute Disease is set as 0 for all such cases.

(2) A COVID-19 positive report is issued to a patient with COVID-19 symptoms and COVID-19 infection in X-Ray. Adopting the same scenario, the value of attribute disease is set as 1 in all such cases.

(3) In case when the patient has no symptoms but COVID- 19 infection is found in his chest X-Ray then the radiologists declare him as COVID-19 patient.

(4) When a patient has no COVID-19 symptoms and non-COVID-19 viral infection is found in his chest X-Ray then the radiologists issue him a report of non-COVID-19 viral.

(5) For a patient having no COVID-19 symptoms but has bacterial infection in chest X-Ray the physician reports that patient has bacterial Pneumonia.

(6) Similarly, if the patient has COVID-19 symptoms but his chest X-Ray is normal then in that case the physician suggest that the patient has initial stage COVID 19.

(7) In cases where the patient has COVID-19 symptoms but his chest X-Ray contains non-COVID-19 viral infection then the radiologist affirms that the patient has non-COVID-19 viral pneumonia.

(8) When the patient has COVID-19 symptoms and his chest X-Ray reflects bacterial infection then bacterial pneumonia is reported.

Based on the above labeling process the Disease attribute is labeled with five classes denoted by 0, 1,2, 3, 4 as shown in Table 1.

The knowledge dataset shown in Table 1 is then used for training the decision tree model. Before the training started, the Disease attribute values are separated from the knowledge dataset and stored in Y_train variable. Similarly, the values of other two attributes, i.e., LG_P and CNN _P are stored in X_train variable. The Decision Tree Classifier is then loaded and the values in X_train and Y_train are passed to its fit function. The trained decision model is then saved using the joblib library.

The decision tree needs the values of three inputs (i.e., LG_P, CNN _P, and Disease) as shown in Table 1 during its training process. However, in case of testing process only the values of two inputs (i.e., LG_P, and CNN _P) need to

**Table 1.** Labeling disease attribute based on domain knowledge.

| LG_P | CNN_P | Disease |
|------|-------|---------|
| 0 | 0001 | 0 |
| 1 | 0100 | 1 |
| 0 | 0100 | 1 |
| 0 | 0010 | 2 |
| 0 | 1000 | 3 |
| 1 | 0001 | 4 |
| 1 | 0010 | 2 |
| 1 | 1000 | 3 |

be provided to the learned decision tree model. The architecture shown in Figure 1 reflect the testing process, that is why Disease attribute is not provided to the decision tree.

The algorithm of our proposed system is described as follows.

Algorithm (patient symptoms, X-Rays)

**Step 1**: Pre-process patient symptoms

   (i) Convert Hebrew language symptoms dataset (S) to English
  (ii) Drop duplicate and null values
 (iii) Convert all features of S into numeric values
 (iv) Store the features and labels in S_feature and S_labels variable respectively
  (v) Split Using train test split library convert to X_train, X_test, Y_train, Y_test = train_test_split (S_feature, S_labels, test_size = 0.2, random_state = 51)

**Step 2**: Pre-process patient X-Ray image

   (i) Reshape X-Rays Dataset images (X) to (150, 150, 3)
  (ii) Random rotation range = 10
 (iii) Horizontal flip = True
 (iv) Zoom Range = 0.2
  (v) Width shift range = 0.2
 (vi) Hight shift range = 0.2
(vii) Validation split = 0.2

**Step 3**: Logistic Regression model training on S

   (i) import the logistic regression function from the sklearn library
  (ii) Fit the X_train and Y_train in logistic regression function and train the model
 (iii) Store the trained model

**Step 4**: CNN Model feature extraction and training on X

   (i) model.add (Convolutional2D 32 filters size 3 by 3 input image size = (150, 150, 3))
  (ii) model.add (Activation('relu'))
 (iii) model.add (MaxPooling2D(pool_size = (2,2)))
 (iv) model.add (Conv2D(64, (3,3)))
  (v) model.add (Activation('relu'))
 (vi) model.add (MaxPooling2D(pool_size = (2,2)))
(vii) model.add (Conv2D(128, (3,3)))

   (viii) model.add (Activation('relu'))
    (ix) model.add (MaxPooling2D(pool_size = (2,2)))
     (x) model.add (Flatten())
    (xi) model.add (Dense(128))
   (xii) model.add (Activation('relu'))
 (xiii) model.add(Dropout(0.5))
 (xiv) model.add(Dense(4))
  (xv) model.add(Activation('sigmoid'))
 (xvi) Train the model with 50 epochs
(xvii) stored the model

**Step 5**: Decision tree training

   (i) load the knowledge dataset
  (ii) store the LG_P, CNN_P Attributes in X_train and Diseases Attribute in Y_train
 (iii) import the Decision tree classifier from sklearn library
 (iv) Fit the X_train and Y_train in Decision Tree Classifier function and train the model
  (v) Store the trained model

**Step 6**: Test the whole system combine three models

   (i) S = patient eight symptom
  (ii) X = Patient X-Ray
 (iii) Load the Logistic Regression pre-trained model
 (iv) LG_P = logistic Regression model predict S
  (v) Load the CNN trained model
 (vi) CNN_P = CNN model predict X
 (vii) Load the Decision tree pre-trained model
(viii) Output = Decision tree model predict LG_P and CNN_P

The parameters used in the logistic regression model, CNN, and decision tree model are shown in Table 2.

## Experimental Results

This section provides detail information regarding the test setup, testing dataset, evaluation metrics, and results' comparison.

 (1) Test Setup and Tuning

**Table 2.** Parameters used in the components of our proposed system.

| Logistic Regression Model Parameters | CNN Parameters | Decision tree parameters |
|---|---|---|
| Training dataset, | Training dataset images, | Training sample, |
| Labels, | steps_per_epoch = 3200/30, | label, |
| C = 1.0, | epochs = 50, | ccp_alpha = 0.0, |
| class_weight = None, | callbacks = [history], | class_weight = None, |
| dual = False, | validation_data = validation_generator, | criterion = 'gini,' |
| fit_intercept = True, | validation_steps = 800/30 | max_depth = None, |
| intercept_scaling = 1, | | max_features = None, |
| max_iter = 100, | | max_leaf_nodes = None, |
| multi_class = 'ovr,' | | min_impurity_decrease = 0.0, |
| n_jobs = 1, | | min_impurity_split = None, |
| penalty = 'l2,' | | min_samples_leaf = 1, |
| random_state = None, | | min_samples_split = 2, |
| solver = 'liblinear,' | | min_weight_fraction_leaf = 0.0, |
| tol = 0.0001, | | presort = 'deprecated,' |
| verbose = 0, | | random_state = None, |
| warm_start = False | | splitter = 'best' |

The online cloud based Jupyter Notebooks provided by google named as Google Colaboratory was used for the development and testing of our proposed system. It provides free GPU facility. It also provides either Python 2 or 3 runtimes pre-configured with the essential machine learning libraries like TensorFlow, Matplotlib, and Keras. The google colaboratory is composed of an Intel Xeon processor with two cores @2.3 GHz and 13 GB RAM. It is also equipped with a NVIDIA Tesla K80 (GK210 chipset), 12 GB RAM, 2496 CUDA cores @560 MHz.

## Testing Dataset

Due to non-availability of patients' chest X-Rays along their respective symptoms, we labeled 800 patients' chest X-ray images and symptoms based on domain knowledge described in section 4–4 and use them as testing dataset. The testing dataset is composed of five classes, i.e., Bacterial, COVID-19, non-COVID-19 viral, initial stage COVID-19 and normal, containing 200, 200, 180, 90, and 130 images along with their respective symptoms.

## Evaluation Metrics

The true positive rate, false positive rate, accuracy, sensitivity, and specificity of our proposed system were computed for each class in the testing dataset using Eq. 1, Eq. 2, Eq. 3, Eq. 4, and Eq. 5, respectively.

$$TPR = \frac{TP}{TP + FP} \tag{1}$$

$$FPR = \frac{FP}{TP + FP} \tag{2}$$

$$Accuracy = \frac{TPR}{TPR + FPR} \tag{3}$$

$$Sensitivity = \frac{TP}{TP + FP} \tag{4}$$

$$Specificity = \frac{TN}{TN + FP} \tag{5}$$

Where TP, FP, TN, TPR, and FPR refer to numbers of true prediction, number of false prediction, number of true negative prediction, rate of correct prediction, and rate of incorrect prediction, respectively.

### *Time and Space Complexity*

The actual training and classification times of our proposed system on CPU and GPU are shown in Table 3.

The time complexity of the regression model used in our proposed system is O(nd), where n, and d refer to the total number of suspected patients, and the number of symptoms of each patient, respectively. The time complexity of the CNN model O(r*m*(ij+jk+kl)), where r, and m denote to the number of epochs, and number of chest X-Rays, respectively. And i, j,k, l represent the number of nodes in each of the four layers used in CNN. The time complexity of the decision tree is and O(p*log(p)*f), respectively. Where p, and f show the number of labeled samples, and number of features, respectively. Using the time complexities of the regression model, CNN, and decision tree the overall time complexity of our proposed system becomes as O(nd+ r*m*(ij+jk+kl)+ p*log(p)*f). The space complexity regression model, CNN and decision tree

**Table 3.** Actual training and classification times of our proposed system.

| Hardware | Model Name | Training time | Classification time |
|---|---|---|---|
| CPU Intel Xeon processor | CNN | 3 h 12 min 49s | 83.3 ms |
| with two cores | Logistic Regression | 12.2 s | 5.3 ms |
| @2.3 GHz and 13 GB | Decision tree | 5.15 ms | 4.12 ms |
| RAM | | | |
| Average | | 3 h 13 min 7.15 ms | 92.71 ms |
| GPU NVIDIA Tesla K80 | CNN | 54 min 23s | 47.3 ms |
| (GK210 chipset), 12 GB | Logistic Regression | 9.86 s | 4.87 ms |
| RAM, 2496 CUDA cores | Decision tree | 4.43 ms | 2.02 ms |
| @560 MHz | | | |
| Average | | 54 min 32s | 54.19 ms |
| | | 90.43 ms | |

are O(d), O(m), and O(nodes), respectively. As our proposed system is composed of these three models, so the space complexity of our proposed system is O(d + m + nodes).

## *Results Evaluation*

It is obvious that the hybrid system, i.e., text (symptoms) and image-based system generates better results than the ones based on either text or images alone. However, due to non-availability of symptoms and image-based COVID-19 detection systems, the results of our proposed system are compared with the symptoms-based system and image-based system as shown in Table 4.

The gradient boosting model proposed in (Ahamad et al. 2020) got 85% prediction accuracy in case of binary classification, but reduced to 34% in five class classification because of two facts. First, it cannot differentiate among COVID-19, viral pneumonia, and bacterial pneumonia. Second, it is incapable to detect initial stage COVID-19. The binary classification of COVID-19 accuracy achieved using CNN, AlexNet, GoogleNet, and SqeezNet are 99% on single modality like chest X-Ray in (Pham 2021). However, it is reduced to 39% because it is incapable to detect COVID-19 in initial stage.

The binary classification accuracy obtained by applying Resnet-101 in (Ardakani et al. 2020) and DeCoVNet in (Zheng et al. 2020) on chest CTScan images are 99%, and 90%, respectively. However, their accuracy is reduced to 39.6%, and 36% when tested on our dataset because of having no potential to detect COVID-19 in initial stage.

The VGG19 is used in (Horry et al. 2020) for binary classification using multi-modality like chest X-Ray, CT-Scan, and Ultrasound images and got 86% accuracy on chest X-Ray images. However, its accuracy is declined to 34.4% due to incapability of detecting COVID-19 in initial stage. The methods proposed in (Brunese et al. 2020; Ohata et al. 2020; Ozturk et al. 2020) applied different machine learning models on the chest X-Ray
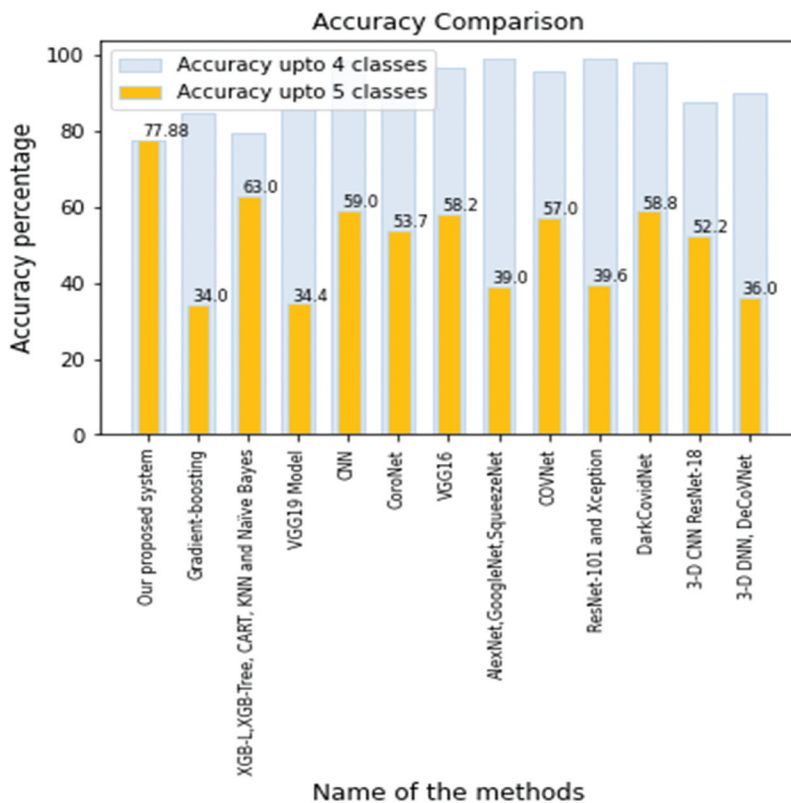
**Table 4.** Results' accuracy of existing systems.

| Method | Accuracy upto 4 classes | Accuracy of 5 classes |
|---|---|---|
| Gradient-boosting (Ahamad et al. 2020) | 85% | 34% |
| Deep learning models (Hussain et al. 2020) | 79.52% | 63% |
| VGG19 Model (Horry et al. 2020) | 86%, | 34.4 |
| CNN (Ohata et al. 2020) | 98.5% | 59% |
| CoroNet (Khan, Shah, and Bhat 2020) | 89.6% | 53.7% |
| VGG16 (Brunese et al. 2020) | 97% | 58.2% |
| AlexNet, GoogleNet, SqueezeNet (Pham 2021) | 99% | 39% |
| COVNet (L. Li et al. 2020) | 96% | 57% |
| ResNet-101 and Xception (Ardakani et al. 2020) | 99% | 39.6% |
| DarkCovidNet (Ozturk et al. 2020) | 98.08% | 58.8% |
| 3-D CNN ResNet-18 (Xu et al. 2020) | 86.7 | 52.02% |
| 3-D DNN, DeCoVNet (Zheng et al. 2020) | 90% | 36% |

**Table 5.** Results' accuracy of our proposed system.

| Diseases | TP | FP | Sensitivity | Specificity | FPR | Accuracy = TPR*100 |
|---|---|---|---|---|---|---|
| Bacterial | 152 | 48 | 0.76 | 0.94 | 0.24 | 76% |
| COVID-19 | 190 | 10 | 0.95 | 0.9 | 0.05 | 95% |
| Non-COVID-19 Viral | 129 | 51 | 0.717 | 0.98 | 0.283 | 71.7% |
| Initial stage COVID-19 | 60 | 30 | 0.667 | 1.0 | 0.333 | 66.7% |
| Normal | 104 | 26 | 0.8 | 0.91 | 0.2 | 80% |
| Average | | | 0.77 | 0.94 | | 77.88% |

images for three class classification and secured 98.5%,89.6%, 97%, and 98% accuracy, respectively. However, the accuracy of these method falls down to 59%, 53.7%, 58.2%, and 58.8%, respectively, when applied on testing dataset because they are incapable to detect initial stage COVID-19. The deep learning models proposed in (Xu et al. 2020) obtained an accuracy of 86.7% on three classes chest CT Scan images. However, its accuracy is dropped in case of testing dataset because it has no potential of initial stage COVID-19 detection. The drop in accuracy of deep learning models 14 is comparatively lower because it can classify four classes of chest X-Ray.



**Figure 2.** Accuracy comparison.

In contrast to the aforementioned deep learning models, our proposed system achieved an average accuracy 77.88% by employing deep learning models on the patients' symptoms and X-Ray images which is far better than the available COVID-19 detection methods. The results' accuracy of our proposed system on testing dataset is shown in Table 5.

The accuracy comparison of our proposed system with the existing ones is shown visually in Figure 2 In summary, the researchers achieved high accuracy results in case of binary classification (Ahamad et al. 2020; Ardakani et al. 2020; Horry et al. 2020; Pham 2021; Zheng et al. 2020) and multi-class classification like (Brunese et al. 2020; Hussain et al. 2020; L. Li et al. 2020; Ohata et al. 2020; Ozturk et al. 2020; Xu et al. 2020). But they have no potential to detect the COVID-19 in initial stage because they use single modality like symptoms or medical images. While our proposed system detects initial stage COVID-19 in addition to bacterial, COVID-19, non-COVID-19 viral, and normal with the help of utilizing both symptoms and medical images of patients. So, due to initial stage COVID-19 detection the average accuracy of our proposed system is better than existing systems (Ahamad et al. 2020; Hussain et al. 2020; Ozturk et al. 2020; Zheng et al. 2020).

In future, we intend to form a dataset containing the symptoms and medical images of each patient and train our proposed system on that dataset.

## Conclusion And Future Work

The available pneumonia classification and COVID-19 detection methods fail to detect COVID-19 in initial stage. The reason behind their failure is the fact that each of them is either using symptoms or medical images of patients for pneumonia classification. None of them utilized both symptoms and medical images of patients for pneumonia classification. By getting inspiration from the fact that the symptoms of COVID-19 appear 3 to 12 before the generation of COVID-19 infection, we proposed a hybrid system for detecting COVID-19 in initial stage. The experimental results proved the superiority of our proposed system in detecting COVID-19 in initial stage. In future, we intend to use microscopic images of nasal swab for COVID-19 detection. We also plan to work on predicting the severity of lungs infected by COVID-19 and predict the infection pattern of COVID-19.

## Disclosure Statement

## ORCID

Lal Hussain ⃝ http://orcid.org/0000-0003-1103-4938
Mamoona Humayun ⃝ http://orcid.org/0000-0001-6339-2257

## References

Ahamad, M. M., S. Aktar, M. Rashed-Al-Mahfuz, S. Uddin, P. Liò, H. Xu, M. A. Summers, J. M. Quinn, and M. A. Moni. 2020. A machine learning model to identify early stage symptoms of SARS-Cov-2 infected patients. *Expert Systems with Applications* 160:113661. doi:10.1016/j.eswa.2020.113661.

Ardakani, A. A., A. R. Kanafi, U. R. Acharya, N. Khadem, and A. Mohammadi. 2020. Application of deep learning technique to manage COVID-19 in routine clinical practice using CT images: results of 10 convolutional neural networks. *Computers in Biology and Medicine* 121:103795. doi:10.1016/j.compbiomed.2020.103795.

Bganglia. (2020). *covid-chestxray-dataset* https://github.com/ieee8023/covid-chestxray-dataset

Brunese, L., F. Mercaldo, A. Reginelli, and A. Santone. 2020. Explainable deep learning for pulmonary disease and coronavirus COVID-19 detection from X-rays. *Computer Methods and Programs in Biomedicine* 196:105608. doi:10.1016/j.cmpb.2020.105608.

Canas, L. S., C. H. Sudre, J. C. Pujol, L. Polidori, B. Murray, E. Molteni, M. S. Graham, K. Klaser, M. Antonelli, and S. Berry. 2021. Early detection of COVID-19 in the UK using self-reported symptoms: A large-scale, prospective, epidemiological surveillance study. *The Lancet Digital Health* 3 (9):e587–e598. doi:10.1016/S2589-7500(21)00131-X.

Chan, J. F.-W., S. Yuan, K.-H. Kok, K. K.-W. To, H. Chu, J. Yang, F. Xing, J. Liu, C. C.-Y. Yip, and R. W.-S. Poon. 2020. A familial cluster of pneumonia associated with the 2019 novel coronavirus indicating person-to-person transmission: A study of a family cluster. *The Lancet* 395 (10223):514–23. doi:10.1016/S0140-6736(20)30154-9.

Chorba, T. 2020. The concept of the crown and its potential role in the downfall of coronavirus. *Emerging Infectious Diseases* 26 (9):2302. doi:10.3201/eid2609.AC2609.

Chunxiaozheng. (2020). *COVID-19-tracker* https://github.com/BDBC-KG-NLP/COVID-19-tracker

Cohen, J. P., P. Morrison, L. Dao, K. Roth, T. Q. Duong, and M. Ghassemi. 2020. Covid-19 image data collection: prospective predictions are the future. *arXiv preprint arXiv:2006.11988*.

Dataset, N., and I. O. H. C. X.-r. (2017). *Random sample of nih chest x-ray dataset* https://www.kaggle.com/nih-chest-xrays/sample

Einstein. (2020). *Diagnosis of Covid-19 and its clinical spectrum* https://www.kaggle.com/einsteindata4u/covid19

Gov, D. (2021). *COVID-19 REPOSITORY* https://data.gov.il/dataset/covid-19

Horry, M. J., S. Chakraborty, M. Paul, A. Ulhaq, B. Pradhan, M. Saha, and N. Shukla. 2020. COVID-19 detection through transfer learning using multimodal imaging data. *IEEE Access* 8:149808–24. doi:10.1109/ACCESS.2020.3016780.

Hussain, L., T. Nguyen, H. Li, A. A. Abbasi, K. J. Lone, Z. Zhao, M. Zaib, A. Chen, and T. Q. Duong. 2020. Machine-learning classification of texture features of portable chest X-ray accurately classifies COVID-19 lung infection. *BioMedical Engineering OnLine* 19 (1):1–18. doi:10.1186/s12938-020-00831-x.

Imon, A. R., M. C. Roy, and S. Bhattacharjee. 2012. Prediction of rainfall using logistic regression. *Pakistan Journal of Statistics and Operation Research* 8 (3):655–67. doi:10.18187/pjsor.v8i3.535.

institute, N. c. (2020). *sars-cov-2* https://www.cancer.gov/publications/dictionaries/cancer-terms/def/sars-cov-2

Jain, R., M. Gupta, S. Taneja, and D. J. Hemanth. 2021. Deep learning based detection and analysis of COVID-19 on chest X-ray images. *Applied Intelligence* 51 (3):1690–700. doi:10.1007/s10489-020-01902-1.

Kaggle. (2018). *RSNA Pneumonia Detection Challenge* https://www.kaggle.com/c/rsna-pneumonia-detection-challenge

Khan, A. I., J. L. Shah, and M. M. Bhat. 2020. CoroNet: A deep neural network for detection and diagnosis of COVID-19 from chest x-ray images. *Computer Methods and Programs in Biomedicine* 196:105581. doi:10.1016/j.cmpb.2020.105581.

Koushik, C., R. Bhattacharjee, and C. S. Hemalatha. 2021. Symptoms based early clinical diagnosis of COVID-19 cases using hybrid and ensemble machine learning techniques. Ed. Eds.,*5th International Conference on Computer, Communication and Signal Processing (ICCCSP)*, IEEE, Chennai, India, 59-64.

Krizhevsky, A., I. Sutskever, and G. E. Hinton. 2012. Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems* 25:1097–105.

Li, F. 2016. Structure, function, and evolution of coronavirus spike proteins. *Annual Review of Virology* 3 (1):237–61. doi:10.1146/annurev-virology-110615-042301.

Li, L., L. Qin, Z. Xu, Y. Yin, X. Wang, B. Kong, J. Bai, Y. Lu, Z. Fang, and Q. Song. 2020. Artificial intelligence distinguishes COVID-19 from community acquired pneumonia on chest CT. *RSNA Radiological Society of North America* 296 (2): 1-16.

Liu, Y.-C., R.-L. Kuo, and S.-R. Shih. 2020. COVID-19: the first documented coronavirus pandemic in history. *Biomedical Journal* 43 (4):328–33. doi:10.1016/j.bj.2020.04.007.

Maladkar, K. (2018). *Overview of convolutional neural network in image classification* https://analyticsindiamag.com/convolutional-neural-network-image-classification-overview/

Marateb, H. R., F. Z. Nezhad, M. R. Mohebian, R. Sami, S. H. Javanmard, F. D. Niri, M. Akafzadeh-Savari, M. Mansourian, M. A. Mañanas, and M. Wolkewitz. 2021. Automatic classification between COVID-19 and non-COVID-19 pneumonia using symptoms, comorbidities, and laboratory findings: the khorshid COVID cohort Study.*Frontiers in Medicine* 8: 1-14.

Medlineplus. (2021). *Pneumonia* https://medlineplus.gov/pneumonia.html

Mooney, P. (2018). *Chest X-Ray images (pneumonia)* https://www.kaggle.com/paultimothy mooney/chest-xray-pneumonia

Nshomron. (2020). *Covidpred* https://github.com/nshomron/covidpred

Ohata, E. F., G. M. Bezerra, J. V. S. Das Chagas, A. V. L. Neto, A. B. Albuquerque, V. H. C. de Albuquerque, and P. P. Reboucas Filho. 2020. Automatic detection of COVID-19 infection using chest X-ray images through transfer learning. *IEEE/CAA Journal of Automatica Sinica* 8 (1):239–48. doi:10.1109/JAS.2020.1003393.

Organization, W. H. 2020. *COVID 19 public health emergency of international concern (PHEIC)*. Global research and innovation forum: towards a research roadmap.

Ozturk, T., M. Talo, E. A. Yildirim, U. B. Baloglu, O. Yildirim, and U. R. Acharya. 2020. Automated detection of COVID-19 cases using deep neural networks with X-ray images. *Computers in Biology and Medicine* 121:103792. doi:10.1016/j.compbiomed.2020.103792.

Pham, T. D. 2021. Classification of COVID-19 chest X-rays with deep learning: New models or fine tuning? *Health Information Science and Systems* 9 (1):1–11. doi:10.1007/s13755-020-00135-3.

Rahman, T. (2021). *COVID-19 radiography database* https://www.kaggle.com/tawsifurrahman/covid19-radiography-database

Rehman, A., M. A. Iqbal, H. Xing, and I. Ahmed. 2021. COVID-19 detection empowered with machine learning and deep learning techniques: A systematic review. *Applied Sciences* 11 (8):3414. doi:10.3390/app11083414.

Sami, R., F. Soltaninejad, B. Amra, Z. Naderi, S. Haghjooy Javanmard, B. Iraj, S. Haji Ahmadi, A. Shayganfar, M. Dehghan, and N. Khademi. 2020. A one-year hospital-based prospective COVID-19 open-cohort in the eastern mediterranean region: the khorshid COVID cohort (KCC) study. *PloS one* 15 (11):e0241537. doi:10.1371/journal.pone.0241537.

Sánchez, J. L., E. G. Ortega, and J. L. Marcos. 2001. Construction and assessment of a logistic regression model applied to short-term forecasting of thunderstorms in león (Spain). *Atmospheric Research* 56 (1–4):57–71. doi:10.1016/S0169-8095(00)00089-2.

Talo, M. (2020). *COVID-19* https://github.com/muhammedtalo/COVID-19

Wong, H. Y. F., H. Y. S. Lam, A. H.-T. Fong, S. T. Leung, T. W.-Y. Chin, C. S. Y. Lo, M. M.-S. Lui, J. C. Y. Lee, K. W.-H. Chiu, and T. W.-H. Chung. 2020. Frequency and distribution of chest radiographic findings in patients positive for COVID-19. *Radiology* 296 (2):E72–E78. doi:10.1148/radiol.2020201160.

Xu, X., X. Jiang, C. Ma, P. Du, X. Li, S. Lv, L. Yu, Q. Ni, Y. Chen, and J. Su. 2020. A deep learning system to screen novel coronavirus disease 2019 pneumonia. *Engineering* 6 (10):1122–29. doi:10.1016/j.eng.2020.04.010.

Zheng, C., X. Deng, Q. Fu, Q. Zhou, J. Feng, H. Ma, W. Liu, and X. Wang. 2020. *Deep learning-based detection for COVID-19 from chest CT using weak label.* MedRxiv, USA.

Zoabi, Y., S. Deri-Rozov, and N. Shomron. 2021. Machine learning-based prediction of COVID-19 diagnosis based on symptoms. *Npj Digital Medicine* 4 (1):1–5. doi:10.1038/s41746-020-00372-6.