



A machine learning-based framework for diagnosis of COVID-19 from chest X-ray images

Jawad Rasheed¹ · Alaa Ali Hameed¹ · Chawki Djeddi² · Akhtar Jamil¹ · Fadi Al-Turjman³

Received: 12 September 2020 / Revised: 5 November 2020 / Accepted: 20 November 2020 / Published online: 2 January 2021
© International Association of Scientists in the Interdisciplinary Areas 2021

Abstract

Corona virus disease (COVID-19) acknowledged as a pandemic by the WHO and mankind all over the world is vulnerable to this virus. Alternative tools are needed that can help in diagnosis of the coronavirus. Researchers of this article investigated the potential of machine learning methods for automatic diagnosis of corona virus with high accuracy from X-ray images. Two most commonly used classifiers were selected: logistic regression (LR) and convolutional neural networks (CNN). The main reason was to make the system fast and efficient. Moreover, a dimensionality reduction approach was also investigated based on principal component analysis (PCA) to further speed up the learning process and improve the classification accuracy by selecting the highly discriminate features. The deep learning-based methods demand large amount of training samples compared to conventional approaches, yet adequate amount of labelled training samples was not available for COVID-19 X-ray images. Therefore, data augmentation technique using generative adversarial network (GAN) was employed to further increase the training samples and reduce the overfitting problem. We used the online available dataset and incorporated GAN to have 500 X-ray images in total for this study. Both CNN and LR showed encouraging results for COVID-19 patient identification. The LR and CNN models showed 95.2–97.6% overall accuracy without PCA and 97.6–100% with PCA for positive cases identification, respectively.

Keywords Artificial neural network · Computer-aided diagnosis · COVID-19 · Image classification · Principal component analysis

1 Introduction

The COVID-19, an infectious virus, emerged in China at the end of December 2019. It has spread to at least 213 countries around the globe as of August 02, 2020. Therefore, on March 11, 2020, the WHO declared the COVID-19 outbreak as pandemic. According to Johns Hopkins University,¹ the tally of confirmed cases has exceeded 25,225,566 while the death toll is over 847,676 until August 30, 2020. It is also confirmed that 17,688,088 people diagnosed with COVID-19 have recovered from the disease so far.

Besides affecting humans, this huge family of viruses, called coronaviruses, mostly affects various animals and other species such as bats, cats, and camels. These animal coronaviruses once infect a human, can then become a cause of horizontal transfer (human-to-human) such as with severe acute respiratory syndrome coronavirus (SARS-CoV), Middle East respiratory syndrome coronavirus (MERS-CoV), and currently with this infectious virus named as COVID-19 (also known as SARS-CoV-2). According to National Institutes of Health,² the SARS-CoV, MERS-CoV, and SARS-CoV-2 viruses are originated in animal reservoir. The history of human coronaviruses started in 1965 when B814 virus was discovered in a human body by [1]. In last 20 years, scientists have discovered more than five new human coronaviruses that have caused substantial rise in mortality and morbidity [2].

✉ Jawad Rasheed
jawadrasheed@ieee.org

¹ Department of Computer Engineering, Istanbul Sabahattin Zaim University, 34303 Istanbul, Turkey

² Department of Mathematics and Computer Science, Larbi Tebessi University, 12018 Tébessa, Algeria

³ Artificial Intelligence Department, Research Center for AI and IoT, Near East University, Nicosia, Mersin 10, Turkey

¹ Johns Hopkins University: COVID-19 Maps—John Hopkins Coronavirus Research Center, <https://coronavirus.jhu.edu/map.html>

² National Institutes of Health: <https://www.nih.gov/health-information/coronavirus>

Until now, scientists could not find any treatment for COVID-19. The symptoms of the virus involve respiratory illness of mild to severe intensity, sore throat, coughing, diarrhoea etc. [3]. This virus spreads in the form of droplets of saliva when someone sneezes or through a physical contact [4]. However, recently, WHO declared that it is also airborne. On the other hand, elder people or humans with chronic medical history or critical diseases like cancer, chronic respiratory problem, cardiovascular disease or diabetes may require special treatment. As mentioned by WHO,³ such people are more vulnerable and develop serious illness when contracted with COVID-19. So far, the only response to this virus is to rely on isolation, quarantine, and infection-control measures to combat and control the COVID-19 outbreak [5].

Early detection of this disease can help in timely isolation of patients and monitor their health status [6]. Machine learning-based approaches can be used that analyze the lungs' X-ray/CT images to identify patients affected by pneumonia due to COVID-19 infection. This technique can be used as an alternative where COVID-19 kits are not available, especially in developing countries where a large population is affected by this virus but no measures could be arranged to confirm the suspects for COVID-19.

Supervised learning techniques have shown great progress for early detection and diagnosis of diseases. For instance, in [7], a decision support system is proposed for prediction of diabetes using machine learning techniques. Three machine-learning algorithms were used that includes CNN, random forest, and support vector machine (SVM). The experiments showed promising results for classification of 768 patients into diabetic and non-diabetic groups. Similarly, [8] compared most popular machine learning techniques commonly used for detection of breast cancer. The methods implemented include random forest, k-nearest-neighbor (kNN) and Naïve Bayes classifiers. In [9], a wavelet transforms and SVM-based methodology is proposed for categorization of brain tumors into two different classes as benign or malignant. Sun et al. [10] measured the diagnostic performance by investigating 15 different classification and several feature selection methods in glioma grading. The results indicated that the combination of feature selection with linear SVM and multilayer perceptron (MLPC) achieved the best performance. A comprehensive assessment of machine learning-based futuristic approaches in medical image analysis can be found in [11].

In spite the success of machine learning approaches, their performance is highly affected by the quality of the hand-engineered features. The hand-engineered features are not optimal and also a time-consuming task [12]. This is the

main drawback of machine learning approach, despite the success, these techniques suffer from serious degradation in performance. The alternative approach, which has become hot research topic in recent era, is to extract automatic optimal features from the input data. These techniques have not only removed the barrier of manual feature extraction but also improved the classification accuracy.

Several techniques based on deep learning (DL) have been put forward to achieve the above goals: automatic feature extraction and improving classification accuracy. With improvements in computer hardware, it has become feasible to train more and more complex models. DL approaches, especially CNNs have shown their efficiency for various computer vision tasks such as object detection, natural language processing, image segmentation and classification. Motivated from the computer vision community, the medical community has also adopted the model to solve many medical image analysis tasks. For instance, in [13], authors investigated two well-known deep neural networks: Inception and VGG16, for diagnosis of pneumonia from images of chest X-ray. It was reported that VGG16 resulted in higher classification accuracy compared to Inception model. Segmentation of lung X-ray images using deep CNNs is proposed in [14] for improving the performance during clinical diagnosis of various diseases in lungs such as lung cancer, tuberculosis, or lung opacities. Van Tulder and De Bruijne [15] used an unsupervised feature learning based on restricted Boltzmann machines for feature extraction with a generative learning objectives. It combines both generative and discriminative learning objectives with convolutional classification for categorization of lung computed tomography (CT) images. Based on histology images, [16] developed a reliable system to enhance the diagnostic quality for identification of breast cancer. Particularly, two machine learning approaches were compared: SVM was used as classifier using handcrafted features while another approach is based on CNNs for automatic feature learning and classification. The results indicated that CNN performed well than the classifier based on handcrafted features.

In [17], various CNN architectures were investigated for classification and detection of interstitial lung disease and thoraco-abdominal lymph node, respectively. Further analysis was performed to evaluate the effect of spatial image context and dataset scale on classifier performance, and the application of transfer learning from pre-trained ImageNet in the domain of image analysis.

Besides conventional down sampling layers, [18] proposed an atrous convolution as an alternative layer in the deep CNNs. Since CNNs mostly exploit down sampling layers to increase the receptive field and gain abstract semantic information; however, the down sampling also decreases the feature maps' spatial dimensions that may not be desirable for semantic segmentation tasks. The atrous convolutions,

³ WHO: World Health Organization, <https://www.who.int/>

on the other hand, can increase the receptive field without changing the feature maps' spatial dimension. A comprehensive analysis of the present day furtherance in machine learning, specifically in DL that helped in distinguishing, categorizing and gauging the patterns in medical images is presented in [19].

Despite the success of DL models, their application for the examination of medical image is particularly a challenging task. The deeper networks usually require large datasets for training; however, large labeled datasets of medical images are not abundant compared to vision related datasets such as ImageNet [20]. In addition, the imbalance datasets and poor representation can make the problem even more complex. Moreover, the privacy and confidentiality concerns related to medical data of patients also limit the access to the data [21].

To overcome this limited availability of labelled data, researchers proposed data augmentation techniques. Data augmentation technique is utilized to elevate the amount of training data by working in various transformation such as rotation, scaling, translations etc. on the original data. It also helps reducing the common problem of overfitting. A detailed review of recent augmentation techniques for DL are described in details in [22].

Data augmentation is also adopted in various medical image analysis applications such as semantically segmenting the various sclerosis lesions using magnetic resonance imaging (MRI) of brain [23], cardiac image enhancement and segmentation or reconstruction [24, 25], mitosis detection in breast cancer histology images [26] and brain tumor segmentation using MRI images [27].

Conversely, researchers have also taken advantage of AI-based internet of things (IoT) in the field of medical sciences for diagnostic purposes. Such as, [28] proposed an application of IoT that detects respiratory motion at real time by monitoring diabetic patient's breathing to diagnose diabetic ketoacidosis (DKA). They used C-band sensing method by exploiting microwave-sensing platform (MSP) as a non-intrusive respiratory monitoring system. They further used peak detection algorithm that acquires respiratory rate for identification of Kussmaul breathing. Similarly, [29] used S-band sensing technique to characterize wandering patterns in patients suffering from dementia. Later, researchers incorporated SVM as a pattern classification algorithm.

In this article, researchers propose a machine learning-based framework for detection of COVID-19 from X-ray images. We used one traditional machine learning approach, LR, which is an efficient and simple method to implement. In addition, for DL, we implemented deep CNN. Moreover, we investigated the application of feature selection technique (PCA) to lower the computational time and enhance the overall accuracy.

The rest of the paper is organized as follows. Section 2 summarizes the details of dataset used in this study and the proposed methodology. Section 3 outlines the experimental setup and presents the results, while Sect. 4 provides the summary and discussion. Finally, the paper concludes at Sect. 5.

2 Materials and methods

This section defines the dataset and methodologies used for prediction of COVID-19 virus. It outlines the experimental dataset, data augmentation technique, feature selection method, DL and machine learning predictive models incorporated in this work. Workflow of system put forward by researchers is shown in Fig. 1, that depicts the preprocessing performed on compiled dataset and various sets of features pertaining high information are extracted via PCA. These extracted sets of features are then fed to proposed DL model and traditional machine learning network for training, and later their performances are evaluated on respective test sets.

2.1 Dataset

The experimental work is performed by acquiring X-rays images of COVID-19 infected persons from COVID-19 images data collection provided by Joseph Paul Cohen,⁴ and normal X-rays images from Chest X-ray Images⁵ (pneumonia) repository as healthy individuals. At the time of this study, COVID-19 images data collection contains only 198 X-rays images pertaining to COVID-19 affected cases. The experimental work is developed with the help of X-ray images of sufferers from COVID-19 data collection by taking 198 X-rays images of virus-affected patients and 210 X-rays images of healthy individuals from Chest X-ray Images (pneumonia) repository. The original size of the compiled dataset varies in range of 1112×624 to 2170×1953 pixels. For this experimental work, all images are scaled to 512×512 pixels. Table 1 summarizes the data set used in this study. The Fig. 2 depicts few samples from both classes of dataset. Data augmentation was applied to add additional data to create class balance. The final version of dataset contains 250 samples for each class. The class label 1 was used to represent COVID-19 positive cases while 0 represents normal healthy cases.

⁴ Joseph Paul Cohen: COVID-19 image data collection, <https://doi.org/2003.11597> (last accessed 2020/07).

⁵ P. Mooney: Chest X-Ray Images (Pneumonia), <https://doi.org/10.17632/rschbjbr9sj.2> <https://www.kaggle.com/paultimothymooney/chest-xray-pneumonia> (last accessed 2020/07).

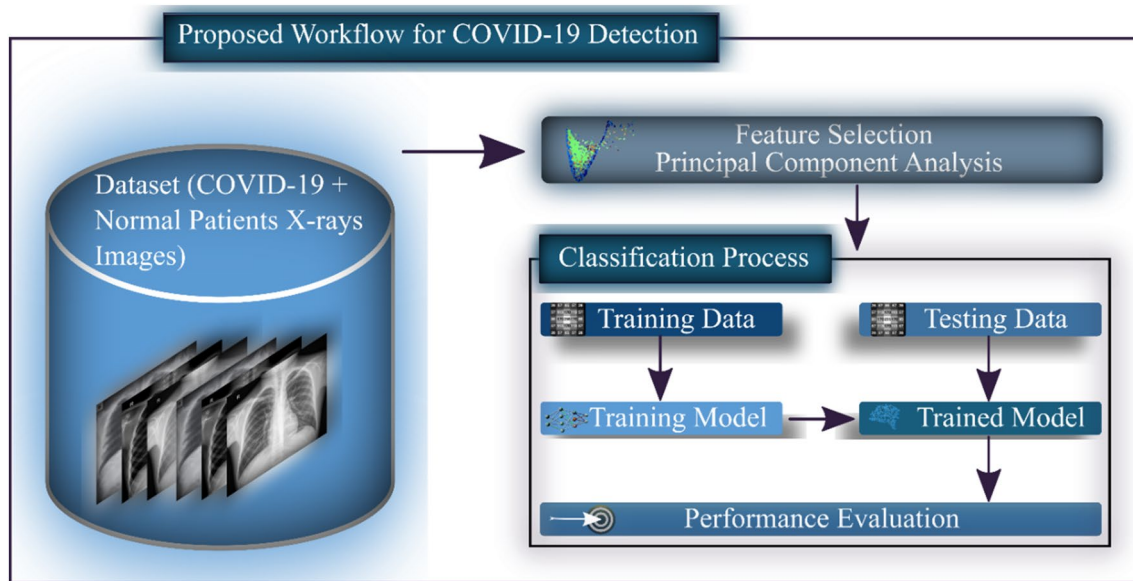


Fig. 1 Workflow of proposed system

Table 1 Dataset information

Cases	Original samples	Augmented samples	Total
COVID-19	198	52	250
Healthy individuals	210	40	250
Total	408	92	500

2.2 Data augmentation with GAN

GAN has been extensively utilized in numerous image generation functions that produces artificial data-like real data. This can help to overcome the issue of small number of training samples thus eliminates class imbalance problem. The GAN architecture consists of multilayer perceptron which has two main elements: generator (G) and discriminator (D) [30]. These two components compete with each

other during training. The generator is trained to produce data similar to the original data and the discriminator should be able to distinguish between fake and actual data.

According to [30], x be the input data. For learning the generator's distribution p_g , a prior is defined on input noise variables $p_z(z)$. $G(z; \theta_g)$, defines the data space mapping, where G corresponds to differentiable function having network parameters θ_g . This differentiable function is implemented as a multilayer perceptron network. Also, another multilayer perceptron is define, $D(x; \theta_d)$ that takes the input from G and produces the final output. $D(x)$ represents the probability that x is original or fake. We trained D to enhance the possibility of assigning the appropriate label to both training instances as well as instances from G . The objective function $V(G, D)$ for both G and D can be combined to train the GAN model as [30]:

$$\min_G \max_D V(G, D) = E_{x \sim p_{\text{data}}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log (1 - D(G(z)))], \quad (1)$$

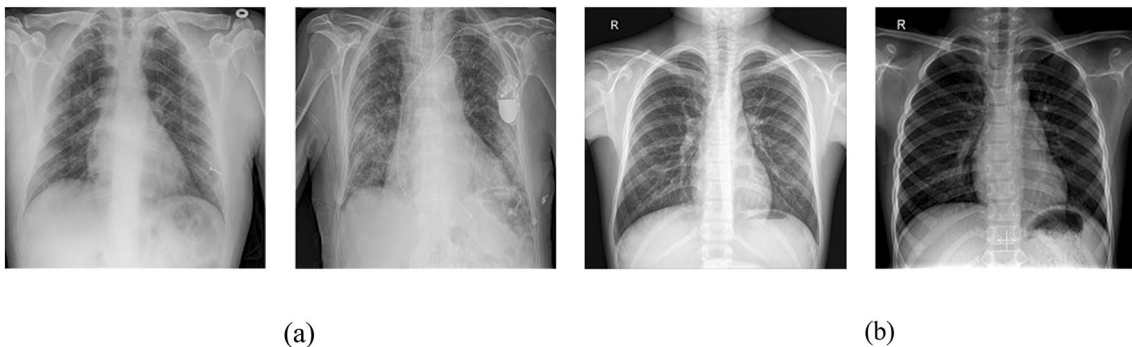


Fig. 2 Samples of X-ray image dataset used for proposed system; **a** images of COVID-19 affected cases, **b** images of healthy individuals

Table 2 GAN parameter settings

Parameter name	Value
Activation function	ReLU
Batch size	4
Learning rate	0.0005
Loss function	Binary cross-entropy
Dropout probability	0.2
Optimizer	Adam
Convolutional layer	2
Kernel size	5×5

where E is the expectation. Equation (1) indicates that we should train the model in such a way that error for G is minimized and D should be maximized so that it should not be able to distinguish between fake and real data.

In this study, the synthetic data were generated using the same architectural setup for both G and D as summarized in Table 2. The data augmentation was achieved using horizontal and vertical shifts, and random r .

2.3 Feature selection

Feature selection plays a crucial role for raw data representation. It is among top of the list hottest research subject in computer vision and machine learning domain. The main aim is to obtain highly discriminative features from the raw data that have potential to enhance the classification accuracy of the classifier.

The explosion of data set size triggers the development of various data dimensionality reduction techniques to boost the performance of data mining and classification systems. Various feature selection and extraction methods have been developed such as random-forest feature selection, PCA, linear discriminant analysis (LDA), forward feature selection and backward feature elimination methods. We employed PCA as feature extraction technique due to its simplicity, efficiency and popularity for being an oldest multivariate technique.

PCA is a multivariate statistical procedure that analyzes dependent and inter-correlated variables in original dataset and extracts the important information by transforming it to a new set of orthogonal variables called principal components [31]. A new set of principal components is attained, each with certain variance, while the first principal component attains the highest variance among others. Amount of information to retain strongly depends on selection of principal components; therefore, maximum

amount of information can be retained by selecting appropriate amount of principal components to reduce data dimensionality.

PCA reduces the 2-D matrix, X (N, M) pertaining images, (where M is total pixels after masking, and N is number of instances such that $N < M$), to smaller matrix $Z(N, L)$, (where L is the number of pixels such that $L < M$), while retaining much information from data, using linear transformation $U(M, L)$ [31, 32].

$$Z = U^T X. \tag{2}$$

It calculates covariance matrix S (L, L) to represent the information as

$$S_Z = \frac{1}{N} Z^T Z. \tag{3}$$

The maximization of covariance yields eigenvector equations with Lagrange multiplier, λ . These eigenvector equations are then decomposed using matrix diagonalization that results S as a product of three matrices:

$$S = PDP^{-1}, \tag{4}$$

where D corresponds to diagonal matrix, consists of Eigen values, and P refers to matrix of eigenvector. Therefore, the sum of eigenvalues corresponds to entire variance of the transformation is

$$\text{Total variance} = \sum_{i=1}^M \lambda_i \tag{5}$$

To project the top L eigenvectors data along the subset of these M vectors, the variance retained is

$$\text{Retained variance} = \sum_{i=1}^L \lambda_i \tag{6}$$

Hence, the amount of information retained is expressed as percentage of the original using

$$\text{Percentage of information retained} = \frac{\sum_{i=1}^L \lambda_i}{\sum_{i=1}^M \lambda_i} \tag{7}$$

Using these equations, PCA first calculates the mean of every dimension of whole dataset, computes the covariance matrix and determines the eigenvector and eigenvalue pairs using matrix diagonalization. Later, it sorts these pairs by decreasing order of eigenvalues. Since the eigenvalues are proportional to variance retained, the selection of top L pairs will retain most of the information while using only fraction of original dimensions.

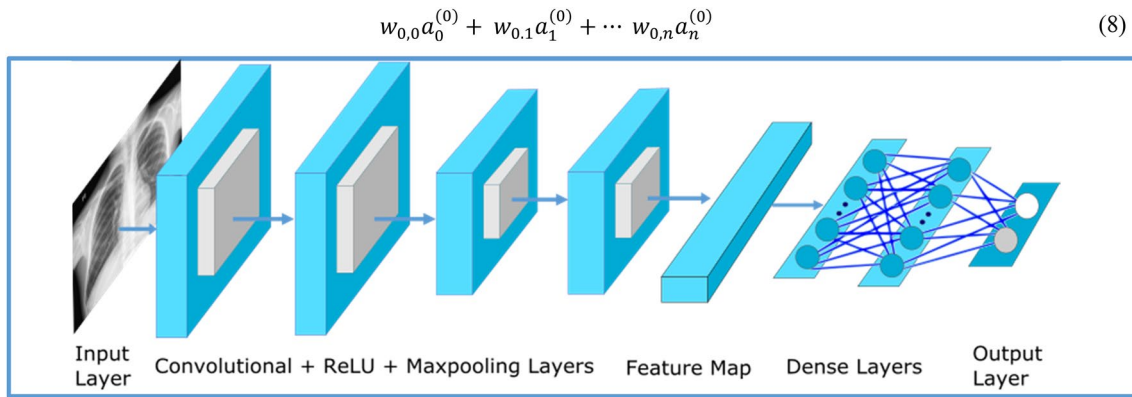


Fig. 3 Architecture of the proposed CNN

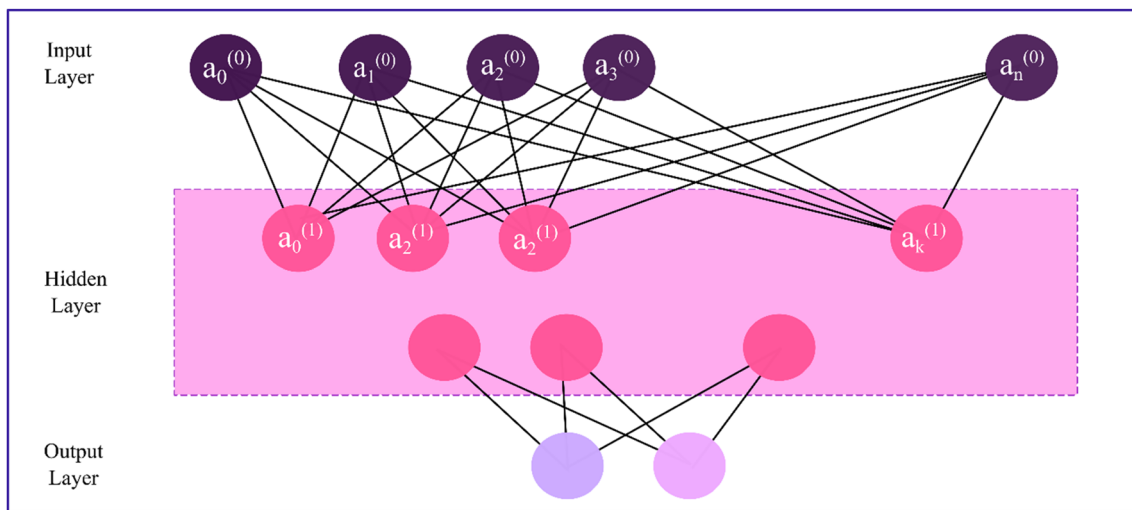


Fig. 4 Visual representation of neurons and weights in CNN

2.4 Classification methods

This section describes machine learning and DL models used for this study to identify patients affected by COVID-19 infection. To accomplish the goal of identifying COVID-19 positive patients among normal healthy individuals, CNN and LR models are used, as described in this section.

2.4.1 Convolutional neural network (CNN)

CNN, an extensively adopted classification algorithm in DL field, consists an architect of successive layers of perceptron connected in a sequence [33]. It distills the input image data by chain of connected layers and predicts the output by transforming it to meaningful representation. Normally, CNN architect has three major types of layer; convolutional layer, pooling layer and dense layer or fully connected layer, as depicted in Fig. 3.

Convolutional layer, which is a pivotal layer comprises of neurons connected to small region of preceding layers with some weights of shared characteristics, is shown in hidden layer of Fig. 4. The network takes the image in form of matrix as input map and assigns a weight ‘w’ to each one of the connections between neurons of input layer and convolutional layer. It computes weighted sums of all activations ‘a’ from the input layer according to these weights:

$$w_{0,0}a_0^{(0)} + w_{0,1}a_1^{(0)} + \dots + w_{0,n}a_n^{(0)}. \quad (8)$$

As computed weighted-sum might be any number, model normalizes these weighted-sum using functions like σ sigmoid that squishes these into range between 0 and 1. It determines the activation of neurons in layer by adding a suitable bias ‘b’ to relevant weighted-sum for meaningful activation. Thus, the first neuron in the second layer is represented by

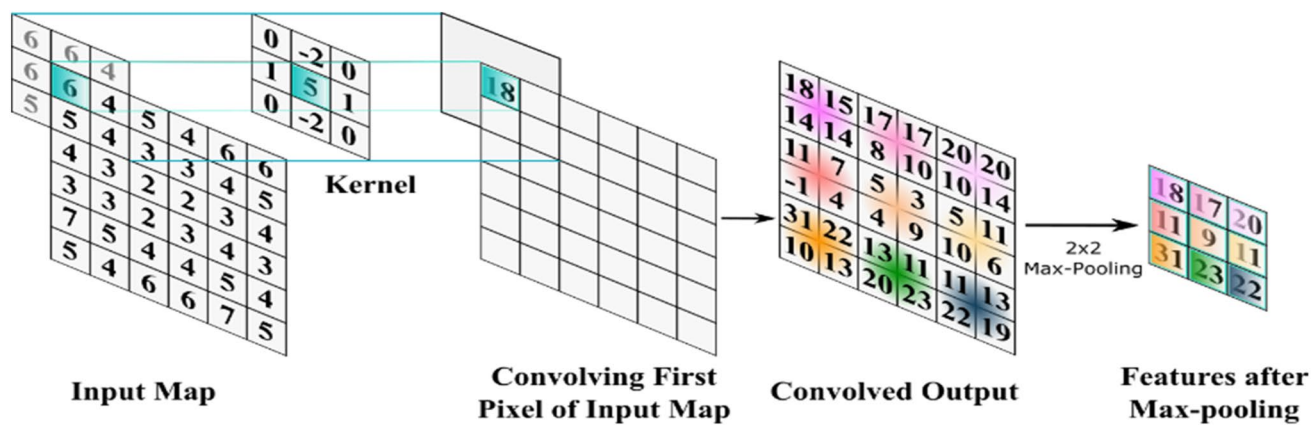


Fig. 5 Visual representation of kernel convolving with input vector and then applying 2 × 2 max-pooling

$$a_0^1 = \sigma(w_{0,0}a_0^{(0)} + w_0, 1a_1^0 + \dots + w_{0,n}a_n^0 + b_0). \tag{9}$$

The representation of all the neurons can be rearranged in form of matrix that corresponds to activations of next layer.

$$\sigma \left(\begin{bmatrix} w_{0,0} & \dots & w_{0,n} \\ \vdots & \ddots & \vdots \\ w_{k,0} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ \vdots \\ b_n \end{bmatrix} \right) = \begin{bmatrix} a_0^1 \\ \vdots \\ a_k^1 \end{bmatrix}. \tag{10}$$

Subsequently model discovers right set of weights and biases for all the layers to solve the problem in best way. In convolutional layers, network calculates the scaler product of kernel with input map, as depicted in Fig. 5, to yield activation maps, which are then stacked to perform final output. A non-linear function, rectified linear unit (ReLU), performs elementwise activation with resultant of convolutional layer.

After convolutional layer, pooling layer acts as fuzzy filter that drastically minimizes computational cost by reducing feature dimensionality. Max-pooling and average-pooling are famous among others for feature dimensionality reduction. Max-pooling reduces dimensionality by picking the feature pixel with maximum value among others in the window as illustrated in Fig. 5, whereas average-pooling takes the average of all the pixels in the window.

Finally, a dense layer gathers all the features extracted by previous layers for classification. Later, output layers classifies the input image accordingly by considering the loss based on soft-max layer probabilities.

2.4.2 Logistic regression (LR)

LR is a statistical prediction technique in machine learning to perform regression analysis when dependent variable (target) is dichotomous or categorical which can be extended to model of several classes of events [34]. It is named over its core method, known as sigmoid function. This S-shaped

curve (shown in Fig. 6) maps real value x in (11) to values between 0 and 1

$$\text{sig}(x) = \frac{1}{1 + e^{-x}}, \tag{11}$$

where x is the actual input numerical number and e is the base of natural logarithms.

LR models the probability of predicted output (y) in (12) by combining input (x) with weight and coefficients where b_0 is bias (intercept term) and b_1 is coefficient. Maximum-likelihood estimates these beta values from training data to minimize the error in probability prediction.

$$y = \frac{e^{b_0+b_1x}}{1 + e^{b_0+b_1x}}. \tag{12}$$

3 Experimental results

Along with the goal of predicting image as positive COVID-19 or negative, the primary aim of the research is to enhance the performance of model by performing feature extraction using PCA technique and analyze its effects. To achieve this, the compiled and augmented dataset comprised of 500 images (250 as COVID-19 positive cases and 250 normal healthy individuals) is randomly divided into training subset and testing subset with a ratio of 75:25 respectively. The first set, i.e. training set consists of 375 samples while testing set had 125 instances as depicted in Table 3.

3.1 Feature extraction using PCA

To enhance the performance of machine learning models, PCA feature extraction technique is employed. As the compiled X-rays images dataset are of different height and width dimensions, so images were scaled and normalized

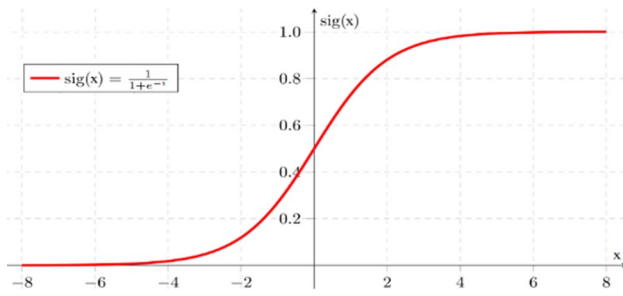


Fig. 6 Core function of LR

Table 3 Dataset split into training and testing set

Classes	Dataset	Training	Testing
COVID-19 cases	250	188	62
Normal cases	250	187	63
Total	500	375	125

as preprocessing step before feeding it to PCA for feature extraction. These normalized images data are then fed to PCA for orthogonal transformation to extract features with the most important data information. Each principal component in new set corresponds to some crucial information as summarized in Table 4. For this experimental setup, only those top features are selected that equates to variance of 1, 0.99, 0.98, 0.97, 0.96, 0.95, 0.90, and 0.85.

3.2 Classification

Once the principal feature sets of dataset are extracted by PCA, each set is then used to train the proposed CNN model and conventional LR network separately. These trained

machine learning and DL models are evaluated on respective testing data, and performance analysis is presented in this section.

3.2.1 CNN

A CNN model is proposed in Fig. 3 for identification of patients as COVID-19-positive case or a COVID-19-negative case by analyzing an input X-ray image. Respective layers in model transform the input to substantial representation and transmit it to following layers for further operations. The introductory layer in Fig. 3 is the input layer with size of $H \times W \times 3$ equivalent to size of input X-ray images. The height ‘H’ and width ‘W’ of input layer varies in each model according to each principal component feature sets. For feature set with variance 1, $H = W = 512$.

Next to input layer, two convolutional layers are created successively with an output shape of $H \times W \times 32$, having a 3×3 kernel, along with padding of same size to assure the size of spatial output. Afterwards, a 2×2 max-pool layer reduces features dimensionality, and dropout is fixed to 0.25. Later, two more convolutional layers are placed with depth of 64 and padding is turned on. A max-pool layer of size 2×2 takes the outcome of these convolutional layers for dimensionality reduction, and dropout of 0.4 is performed. In all the convolutional layers, ReLU is used as activation function. The resultant features are flattened and passed to dense layer of 128 depth, followed by a dropout of 0.5. A fully connected layer puts together all the features of preceding layers and the final output of fully connected layer is normalized by shaping probability outcomes into final prediction with a *Softmax* activation function. The various sets of extracted features are then used to train this CNN model with *Adam* optimizer having learning rate of 0.001, beta_1

Table 4 Pixel values of complete dataset. (A) Normalized pixel value of original dataset images, (B) extracted feature values of top 147 principal components after applying PCA with variance 0.99

Sr. No	Original dataset normalized pixel values								
	Pixel 0	Pixel 1	Pixel 2	Pixel 3	...	Pixel 786429	Pixel 786430	Pixel 786431	
0	0.156863	0.156863	0.156863	0.176471	...	0.047059	0.047059	0.047059	
1	0.129412	0.129412	0.129412	0.137255	...	0.003922	0.011765	0.011765	
2	0.482353	0.482353	0.482353	0.466667	...	0.000000	0.000000	0.000000	
3	0.000000	0.000000	0.000000	0.019608	...	0.062745	0.290196	0.290196	
(A)									
Sr. No	147-Principal component values								
	PC 1	PC 2	PC 3	PC 4	...	PC 145	PC 146	PC 147	
0	- 27.311999	- 36.931316	11.087628	22.476118	...	- 3.152277	2.403450	- 1.161468	
1	- 72.68627	55.72737	5.388810	- 3.205570	...	0.287251	- 1.940857	3.880443	
2	- 30.392400	- 16.253388	- 0.724725	- 18.254980	...	- 1.018936	0.819682	2.595462	
3	- 7.005396	- 6.036791	- 22.227564	- 0.651227	...	0.315235	0.417503	0.745322	
(B)									

Table 5 Network topology of proposed CNN model

Dataset	Convolution layers				Fully connected layers		
	Kernel Size	ReLU	Max-pooling	Dropout	No. of neurons	Function	Dropout
Covid-19	3 × 3 × 32	Yes	No	No			
	3 × 3 × 32	Yes	2 × 2	Yes (25%)			
	5 × 5 × 64	Yes	No	No			
	5 × 5 × 64	Yes	2 × 2	Yes (40%)	128	ReLU	Yes (50%)
					2	Softmax	No
Other parameters							
Batch size	Steps (iterations)	Epochs	Learning rate	Optimizer			
8	100	25	0.001	Adam			

of 0.9, beta_2 of 0.999 while *epsilon* is set to none and *amsgrad* to false. CNN network is trained with online augmentation (augmentation on the fly technique) over 25 epochs with a batch size of 8. Table 5 summarizes the network topology and parameters used in proposed CNN model. Fig. 7 shows the accuracy and loss curves of model against various principal components sets, while the evaluation performance of proposed CNN classification network with different principal components are summarized in Table 6.

3.2.2 Linear regression

Beside DL-based CNN model, the extracted features are also fed to conventional LR model for training and testing. All the sets consisting of principal components of various variances are fed to separate LR networks and performance is recorded for each LR model against each set, which is depicted in receiver operating characteristics (ROC) curves, shown in Fig. 8. Table 7 summarizes the classification performance of LR model over selection of various PCA components.

3.3 Performance analysis

The classification accuracies of trained models (CNN and LF) corresponding to different PCA sets are summarized in Table 8 along with the time taken by model for training. The experimental results show that without feature extraction, the proposed CNN and LR models achieved accuracy of 97.6% and 95.2%, respectively. It is prominent by results in Table 8 that proposed CNN framework, trained with dataset of extracted features having 0.99 variance, outperformed other trained models by achieving 100% accuracy with just 233 ms of training. CNN model trained on dataset of 0.99 variance, predicted all the testing data correctly, while other models misclassified either a healthy person as COVID-19-infected patient or a COVID-19-infected patient as healthy person (see Fig. 9).

All the experimental work is carried out using Python 3.6 and Keras packages on Jupyter Notebook (version

6.0.3) having 32 GB system RAM with Intel® processor (Xeon® CPU E3-1231 3.40 GHz). Additionally, the code is implemented with help of graphical processing unit (GPU) NVIDIA Quadro K500.

4 Summary and discussion

The clinical and radiological COVID-19 images compiled by Dr. Joseph P. Cohen are used by research community to develop models for accurate identification of patients affected with COVID-19 infectious disease. Table 9 summarizes the attempts made by researchers to identify the patients affected by COVID-19 using combinations of various artificial intelligence methods. Togacar et al. [35] introduced DL frameworks based on different CNN architectures that uses X-ray images to diagnose COVID-19. They combined SVM classifier with MobileNetV2 and SqueezeNet, and exploited social mimic optimization method while incorporating fuzzy color scheme with stacking approaches. 99.27% accuracy and 98.58% F1 score was attained by the researchers on an X-ray images dataset comprising of 295 of COVID-19 positive cases and 163 of healthy or other pneumonia affected patients.

Tuncer et al. [36], with the help of residual exemplar local binary pattern (ResExLBP), diagnosed corona virus disease in digital lungs X-ray images by extracting the features, drew on iterative relief (iRF) and fed the resultant to six different machine learning classifiers. The dataset used for the study consists of 87 X-ray images related to patients affected by COVID-19 while 234 X-ray images of normal patients. They achieved a maximum accuracy of 99.69% with an overall sensitivity of 98.85%.

A combination of transfer learning with various CNN frameworks is adopted by Apostolopoulos and Mpesiana [37] to facilitate the detection of this pandemic disease. They inspected 224 confirmed coronavirus disease, 714 bacterial pneumonia, and 504 ordinary X-ray images and observed that MobileNetV2 outperformed other frameworks by

Fig. 7 Training and testing loss/accuracy graphs of CNN network when **a** variance=1, **b** variance=0.99 and components=147, **c** variance=0.98 and components=126, **d** variance=0.97 and components=108, **e** variance=0.96 and components=96, **f** variance=0.95 and components=84, **g** variance=0.90 and components=48, **h** when variance=0.85 and components=30

securing an accuracy and sensitivity of 97.40%, and 99.10% respectively. Brunese et al. [38] tailored the transfer learning technique with fine-tuned visual geometry group (VGG-16)

model that examined 6523 chest X-rays (250 images of patients affected by COVID-19, 2753 images corresponds to pulmonary diseases patients and 3520 images related to healthy persons). It first segregated between pulmonary diseases patients and healthy person, then classified the X-ray image of patient with former disease either as coronavirus disease positive or other pneumonia-affected person. Thus accomplished a success rate of 97.0% and F1 score of 92.0%.

Jaiswal et al. [39] blended DL with deep transfer learning approach by employing a modified pre-trained DenseNet201

Table 6 Performance evaluation of CNN model

Variance Retained	CNN								
	Precision			Recall			F measure		
	Corona	Normal	Overall	Corona	Normal	Overall	Corona	Normal	Overall
1	0.952	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.99	1	1	1	1	1	1	1	1	1
0.98	0.952	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.97	0.952	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.96	0.952	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.95	0.952	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.90	0.952	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.85	0.910	1	0.955	1	0.913	0.957	0.952	0.955	0.945

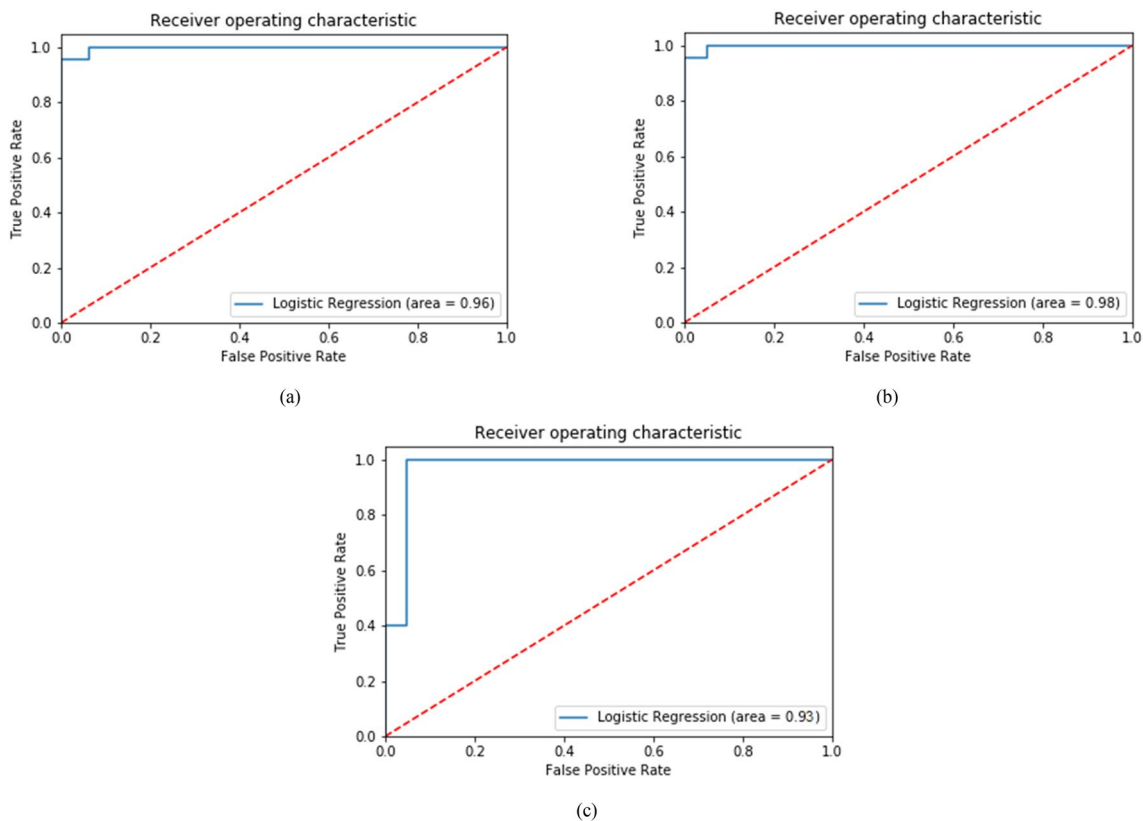


Fig. 8 ROC curves for logistic regression model when **a** variance=1, 0.99, 0.98, 0.96, 0.95 and 0.90, **b** variance=0.97, **c** variance=0.85

Table 7 Performance evaluation of linear regression model

Variance retained	LR								
	Precision			Recall			F-measure		
	Corona	Normal	Overall	Corona	Normal	Overall	Corona	Normal	Overall
1	0.91	1	0.955	1	0.913	0.957	0.952	0.955	0.945
0.99	0.91	1	0.955	1	0.913	0.957	0.952	0.955	0.945
0.98	0.91	1	0.955	1	0.913	0.957	0.952	0.955	0.945
0.97	0.95	1	0.976	1	0.955	0.977	0.975	0.977	0.976
0.96	0.91	1	0.955	1	0.913	0.957	0.952	0.955	0.945
0.95	0.91	1	0.955	1	0.913	0.957	0.952	0.955	0.945
0.90	0.91	1	0.955	1	0.913	0.957	0.952	0.955	0.945
0.85	0.90	0.95	0.930	0.95	0.913	0.932	0.927	0.933	0.930

Table 8 Classification accuracy and computational time to fit proposed models after PCA with different fractions of variance retained

Variance retained	No. of components	Time (ms)		Accuracy (%)	
		CNN	LR	CNN	LR
1	786432	858000	17600	97.6	95.2
0.99	147	233	1527	100	95.2
0.98	126	230	1521	97.6	95.2
0.97	108	217	1519	97.6	97.6
0.96	96	219	1508	97.6	95.2
0.95	84	216	1508	97.6	95.2
0.90	48	211	1499	97.6	95.2
0.85	30	211	1493	95.2	92.8

model for COVID-19 screening among CT images. They used 1262 CT images of COVID-19-positive cases and 1230 CT images of normal patients to secure an accuracy of 96.25% while F1 score of 96.29% with addition of data augmentation technique. Similarly, Sharma [40] designed a customized CNN-based ResNet50 architect as an automated tool for coronavirus disease diagnosis using 2200 computed tomography images of lungs (800 belongs to COVID-19 positive while 1400 of other pneumonia and normal patients) that attained 91.0% accuracy with sensitivity of 92.1%.

A SqueezeNet demonstrated by Ucar and Korkmaz [41] is tuned with Bayesian optimization additive for COVID-19 diagnosis. They gained an overall accuracy of 98.30% on dataset of 5949 chest radiography images that consists of 76 COVID-19 infected cases. Bai et al. [42] incorporated long short term memory (LSTM) with multi-layer perceptron (MLP) for classification of COVID-19 disease. They used 133 samples for training the proposed model and achieved 89.10% performance rate.

In this research, researchers analyzed the effect of PCA on proposed DL-based CNN model and LR network for detection of COVID-19 in X-ray images. For experimental work, 500 X-ray images are used, among which 250 images corresponds to patients affected by COVID-19 and 250 belongs to normal patients. By practicing PCA as feature extraction technique with proposed CNN model, we accomplished 100% accuracy with less computational time. While same extracted features when fed to LR network, model secured an accuracy of 97.6%. Evidently, our proposed method (PCA + CNN) outperformed other studies mentioned in Table 9.

Due to limited COVID-19 chest radiography images available publicly, which is constantly updating, most of the studies mentioned above either used 25, 50 or around 100 images related to COVID-19-affected patients in their

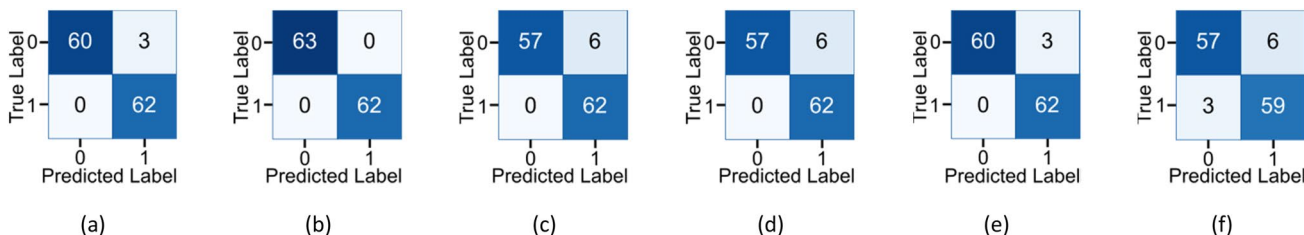


Fig. 9 Confusion matrix for **a** CNN model at variance of 1, 0.98, 0.97, 0.96, 0.95, and 0.90, **b** CNN model at variance of 0.99, **c** CNN model at variance of 0.85, **d** LR model at variance of 1, 0.99, 0.98, 0.96, 0.95 and 0.90, **e** LR model at variance of 0.97, **f** LR model at variance of 0.85

Table 9 Comparison of current AI diagnostic techniques with proposed methods using different medical radiology images

Ref.	Design/technique title	Images Type	Accuracy	Precision	Sensitivity	F1-score
[35]	DL frameworks (SqueezeNet, and MobileNetV2) with additional SMO method and Support Vector Machine classifier	X-ray	99.2	98.8	98.3	98.5
[36]	Five various ML-based classifiers with iRF feature extraction	X-ray	99.6	–	98.8	–
[37]	Pre-trained DL-based frameworks (Inception ResNetV2, Inception, VGG19, Xception, and MobileNet v2)	X-ray	97.4	–	99.1	–
[38]	VGG-16	X-ray	97.0	–	92.0	92.0
[39]	Deep CNN with transfer learning based DenseNet201	CT	96.2	96.2	96.2	96.2
[40]	Customized CNN-based ResNet50	CT	91.0	–	92.1	–
[41]	Deep Bayes-SqueezeNet	X-ray	98.3	–	–	98.3
[43]	Convolutional CAPSNET	X-ray	97.2	97.0	97.4	97.2
[44]	CSSA-EfficientNet-B0 based on deep learning 2D curvelet transform	X-ray	99.6	99.6	99.4	99.5
Proposed	PCA + CNN	X-ray	100	100	100	100

dataset for developing artificial intelligence-based classification system. However, in this proposed study, we used 198 images of COVID-19-positive cases, and employed GAN data augmentation technique to avoid overfitting. At large, studies conducted prior to this study generally used ResNet50 as classification technique without implementing any feature extraction method, while this study is based on PCA as feature extraction technique that drastically reduces the computational time as well as enhances the performance of proposed CNN model.

5 Conclusion

The study inspected the implementation of machine learning methods for identification of patients affected by COVID-19 with the help of X-ray images. Suggested approach is especially useful to quickly identify the patients so that necessary medical care can be provided in timely manner. In addition, it provides an alternative cheap solution for the situation in developing countries where testing kits are not available or expensive to conduct standard testing at mass scale. Data augmentation technique using GAN was also employed for increasing the amount of the data and to maximize the classification accuracy of proposed classifiers by reducing the risk of overfitting. The main advantage of the DL method was to reduce the do-it-yourself characteristics which are painstaking as well as labor-intensive, thus it improved the classification accuracy based on data-driven feature learning approach. The proposed method achieved high accuracy of 100% using CNN + PCA when variance of 0.99 was used. Moreover, the training time was also drastically reduced with the incorporation of PCA while achieving 100% accuracy on testing data. The proposed CNN + PCA technique eclipsed the ultra-modern and advanced approaches as it attained highest accuracy.

The researchers of this study plan to coach the network on relatively substantial dataset and then apply the network to improve the reliability of decision and overall accuracy. Moreover, the architecture of the CNN network can be made more complex by adding additional layers, for larger data set, to enhance the system capabilities of learning highly abstract features, thus making it less prone to overfitting. In addition to the image data, other features such as body temperature, information about presence of chronic diseases such as diabetes, heart etc. will be integrated with the system to make it more robust and reliable to support healthcare actioners. The proposed work can also be adopted for other medical applications such as breast cancer detection, tumor detection etc.

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

Compliance with ethical standards

Conflict of interest The authors declare that they have no conflict of interest.

Informed consent For this type of study, formal consent is not required.

Human and animal rights This paper does not contain any studies with human participants or animals performed by any of the authors.

References

1. Tyrrell DA, Bynoe M (1966) Cultivation of viruses from a high proportion of patients with colds. *Lancet* 287:76–77. [https://doi.org/10.1016/S0140-6736\(66\)92364-6](https://doi.org/10.1016/S0140-6736(66)92364-6)
2. Kahn JS, McIntosh K (2005) History and recent advances in coronavirus discovery. *Pediatr Infect Dis J* 24:S223–S227. <https://doi.org/10.1097/01.inf.0000188166.17324.60>

3. Jain V, Yuan J-M (2020) Predictive symptoms and comorbidities for severe COVID-19 and intensive care unit admission: a systematic review and meta-analysis. *Int J Public Health* 65:533–546. <https://doi.org/10.1007/s00038-020-01390-7>
4. Ren Y, Li L, Jia Y (2020) New method to reduce COVID-19 transmission: the need for medical air disinfection is now. *J Med Syst* 44:119. <https://doi.org/10.1007/s10916-020-01585-8>
5. Fisher D, Heymann D (2020) Q and A: the novel coronavirus outbreak causing COVID-19. *BMC Med* 18:57. <https://doi.org/10.1186/s12916-020-01533-w>
6. Seshadri DR, Davies EV, Harlow ER, Hsu JJ, Knighton SC, Walker TA, Voos JE, Drummond CK (2020) Wearable sensors for COVID-19: a call to action to harness our digital infrastructure for remote patient monitoring and virtual assessments. *Front Digit Health* 2:8. <https://doi.org/10.3389/fgdth.2020.00008>
7. Yahyaoui A, Jamil A, Rasheed J, Yesiltepe M (2019) A decision support system for diabetes prediction using machine learning and deep learning techniques. In: 2019 1st international informatics and software engineering conference (UBMYK). IEEE. <https://doi.org/10.1109/UBMYK48245.2019.8965556>
8. Nallamala SH, Mishra P, Koneru SV (2019) Breast cancer detection using machine learning way. *Int J Recent Technol Eng* 8:1402–1405
9. Gurbina M, Lascu M, Lascu D (2019) Tumor detection and classification of MRI brain image using different wavelet transforms and support vector machines. 2019 42nd international conference on telecommunications and signal processing. TSP. <https://doi.org/10.1109/TSP.2019.8769040>
10. Sun P, Wang D, Mok VC, Shi L (2019) Comparison of feature selection methods and machine learning classifiers for radiomics analysis in glioma grading. *IEEE Access* 7:102010–102020. <https://doi.org/10.1109/access.2019.2928975>
11. de Bruijne M (2016) Machine learning approaches in medical image analysis: From detection to diagnosis. *Med Image Anal* 33:94–97. <https://doi.org/10.1016/j.media.2016.06.032>
12. Nanni L, Ghidoni S, Brahnam S (2017) Handcrafted vs non-handcrafted features for computer vision classification. *Pattern Recognit* 71:158–172. <https://doi.org/10.1016/j.patcog.2017.05.025>
13. Ayan E, Ünver HM (2019) Diagnosis of pneumonia from chest X-ray images using deep learning. 2019 scientific meeting on electrical-electronics and biomedical engineering and computer science. EBBT. <https://doi.org/10.1109/EBBT.2019.8741582>
14. Huynh HT, Anh VNN (2019) A deep learning method for lung segmentation on large size chest x-ray image. RIVF 2019. *Proceed IEEE-RIVF Int Conf Comput Commun Technol*. <https://doi.org/10.1109/RIVF.2019.8713648>
15. Van Tulder G, De Bruijne M (2016) Combining generative and discriminative representation learning for lung CT analysis with convolutional restricted Boltzmann machines. *IEEE Trans Med Imaging* 35:1262–1272. <https://doi.org/10.1109/TMI.2016.2526687>
16. Bardou D, Zhang K, Ahmad SM (2018) Classification of Breast Cancer Based on Histology Images Using Convolutional Neural Networks. *IEEE Access* 6:24680–24693. <https://doi.org/10.1109/ACCESS.2018.2831280>
17. Shin H, Roth HR, Gao M, Lu L, Xu Z, Nogues I, Yao J, Mollura D, Summers RM (2016) Deep convolutional neural networks for computer-aided detection: CNN architectures, dataset characteristics and transfer learning. *IEEE Trans Med Imaging* 35:1285–1298. <https://doi.org/10.1109/TMI.2016.2528162>
18. Zhou X.-Y, Zheng J.-Q, Li P, Yang G.-Z. (2020) ACNN: a Full Resolution DCNN for Medical Image Segmentation. In: 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE <https://doi.org/10.1109/ICRA40945.2020.9197328>.
19. Sudheer Kumar E, Shoba Bindu C (2019) Medical image analysis using deep learning: a systematic literature review. *Commu Comput Informat Sci*. https://doi.org/10.1007/978-981-13-8300-7_8
20. Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg AC, Fei-Fei L (2015) ImageNet large scale visual recognition challenge. *Int J Comput Vision (IJCV)* 115:211–252. <https://doi.org/10.1007/s11263-015-0816-y>
21. Salehinejad H, Colak E, Dowdell T, Barfett J, Valaee S (2019) Synthesizing chest X-Ray pathology for training deep convolutional neural networks. *IEEE Trans Med Imaging* 38:1197–1206. <https://doi.org/10.1109/TMI.2018.2881415>
22. Shorten C, Khoshgoftaar TM (2019) A survey on image data augmentation for deep learning. *Journal of Big Data* 6:60. <https://doi.org/10.1186/s40537-019-0197-0>
23. Zhang C, Song Y, Liu S, Lill S, Wang C, Tang Z, You Y, Gao Y, Klistorner A, Barnett M, Cai W (2018) MS-GAN: GAN-based semantic segmentation of multiple sclerosis lesions in brain magnetic resonance imaging. In: 2018digital image computing. Tech Appl (DICTA). <https://doi.org/10.1109/DICTA.2018.8615771>
24. Oktay O, Ferrante E, Kamnitsas K, Heinrich M, Bai W, Caballero J, Cook SA, de Marvao A, Dawes T, O'Regan DP, Kainz B, Glocker B, Rueckert D (2018) Anatomically constrained neural networks (ACNNs): application to cardiac image enhancement and segmentation. *IEEE Trans Med Imaging* 37:384–395. <https://doi.org/10.1109/TMI.2017.2743464>
25. Schlemper J, Caballero J, Hajnal JV, Price AN, Rueckert D (2018) A deep cascade of convolutional neural networks for dynamic MR image reconstruction. *IEEE Trans Med Imaging* 37:491–503. <https://doi.org/10.1109/TMI.2017.2760978>
26. Albarqouni S, Baur C, Achilles F, Belagiannis V, Demirci S, Navab N (2016) AggNet: deep learning from crowds for mitosis detection in breast cancer histology images. *IEEE Trans Med Imaging* 35:1313–1321. <https://doi.org/10.1109/TMI.2016.2528120>
27. Pereira S, Pinto A, Alves V, Silva CA (2016) Brain tumor segmentation using convolutional neural networks in MRI images. *IEEE Trans Med Imaging* 35:1240–1251. <https://doi.org/10.1109/TMI.2016.2538465>
28. Yang X, Fan D, Ren A, Zhao N, Alam M (2019) 5G-based user-centric sensing at C-band. *IEEE Trans Industr Inf* 15:3040–3047. <https://doi.org/10.1109/TII.2019.2891738>
29. Yang X, Shah SA, Ren A, Zhao N, Fan D, Hu F, Ur Rehman M, von Deneen KM, Tian J (2018) Wandering pattern sensing at S-band. *IEEE J Biomed Health Inform* 22:1863–1870. <https://doi.org/10.1109/JBHI.2017.2787595>
30. Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y (2014) Generative adversarial nets. In: *proceedings of the 27th international conference on neural information processing systems* 2, 2672–2680. MIT Press, Cambridge, MA, USA
31. Abdi H, Williams LJ (2010) Principal component analysis. *Wiley Interdiscipl Rev Comput Statist* 2:433–459. <https://doi.org/10.1002/wics.101>
32. Kroonenberg PM, de Leeuw J (1980) Principal component analysis of three-mode data by means of alternating least squares algorithms. *Psychometrika* 45:69–97. <https://doi.org/10.1007/BF02293599>
33. Albawi S, Mohammed TA, Al-Zawi S (2017) Understanding of a convolutional neural network. *Int Conf Eng Technol (ICET)*. <https://doi.org/10.1109/ICEngTechnol.2017.8308186>
34. Peng C-YJ, Lee KL, Ingersoll GM (2002) An introduction to logistic regression analysis and reporting. *J Educ Res* 96:3–14. <https://doi.org/10.1080/00220670209598786>
35. Toğaçar M, Ergen B, Cömert Z (2020) COVID-19 detection using deep learning models to exploit social mimic optimization and

- structured chest X-ray images using fuzzy color and stacking approaches. *Comput Biol Med.* <https://doi.org/10.1016/j.compbiomed.2020.103805>
36. Tuncer T, Dogan S, Ozyurt F (2020) An automated residual exemplar local binary pattern and iterative relief based corona detection method using lung X-ray image. *Chemomet Intell Lab Syst* 203:104054. <https://doi.org/10.1016/j.chemolab.2020.104054>
 37. Apostolopoulos ID, Mpesiana TA (2020) Covid-19: automatic detection from X-ray images utilizing transfer learning with convolutional neural networks. *Phys Eng Sci Med* 43:635–640. <https://doi.org/10.1007/s13246-020-00865-4>
 38. Brunese L, Mercaldo F, Reginelli A, Santone A (2020) Explainable deep learning for pulmonary disease and coronavirus COVID-19 detection from X-rays. *Comput Methods Programs Biomed* 196:105608. <https://doi.org/10.1016/j.cmpb.2020.105608>
 39. Jaiswal A, Gianchandani N, Singh D, Kumar V, Kaur M (2020) Classification of the COVID-19 infected patients using DenseNet201 based deep transfer learning. *J Biomol Struct Dyn.* <https://doi.org/10.1080/07391102.2020.1788642>
 40. Sharma S (2020) Drawing insights from COVID-19-infected patients using CT scan images and machine learning techniques: a study on 200 patients. *Environ Sci Pollut Res.* <https://doi.org/10.1007/s11356-020-10133-3>
 41. Ucar F, Korkmaz D (2020) COVIDiagnosis-Net: deep bayes-SqueezeNet based diagnosis of the coronavirus disease 2019 (COVID-19) from X-ray images. *Med Hypotheses* 140:109761. <https://doi.org/10.1016/j.mehy.2020.109761>
 42. Bai X, Fang C, Zhou Y, Bai S, Liu Z, Xia L, Chen Q, Xu Y, Xia T, Gong S, Xie X, Song D, Du R, Zhou C, Chen C, Nie D, Qin L, Chen W (2020) Predicting COVID-19 malignant progression with AI techniques. *SSRN Elect J.* <https://doi.org/10.2139/ssrn.3557984>
 43. Toraman S, Alakus TB, Turkoglu I (2020) Convolutional capnet: A novel artificial neural network approach to detect COVID-19 disease from X-ray images using capsule networks. *Chaos Solitons Fractals.* 140:110122. <https://doi.org/10.1016/j.chaos.2020.110122>
 44. Altan A, Karasu S (2020) Recognition of COVID-19 disease from X-ray images by hybrid model consisting of 2D curvelet transform, chaotic salp swarm algorithm and deep learning technique. *Chaos, Solitons Fractals.* <https://doi.org/10.1016/j.chaos.2020.110071>