

Lung Disease Prediction Using Deep Learning



Debasree Mitra, Pranati Rakshit, Anjali Jha, Dristi Dugar,
and Kamran Iqbal

Abstract The evolution of deep learning has enhanced the technique of identifying and classifying lung diseases into various categories using medical images. This project aims to build lung disease detection models using deep learning to identify future potential and thus to efficiently observe and visualize the recent and upcoming trends in this domain. Identifying and discovering lung disease at an early stage has become a vital part of the medical domain because this would facilitate patient's subsequent clinical management. The project primarily focuses on pneumonia as well as considering the breathing problems of patients. Deep learning and machine learning have served the utmost significance in detecting such lung diseases at a prior stage. This enhancement has contributed much to the doctors and medical systems to provide early treatment to patients. In this project, convolutional neural network (CNN) is used to predict lung disease (pneumonia) from chest X-ray images using machine learning and deep learning frameworks.

Keywords Lung disease prediction · Deep learning · Machine learning · Chest X-ray images · Convolutional neural networks (CNNs) · Pneumonia

1 Introduction

Medical X-ray images provide a gateway to diagnose body parts such as the chest, bones, skull and teeth. Medical pioneers have used this same approach for several decades to investigate and analyse fractures or abnormalities in body organs.

Machine learning and deep learning play a pivotal role in performing the algorithms accurately to categorize the chest X-ray peculiarities besides providing a large window for the further prediction of deadly lung diseases. Over the past few

D. Mitra (✉) · P. Rakshit · A. Jha · D. Dugar · K. Iqbal
Department of Computer Science and Engineering, JIS College of Engineering, MAKAUT,
Kolkata, India
e-mail: debasree.mitra2005@gmail.com

P. Rakshit
e-mail: pranati.rakshit@jiscollege.ac.in

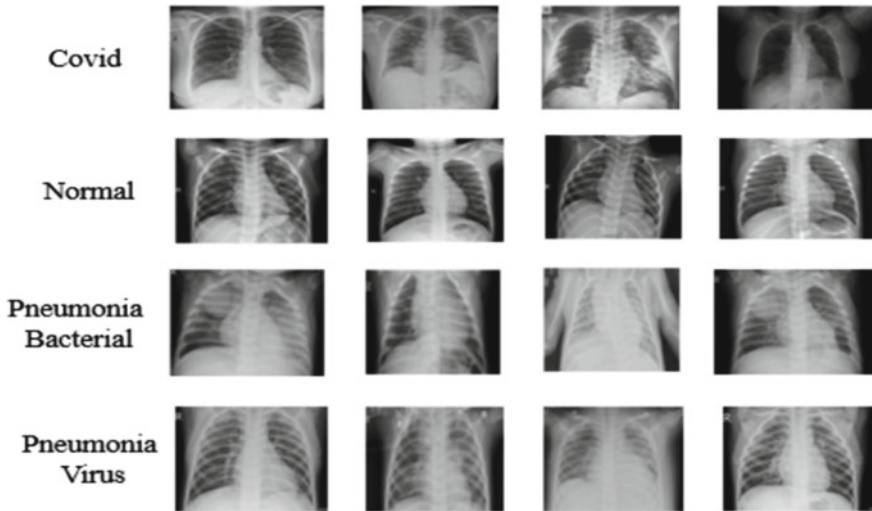


Fig. 1 Different types of lung diseases

years, digital technology has gained the utmost importance globally. Systems based on deep learning have been put into play to elevate the exactness of image classification. Deep networks showed extraordinary accuracies in such predictions. This success led the researchers to enforce the deep networks to chest X-ray images for lung disease prediction, and the classification paving their way out showed these deep networks masterfully extract beneficial features that distinguish different categories of images provided to the model. The most conventional deep learning framework used is the convolutional neural network (CNN). CNN provides a unique power to extract features from images at different levels.

This research paper may provide other doctors and researchers the right path for bringing to light the lung diseases such as pneumonia, asthma, tuberculosis, SARS, ARDS and COPD with the help of deep learning methodology as shown in Fig. 1 [1]. Several works regarding the research have been carried out on the diagnosis of lung diseases using artificial intelligence methodologies. Neural networks and artificial immune systems are already being used to diagnose fatal diseases like pneumonia and chronic obstructive pulmonary disease (COPD).

2 Literature Survey

The most severe chronic lung diseases involve chronic obstructive pulmonary diseases (COPD), pulmonary sarcoidosis, asthma, pneumonia and interstitial lung diseases, India, holding a second rank in terms of population with 1.3 billion people

residing in 29 states and seven union territories often vary widely in terms of demography, economy and ecology, through which respiratory health gets affected [2]. Table 1 [3] shows the ten most common causes of death in 2008. In India, the estimation indicates the cases of pneumonia have reduced approximately by 41% between the years 2000 and 2015, which in turn is greater than 22% reduction in pneumonia cases observed globally during the same period. Nevertheless, pneumonia still remains a crucial cause of morbidity among children country-wise [4].

Numerous researches have faced the difficulty of classifying the images through such a high degree of accuracy. Here are a few citations from our paper regarding the CNN model.

Rubin et al. [5] created a model to identify basic thoracic illness in frontal chest X-ray images. To accomplish large-scale automatic recognition of these pictures, the MIMIC-CXR dataset was employed. The dataset was divided into three sections: testing, training and validation, with a piece accounting for 20%, 70% and 10% of the entire lot. To improve overall performance, normalization of pixel was utilized.

Lakhani and Sundaram [6] developed such model to classify pulmonary tuberculosis. Chest X-ray images were also classified using transfer learning models such as AlexNet and GoogleNet. About 68%, 17.1% and 14.9% of the dataset were split into training set, validation sets and testing sets, respectively. To get the ultimate model with an AUC of 0.99, techniques like augmentation of data as well as pre-processing were used. The model’s recall and veracity were 97.3% and 100%, respectively.

Guan et al. [7] build an AG-CNN model to detect diseases related to the thorax. For the detection of such diseases from X-ray images of the chest, the chest X-ray14 dataset was applied. For categorization, a local and global branch attention-guided CNN was utilized. With an AUC of 0.868, their model outperformed the other models.

To classify chest X-ray images into pneumonia and 14 other diseases, Rajpurkar et al. [8] used chest X-ray14 dataset to train the model. They compared their ChXNet model (121-layer model) with academic radiologists in practice. Their model has an F1 score of 0.435, which was higher than radiologists’ F1 score of 0.387.

Table 1 Statistics of world health 2011

| Reasons for deaths | Worldwide data | WHO European region |
|------------------------------|---------------------|----------------------|
| Ischaemic heart disease | 7.3 million (12.8%) | 2.40 million (24.7%) |
| Cerebrovascular disease | 6.2 million (10.8%) | 1.40 million (14.0%) |
| Lower respiratory infections | 3.5 million (6.1%) | 0.23 million (2.3%) |
| COPD | 3.3 million (5.8%) | 0.25 million (2.5%) |
| Diarrhoeal diseases | 2.5 million (4.3%) | 0.03 million (0.3%) |
| HIV/AIDS | 1.8 million (3.1%) | 0.08 million (0.8%) |
| Trachea/bronchus/lung cancer | 1.4 million (2.4%) | 0.38 million (3.9%) |
| Tuberculosis | 1.3 million (2.4%) | 0.08 million (0.8%) |
| Diabetes mellitus | 1.3 million (2.2%) | 0.17 million (1.7%) |
| Road traffic accidents | 1.2 million (2.1%) | 0.12 million (1.2%) |

Krizhevsky et al. [9] trained this sort of model with five convolutional sheets, some of which were followed by max pooling layers and three fully connected layers. Such networks contained parameters over 60 million. This model earned an error rate of 17% while using dropout.

Simonyan and Zisserman [10] created a supreme model that used many tiny kernel-sized filters to reach an exactness of 92.7%. The ImageNet dataset was used to train this model and then proposed to the ILSVRC 2014 competition.

Xu et al. [11] created a convolutional neural network for brain tumour MRI classification and segmentation. This model implemented a plethora of techniques, such as data augmentation, feature selection and pooling. This model achieved a validation accuracy of 97.5% for classification and 84% for segmentation on 256 256 pixels sized frontal chest radiographs that were fed to a deep convolutional neural network to detect anomalies.

Anthimopoulos et al. [12] developed this replica with five convolution sheets, leaky ReLU, average pooling and three fully connected layers to detect interstitial lung disease patterns in a dataset comprising 14,696 images. This model had an accuracy of 85.5%.

3 Chest X-ray Data Analysis

The datasets used in the project are reported in this section. This is done to provide us with relevant information on the datasets. Dataset description and view positions of the chest X-ray images are discussed briefly in this section.

3.1 Dataset Description

As a matter of concern, public datasets have been chosen because of their availability to the public, whereas private datasets are not accessible without permission. Here, the model is exposed to two different datasets to protract a conclusion over the CNN model while visualizing the difference in results of pneumonia prediction.

- The first dataset contains 5933 images with resolutions varying from 1024×1024 to 2444×1800 .
- The second dataset contains a total of 5357 images with resolution varying from 1114×928 to 2290×2066 .

The entire dataset, one and two contain numerous classes of images with the name normal, pneumonia (virus and bacteria), syndromes such as severe acute respiratory (SARS), acute respiratory distress (ARDS), COVID-19 positive and negative, CT scan. For better visualization and accurate analysing the chest X-ray images, the diagnosis can be carried out in a more time-efficient manner by medical staff and

doctors to monitor the patient’s health accordingly. Moreover, by creating applications using intelligent machines, physicians could analyse the condition of patients better.

3.2 View Position

1. Posterior-anterior (PA) projection: The posterior-anterior (PA) projection in Fig. 2a is considered a standard position for determining matured chest radiograph as the light beams pass through the chest part of the torso from a posterior position to anterior position. The image was captured, while the patient was standing still against a film, holding his arms straight up in the air and having few deep breaths throughout the process to get good and clear pictures of both heart and lungs. These techniques ensure precise higher quality images shown in (b) and access the big heart size than the anterior–posterior (AP) position.
2. Anterior-posterior (AP) projection: The anteriorposterior (AP) projection in Fig. 3a is less considered because the light beams are restricted to pass through the chest wall. This condition occurs when the patients are unwell or too weak to even stand straight against a film. Hence, the projection is passed through the anterior to posterior position. This results in a lower quality image of the chest X-ray shown in (b), due to which capturing the figure of a clear and focussed heart becomes a tough job.

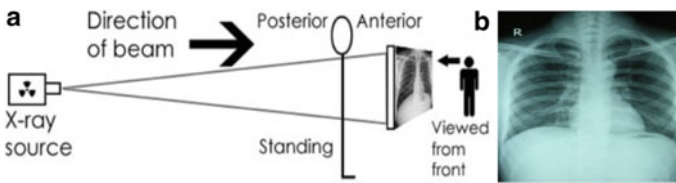


Fig. 2 a PA projection. b PA chest X-ray

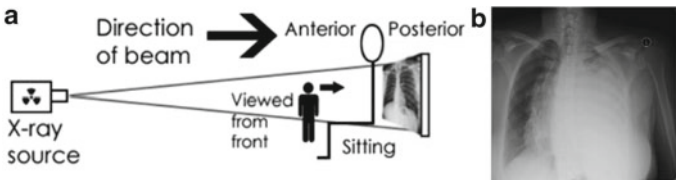


Fig. 3 a AP Projection. b AP chest X-ray

A comparative study demonstrates that it is pretty easy to find lung diseases in middle aged to elderly humans comparable to young ones. Any vague disease found in young people is easy to diagnose through primary treatments, nonetheless, the discovery of high-rated infections in middle-aged groups is a matter of concern [13, 14].

4 Description of Algorithm Used

Convolutional neural networks (CNNs) were primarily introduced by Haffner, Bengio, Bottou and Le Cun, in the late 1998s. The arrangement and anatomy of a mammal brain have played an important role towards inspiring the evolution of algorithms of machine learning, one of them being convolutional neural network (CNNs) in Fig. 4 [15], an unsupervised single layered that uses a hierarchical approach to build a neural network, a funnel-like structure that further results in a network where neurons are connected to form a fully connected network leading to training a model layer by layer to process the output. Several procedures need to be followed while classifying an image using CNN.

Firstly, a filter or kernel is made to slide over the input data (layer) pixel by pixel, the stride is considered to be one(by default) and multiplication of matrices is performed to obtain a feature map (matrix). The size of the filter may vary. This process is carried out on each convolution until a final feature map is achieved.

Secondly, the images are resized before the model undergoes training. This is done with the help of pooling and flattening layers. Pooling refers to continuous reduction in dimensionality to eliminate the computation in the network and the number of parameters. This leads to handling overfitting with less training time. Max pooling is the most frequent type which calculates the maximum value in each and every window. Hence, decreasing the size of the feature map while retaining the chief information. Flattening helps in converting original data into a 1-D array which would further be used as an input layer in the next step. This process is carried out to

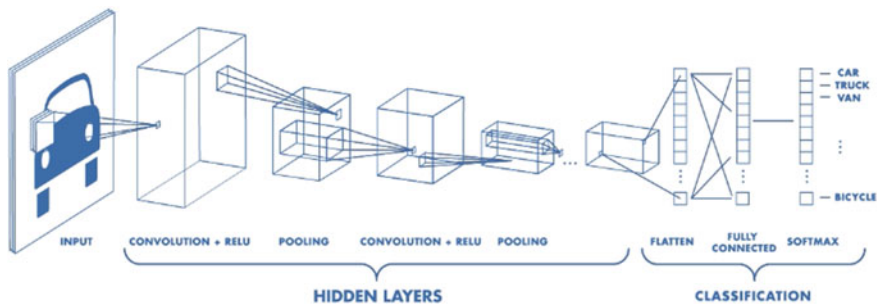


Fig. 4 Architecture of a CNN

extract significant features. Flattening assists in getting an output of fully connected neural networks after the complete convolution is done.

Lastly, to prevent overfitting of the model, batch normalization is considered to be the most efficient as it regularizes the output of the previous layers which in turn makes the output layers learn independently. Here, activation functions highly contribute when it comes to scale the input. Most widely used activation functions in neural networks are ReLu, Leaky ReLu, softmax Function, sigmoid function and so on. The choice of using a particular activation function depends on the features provided by a particular activation function. In this paper, the activation functions used are rectified linear unit (ReLu) and softmax activation. The advantage of implementing ReLu over any activation function is that it prevents the problem of vanishing gradient. The former function does not activate all the neurons at the same time, hence, it performs the nonlinear transformation of the input data. Hence, the output layer achieved is a fully connected layer. Every neuron of the latest layer is attached to all the neurons in previous film, thus forming a densely connected layer.

5 Methodology

5.1 Data (Image) Pre-processing Phase

To detect any lung disease, image dataset is considered to be the most relevant when it comes to using deep learning. These include medical photocopies of CT scans, X-ray images and many more. These medical images are available in the raw form which may not be understandable by machines. So, in order to make it machine readable, these images need to be pre-processed before moving on towards training.

In order to achieve a high training phase, these X-ray images are rescaled and the size is reduced by 255 NumPy array matrices. The images are then transformed from greyscale image to RGB for training purposes. Furthermore, various modifications are made on the resultant images and after carrying out data augmentation on this dataset to gain more relevant data out of it. Unnecessary data are removed, and the complete dataset is filtered and pre-processed. Feature extraction is then performed on this processed data to draw out the most cardinal features so as to prepare the data finally to be trained on the chosen model.

5.2 Training Phase

Since, in this paper, the dataset used for predicting lung disease is X-rays and CT scan images, deep learning algorithms fit best when it comes to training image datasets. There are many deep learning algorithms available that perform differently for different datasets. To compute and train large amounts of data, deep learning

introduces many multiple layers of neural networks. Few well-known deep learning algorithms of neural networks are convolutional (CNNs), recurrent (RNNs), multi-layer perceptron (MLPs), deep belief (DBNs), long short-term memory (LSTMNs) and so on. These algorithms are chosen accordingly based on the types of data to be dealt with and which model fits best to give desired output to the given problem statement. When the dataset contains time and date series, RNNs and LSTMNs are the algorithms that are considered the most. RNNs remember previous output and use it as current input while having hidden layers in between them. LSTM networks are a type of RNNs that works on sequence prediction problems. MLP networks can be used to deal with data that are nonlinearly separable, and in computational neuroscience, it contains multiple layers of neurons, and the use of back-propagation is a major factor here. To deal with image and video sequences of data, deep belief networks are mostly preferred since it uses a greedy algorithm approach. CNNs are used to get acquainted with image datasets as it takes tensors as input, and after passing filters over multiple layers of neurons and hidden layers, it succeeds in finding relationships between each pixel in the images. The datasets used in this paper are grey-scaled images that minutely need to be spatially connected among the pixels. Hence, CNN is used in this regard for generating the best confusion matrix. During training, a technique named dropout selects neurons randomly and drops them, thus regularizing deep neural networks. There is a fundamental difference between application of dropout to fully connected layers and to convolutional layers. To avoid overfitting, regularization techniques are used to appropriately fit the model by reducing the error. With the intent of standardizing the input to a layer, batch normalization is applied to the deep learning network. This aims to put a stop to the overfitting of model by reducing the number of training epochs. The stabilization attained by the method early stopping targets to improve model performance by stopping the training and specifying the increased number of training epochs.

5.3 Classification Phase

The problem statement of this paper is to predict lung disease and to determine whether a patient is suffering from pneumonia or not. The chosen algorithm, convolutional neural network, is a classification algorithm that predicts the output classifying it, corresponding to which it gives high accuracy. The two classes are normal (the value assigned is '0') and pneumonia (the value assigned is '1'). The extracted features are arranged into layers, and according to the working principle of CNN, filters are supposed to slide over these input layers as well as some hidden layers until it reaches and becomes a fully connected network. Learning rate, being a configurable hyperparameter, determines the step size. It serves a paramount role in optimizing and training the model. If the learning rate of a neural network model is too high, then the loss function may show undesirable divergent behaviour of the loss function. A desirable learning rate could be a prime difference between a model that presents

accurate results and a model that does not learn. Hence, reduction of learning rate becomes a momentous step in optimizing the model's performance. Therefore, an ideal learning rate should be in the range of 0.1–0.001.

The model is then fitted and has used a certain number of epochs until the desired accuracy is reached. Then, the model is tested over a certain amount of testing data and checked for how the model performs after getting trained on the selected features. This process results in classifying images into the respective classes it belongs to. As an outcome, the classified images are tagged with accuracy percentage which supports their belonging to that particular class (normal or pneumonia).

6 Result and Discussion

The overall implementation of the proposed CNN model is tested and presented in this section.

6.1 *Model Validation and Evaluation*

After training our model, a validation set was used for estimation purposes. The performance of the model was visualized by plotting various graphs to observe the general behaviour that our model learnt. The graphs between training loss versus training validation loss and training accuracy versus validation accuracy are represented for both datasets. The training of the model was performed in two ways on two different datasets.

The first, without the use of data augmentation. Figure 5a, b depicts the graphs without the use of data augmentation for dataset 1; (c) and (d) show the graphs with the use of data augmentation for dataset 2.

The second, with the use of data augmentation; (e) and (f) depict the graphs with the use of data augmentation for the dataset 1; (g) and (h) depict the graphs with the use of data augmentation for the dataset 2.

6.2 *Justification*

For training and validation of the model, two different datasets (chest X-rays images) were taken into consideration to demonstrate a comparative study which ensued in a better visualization on how the chosen algorithm (CNN) works. The model trained without data augmentation has resulted into overfitting in Table 2.

This model when introduced with data augmentation as it is about fabricating more data from the data we actually got, adding variance without losing the information the data carries. Doing this reduces the risk of overfitting, and generally, the accuracy on unseen data can be improved [16] in Table 3.

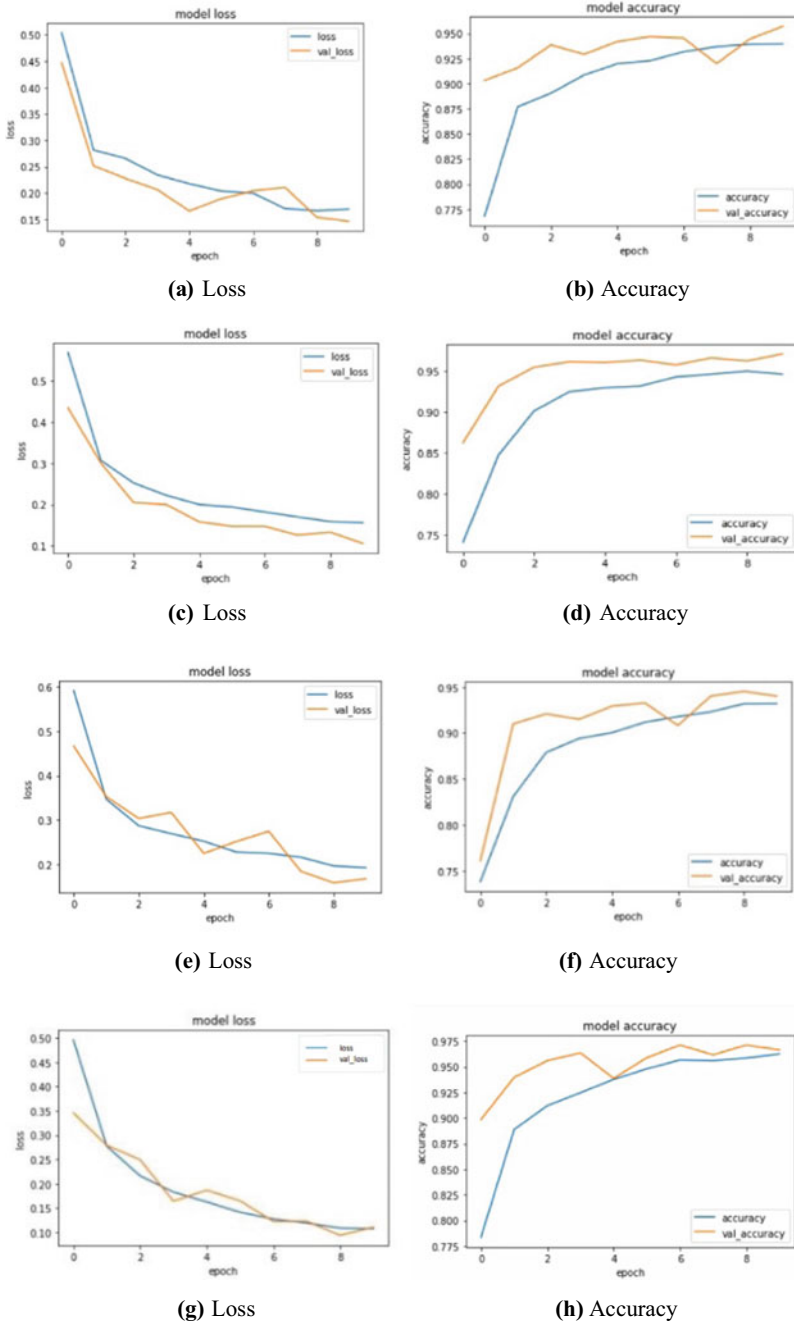


Fig. 5 a Loss. b Accuracy. c Loss. d Accuracy. e Loss. f Accuracy. g Loss. h Accuracy

Table 2 Model performance without data augmentation

| S. No. | Name | Train accuracy (%) | Train loss (%) | Test accuracy (%) | Test loss (%) | Validation accuracy (%) | Validation loss (%) |
|--------|-----------|--------------------|----------------|-------------------|---------------|-------------------------|---------------------|
| 1 | Dataset 1 | 94.15 | 15.85 | 95.71 | 14.6 | 95.71 | 14.60 |
| 2 | Dataset 2 | 94.65 | 15.77 | 97.03 | 10.5 | 97.03 | 10.57 |

Table 3 Performance of the model with data augmentation

| S. No. | Name | Train accuracy (%) | Train loss (%) | Test accuracy (%) | Test loss (%) | Validation accuracy (%) | Validation loss (%) |
|--------|-----------|--------------------|----------------|-------------------|---------------|-------------------------|---------------------|
| 1 | Dataset 1 | 92 | 19.90 | 95.74 | 14.67 | 95.75 | 14.67 |
| 2 | Dataset 2 | 96.07 | 10.90 | 96.64 | 11.05 | 96.65 | 11.05 |

6.3 Free-Form Visualization

In this paper, to evaluate the trained model, we tested it on five hundred random images from the test dataset. The rationale behind is to interrogate how our model performs when introduced to unseen dataset (images). The final results from both the datasets are below as Fig. 6a–d.

7 Conclusion

In the current paper, pre-training and the detection of lung disease became a great aid using convolutional neural network (CNNs). CNNs have three layers designed for the ease of detecting diseases at medical diagnostics. Hence, sliding filters over the input multiple times make the process worth determining the results earlier. According to the research, a lot of papers have been issued on lung diseases. The emergence of various lung diseases such as asthma, COPD and infections like influenza, pneumonia, lung cancer and so on has made it essential for the improvement and amplification of the technology and designing such early disease prediction models using deep learning.

TensorFlow and Keras have also played an influential role in supporting backend of the CNN model. Since, raw medical images (X-ray images) contain a lot of trivial data; hence, best algorithms were used for feature extraction and data pre-processing. The use of two datasets gave proper understanding of how the results vary on a large scale such as hospitals. The choice of suitable algorithm made the task more efficient as CNN is considered to be one of the best algorithms that targets to achieve best results in detecting lung diseases. Within less iterations, CNN accomplished a good precision model also with less number of epochs. The results so far obtained

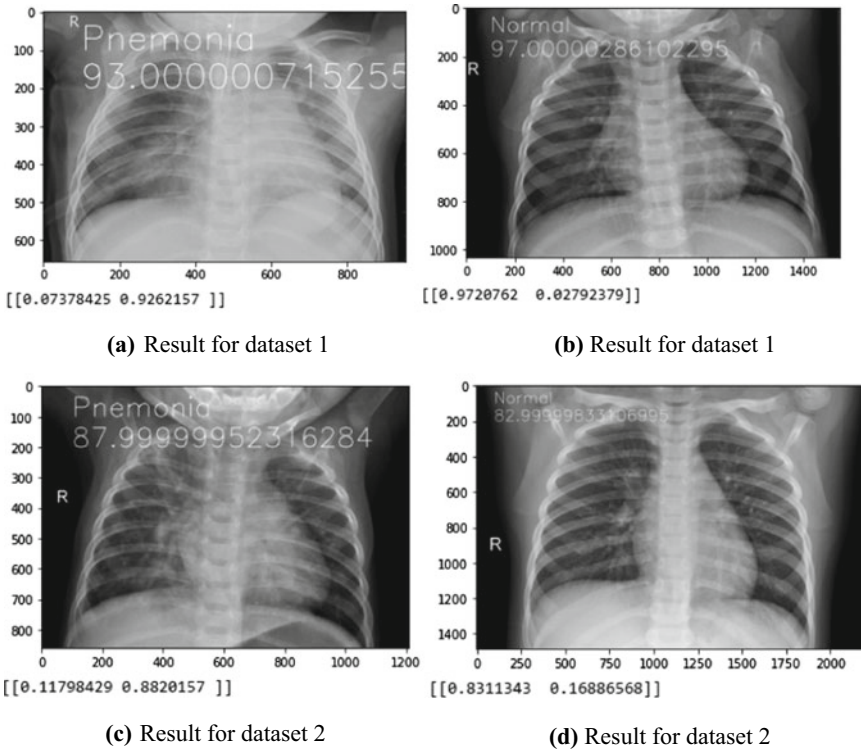


Fig. 6 **a** Result for dataset 1. **b** Result for dataset 1. **c** Result for dataset 2. **d** Result for dataset 2

demonstrate a fair performance of the CNN model used for this project. So, as a means to make the proposed CNN model through this paper to be actively used in hospitals and other medical diagnostics centres, further improvements can be made to achieve the accuracy to reach a peak point and also to detect lung diseases at a more precise level. This would serve a lot and is definitely going to make a difference in early diagnosis of lung diseases through its applications.

References

1. <https://medium.com/analytics-vidhya/classification-of-chest-xrays-using-pytorch-d50edd9ebb0>
2. [https://www.thelancet.com/journals/langlo/article/PIIS2214-109X\(18\)30409-1/fulltext](https://www.thelancet.com/journals/langlo/article/PIIS2214-109X(18)30409-1/fulltext)
3. <https://www.erswhitebook.org/chapters/the-burden-of-lung-disease/>
4. [https://www.thelancet.com/pdfs/journals/lanchi/PIIS2352-4642\(20\)30129-2.pdf](https://www.thelancet.com/pdfs/journals/lanchi/PIIS2352-4642(20)30129-2.pdf)
5. Rubin J, Sanghavi D, Zhao C, Lee K, Qadir A, Xu-Wilson M (2018) Large scale automated reading of frontal and lateral chest X-Rays using dual convolutional neural networks. arXiv preprint [arXiv:1804.07839](https://arxiv.org/abs/1804.07839)

6. Lakhani P, Sundaram B (2017) Deep learning at chest radiography: automated classification of pulmonary tuberculosis by using convolutional neural networks. *Radiology* 284(2):574–582
7. Guan Q, Huang Y, Zhong Z, Zheng Z, Zheng L, Yang Y (2018) Diagnose like a radiologist: attention guided convolutional neural network for thorax disease classification. arXiv preprint [arXiv:1801.09927](https://arxiv.org/abs/1801.09927)
8. Rajpurkar P, Irvin J, Zhu K, Yang B, Mehta H, Duan T, Ding D, Bagul A, Langlotz C, Shpan-skaya K, Lungren MP (2017) Chexnet: radiologist-level pneumonia detection on chest X-rays with deep learning. arXiv preprint [arXiv:1711.05225](https://arxiv.org/abs/1711.05225)
9. Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*, pp 1097–1105
10. Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. arXiv preprint [arXiv:1409.1556](https://arxiv.org/abs/1409.1556)
11. Xu Y, Jia Z, Ai Y, Zhang F, Lai M, Eric I, Chang C (2015) Deep convolutional activation features for large scale brain tumor histopathology image classification and segmentation. In: *2015 International conference on acoustics, speech and signal processing (ICASSP)*, pp 947–951
12. Anthimopoulos M, Christodoulidis S, Ebner L, Christe A, Mougiakakou S (2016) Lung pattern classification for interstitial lung diseases using a deep convolutional neural network. *IEEE Trans Med Imaging* 35(5):1207–1216
13. https://www.physio-pedia.com/Chest_X-Rays
14. https://www.radiologymasterclass.co.uk/tutorials/chest/chest_quality/chest_xray_quality_projection
15. <https://medium.com/@RaghavPrabhu/understanding-of-convolutional-neural-network-cnn-deep-learning-99760835f148>
16. <https://towardsdatascience.com/how-i-got-1-better-accuracy-by-data-augmentation-2475c509349a>