

# **COVID-19 Detection Using Radiology - CXR Images**

*Author: Raghwendra Dey, IIT Kharagpur*

## **Introduction:**

Since December, 2019 novel-coronavirus has wreaked havoc on the whole of mankind. Since then, it has approximately 4,787,666 confirmed cases, and 315,982 deaths as of 18th May, 2020, with 90,893 deaths in the United States alone. It was officially declared to be a global pandemic by the World Health Organization (W.H.O.) on 11 March, 2020. Experts say, until the vaccine of coronavirus is developed it can only be controlled by social distancing and extensive testing and isolating affected people, but the problem with testing of coronavirus is that its testing(PCR test) requires sophisticated labs and at least 12 hrs of time to get the results by the *swab test*. Even though rapid test kits have been developed but it is still way behind in production than its need and also is yet unavailable to the third world countries. Amidst this pandemic, Deep Learning shows rays of hope in helping to fight the virus. It's still a constant area of research to find ways to speed up and cheapen down the testing process, and screening through deep learning models by using the chest X-ray image of the subject seems promising (at least promising than the present methods of screening through temperature measurements, etc). Chest X-ray is not only cheap but X-ray machines are also widely available across the world.

## **About:**

This project was a part of the annual project for the course *Deep Learning Foundations and Applications* (spring-2020), IIT Kharagpur. A [dataset](#) was given which is a collection of many dataset sources. Task was to classify subjects into *Non-Pneumonia*, *Other Pneumonia* and *COVID-19*, by using Chest X-rays images of the subjects.

## **About Dataset:**

Dataset contains chest X-ray images from different sources. You can read more about the particulars of the dataset [here](#). I didn't use the Source-3 dataset due to data restrictions and also Source-3 mainly contained images of non-Covid cases which is

already outnumbered than the covid confirmed ones. The dataset was highly class imbalanced. To give you a hint The initial dataset with Source-3 data contained:

Disease	Number of Examples
Non-Pneumonia	199958
Other Pneumonia	9965
COVID-19	245

Even after not using Source-3 as well Number of Non-pneumonia and Other pneumonia remained as high as 27000 and 6000(approx.), but the number of covid examples remained the same i.e. 245. Also the dataset contained many types of noises like the unknown samples were marked with a '-1' in the first and the second column, for the Source-6 images, the names in the *Test\_Combined.csv* file didn't contained the extensions of the images unlike, other sources, some images was of single channel, unlike most of three channel images, etc, to name a few. Firstly, I need to clean all this noise off the dataset. In Fact, the dataset cleaning was one of the main hurdles to overcome.

### **About Model and Training:**

I first tried the dataset on a self made LeNet model with three conv layers, and two fully connected layers on top of it. It gave a test set accuracy of about 64.9%. Then, I tried transfer learning on the standard *ResNet18* pretrained Model by modifying the top fully connected layer to the required number of classes, which gave a test accuracy of 66%. For stopping criterion of training, validation loss was observed and when validation loss didn't decrease for consecutively, *patience* (a variable) number of epochs, It halts training and returns the last best saved model. For coping with the class imbalance problem, I used a weighted cross entropy loss function. The weights of a class *c* are given as the ratio of the number of examples in the majority class and the number of examples in class *c*. This penalizes the misclassification of minority classes heavily thus reducing the effect of class imbalance. In case of this dataset, other methods of removing class imbalances like undersampling or oversampling won't work since undersampling could have led to huge loss in training data, since the difference between the minority classes and majority classes is huge, and undersampling might have led to overfitting to the minority class examples.

### **Problems Faced:**

The main problems faced while doing this project was the noisy dataset, more than half of the time was used for cleaning up the dataset. Also since, there's an unavailability of training resources, with no gpu the training time was dead slow, it was taking around 20 seconds for each batch, i.e.  $20 \times 47 = 15$  mins (approx.) for each epoch. Due to which, I could only test the approach for 15 epochs, but it seems like the model performance could be improved further by training for more number of epochs and using source-3 data as well, and training few more layers of *Resnet18* from the top to better fit the data. Also the whole dataset was large enough to fit in RAM all in once, thus making it necessary for me to read it each time from the hard-disk by using its name given in the labels csv file, and since read operations from the hard disk is slow, this might also be slowing up the training process even further.

*Inference Code and trained models are hosted on [github](#), for community use.*

### **Results:**

- *LeNet* : Testing accuracy: 64.9%
- *ResNet18* : Testing accuracy: 66%

### **Future Prospects:**

- Training this model further for more number of epochs and with more top layers of *ResNet18*.
- Training with Source-3 data first on the given symptoms recognition task then using that as a helper model to train the whole dataset with all the images.
- Trying out more sophisticated and more varieties of architecture and observing the accuracies for them.
- Trying to take help from the region of interest detection model from the [kaggle challenge](#) and only focussing on those regions to drive the accuracy further.
- Implementing a five fold cross-validation based training.

### **References:**

- ❖ He, Kaiming, et al. "Deep residual learning for image recognition." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- ❖ <https://twitter.com/ChestImaging/status/1243928581983670272>
- ❖ <https://www.sirm.org/category/senza-categoria/covid-19/>
- ❖ <https://www.kaggle.com/c/rsna-pneumonia-detection-challenge>
- ❖ <https://github.com/agchung/Figure1-COVID-chestxray-dataset>

- ❖ <https://github.com/ieee8023/covid-chestxray-dataset>
- ❖ <https://pubs.rsna.org/doi/full/10.1148/ryct.2020200028>
- ❖ <https://pytorch.org/>