



CAPSTONE PROJECT

Project Title: X-Ray Pneumonia Prediction

PGA-02

Abstract

The goal of the project is to predict whether a patient is diagnosed with Pneumonia disease or not. The data that is used consists of X-Ray images that are divided into 2 folders train and test. Several deep learning models are trained using train dataset & applied to test dataset in order to evaluate the model performance. These performance measures are then compared to determine which model is best in prediction of Pneumonia disease in patients.

Submitted by - Begum Zubeda

Table of Content

1.	Introduction	2
2.	Normalizing Data & Image Augmentation	3
3.	Data Visualization	5
4.	Handling Imbalanced Data	7
5.	Convolutional Neural Network (CNN)	8
6.	Optimization	10
7.	Steps Performed	11

Capstone Project – PGA-02

X-Ray Pneumonia Prediction

1. Introduction

Chest X-ray images (anterior-posterior) were selected from retrospective cohorts of pediatric patients of one to five years old from Guangzhou Women and Children's Medical Center, Guangzhou. All chest X-ray imaging was performed as part of patient's routine clinical care.

For the analysis of chest x-ray images, all chest radiographs were initially screened for quality control by removing all low quality or unreadable scans. The diagnoses for the images were then graded by two expert physicians before being cleared for training the AI system. In order to account for any grading errors, the evaluation set was also checked by a third expert.

The Dataset

The dataset is organized into 2 folders (train, test) and contains subfolders for each image category (Pneumonia/Normal). There are 5,863 X-Ray images (JPEG) and 2 categories (Pneumonia/Normal).

2. Normalizing Data & Image Augmentation

The X-Ray Images are converted into pixels ranging between 0-255 of shape (height, width, color channel (RGB)) so that it can be processed by the model.

```
train_datagen = ImageDataGenerator(rescale=1./255, horizontal_flip=True, zoom_range=0.2, shear_range=0.2,  
                                   fill_mode='nearest') #Divide pixels by 255(where pixel value range  
test_datagen = ImageDataGenerator(rescale=1./255)
```

The process of standardizing features present in the data in a fixed range or same scale is referred to as normalization. It is performed during the data pre-processing to handle highly varying magnitudes or values or units. In this case, Normalization of data is done by scaling image pixels. Here, each pixel is divided by 255 pixel value such that pixel values ranges between 0 and 1, because 0-255 RGB range pixel values would be too high for our models to process (given a typical learning rate).

Image Augmentation is referred as a process of applying various random transformation on images, which can be applied for following reasons:

- Make the most of our few training examples.
- Helps prevent overfitting and helps the model generalize better. Overfitting happens when a model exposed to too few examples learns patterns that do not generalize to new data.

Several Image Augmentation that can be applied are:

- 'horizontal_flip' is for randomly flipping half of the images horizontally.



Figure 2.1

- ‘shear_range’ is for randomly applying shearing transformations.

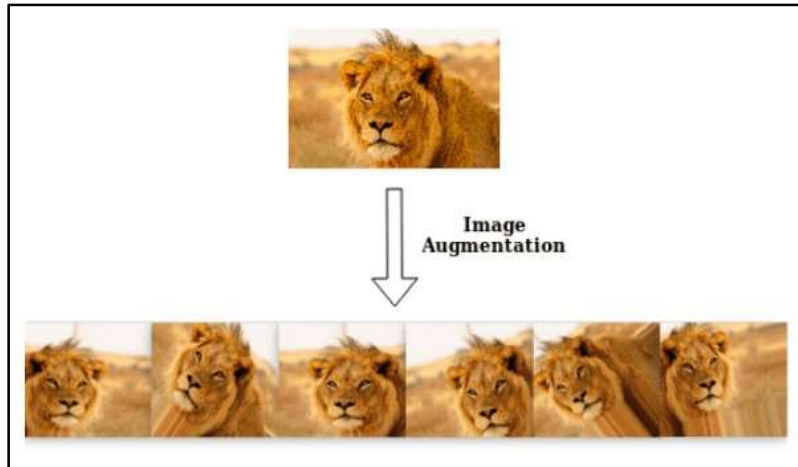


Figure 2.2

- ‘zoom_range’ is for randomly zooming inside pictures.

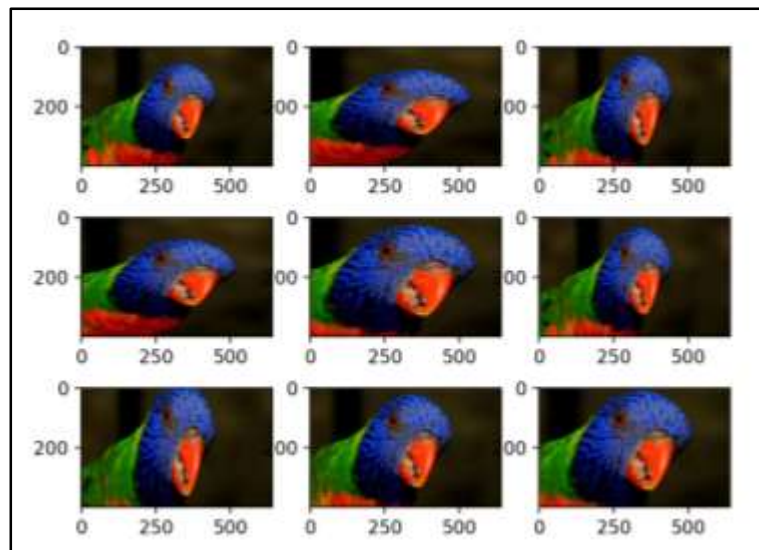


Figure 2.3

- ‘fill_mode’ is the strategy used for filling in newly created pixels, which can appear after a rotation or a width/height shift.

3. Data Visualization

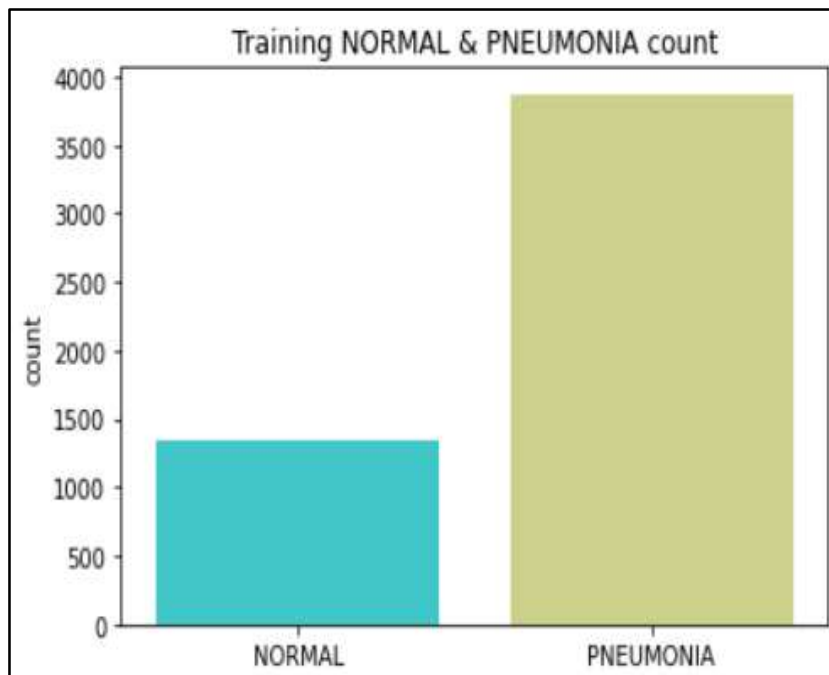


Figure 3.1

The graph on the left represents the total number of Normal and Pneumonia patients in the training dataset. We can clearly see that there are more number of Data/X-Ray Images for Pneumonia patients than Normal patients for training the model, that results in the problem of Imbalanced Data.

The graph on the right represents the total number of Normal and Pneumonia patients in the test dataset to which the trained model is going to be applied. We can clearly see that there are more number of Data/X-Ray Images for Pneumonia patients than Normal patients.

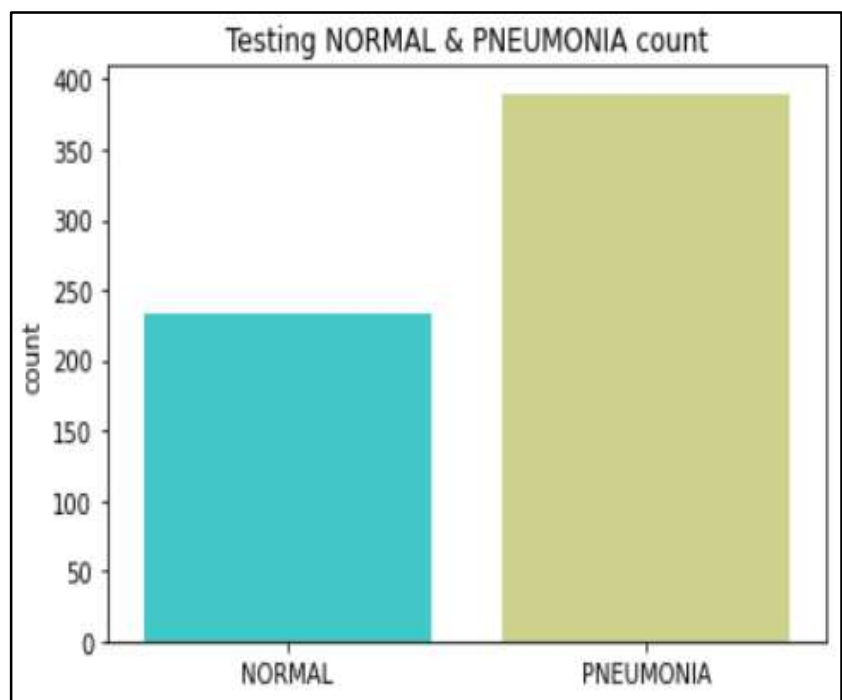


Figure 3.2

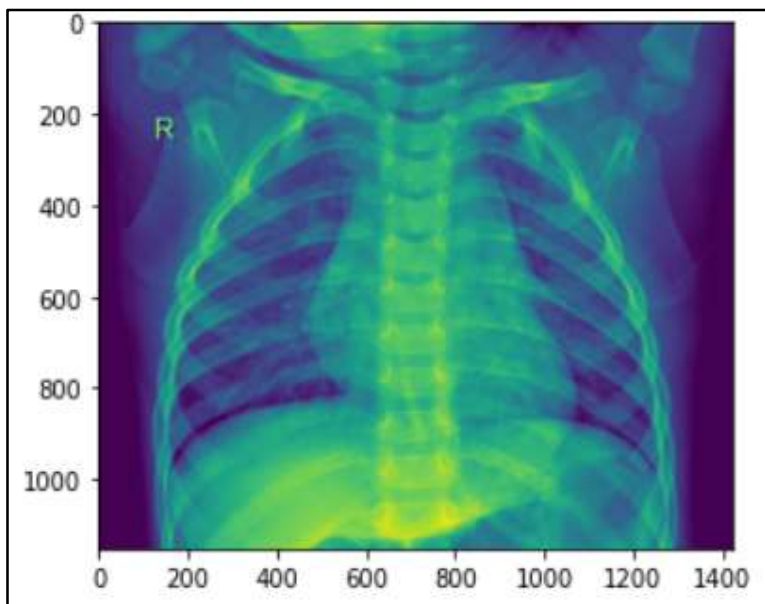


Figure 3.3

The image plot on the left represents X-ray Image for Normal patient. This normal chest X-ray depicts clear lungs without any areas of abnormal opacification in the image.

The image plot on the right represents X-ray Image for Pneumonia patient. We can see that the image manifests with a more diffuse "interstitial" pattern in both lungs.

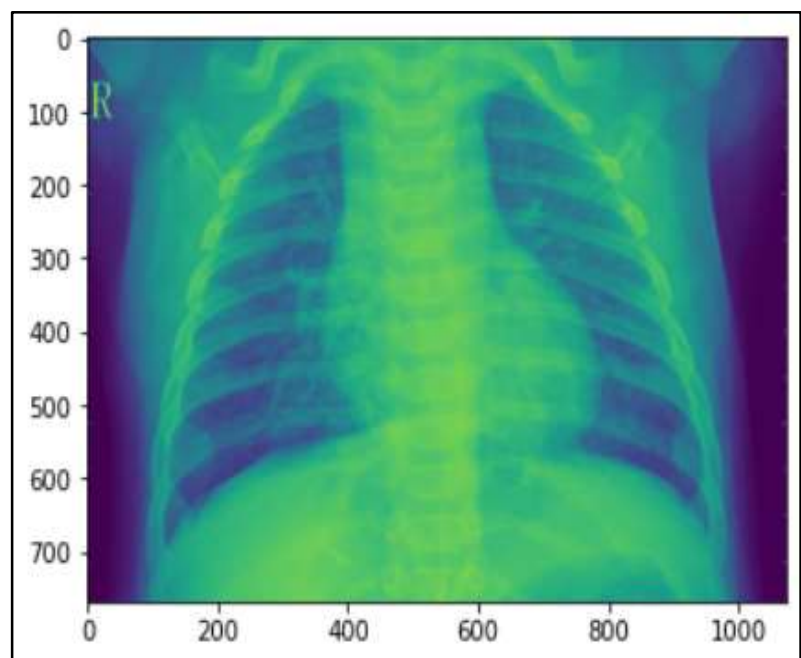


Figure 3.4

4. Handling Imbalanced Data

In a classification problem, a good model is the one that can predict or classify data well.

Imbalanced data is referred to a scenario where the dataset has unequal number of observations or data for each class or in other words difference in the number of observations or data for classes is very huge like we have seen in figure 3.1. This problem can lead to biasness for majority class and when the model is trained with such data it cannot generalize the new data well where bias refers to difference between average prediction of a model & the correct value which it is trying to predict. Therefore, when the optimization is applied while training it basically tries to reduce the misclassification error & gives more priority to the majority class data so that the total error is minimum & the accuracy is maximum, but for minority class data the performance is not so good.

In order to solve this problem we can apply various techniques like undersampling the majority class data, oversampling the minority class data and the technique we used is setting the class weights where we assign higher weights for minority class using the imbalanced ratio while training our model with train dataset, and this is the most efficient way where we can keep the data as it is without losing or adding extra data.

```
Normal:1341
Pneumonia:3875
Imbalance Ratio: 2.89

Using class_weights as:
Normal:2.89
Pneumonia:1.00
```

Figure 4.1

5. Convolutional Neural Network (CNN)

Neural Network is basically a type of machine learning algorithm that is inspired by human biological brain. The goal is to solve the problem in a similar way as the human brain does. A Convolutional Neural Network (CNN) is a type of neural network that specializes in Image Recognition & Computer Vision tasks.

CNN has several types of layers:

- **Input Layer:** This layer consists of Image in the form of pixels of shape (height, width, color channels (RGB)).

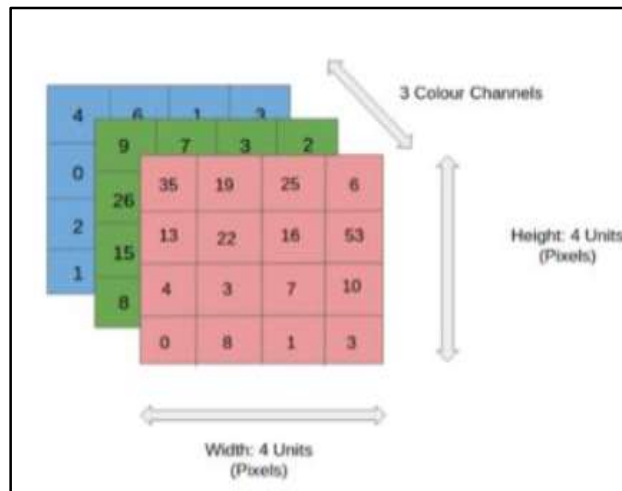


Figure 5.1

- **Convolutional Layer:** A filter/kernel that is in the form of matrix scans few pixels at a time creates a feature map that predicts the class that each feature belongs to.

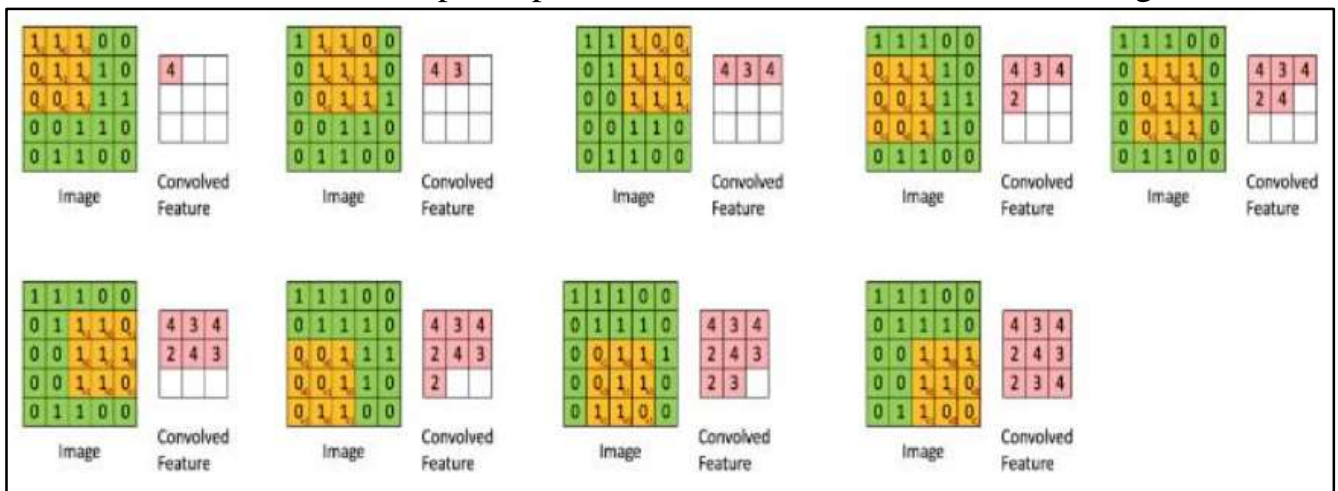


Figure 5.2

- **Pooling Layer:** This layer reduces the amount of information in each feature that was obtained in the convolutional layer while maintaining the most important information. In this project we used 2 types of pooling, Max pooling that calculates

the maximum value for each patch of the feature & Global Max pooling that calculates the overall maximum value for feature map.

12	20	30	0
8	12	2	0
34	70	37	4
11	10	25	12
2	0		

Max pooling			
20	30		
112	37		

- **Batch Normalization Layer:** This layer helps to coordinate the update of multiple layers in the model. It scales the output of the layer, specifically by standardizing the activations of each input variable per mini-batch, such as the activations a node from the previous layer which means that assumptions the subsequent layer makes about the distribution of inputs during the weight update will not change, at least not dramatically. This has the effect of stabilizing and speeding-up the training process of deep neural networks.

Activations or activations functions are mathematical gates between input feeding the current neuron and output going to the next layer. Some of the activations are ReLU or rectified linear activation function is a linear function that will output the input directly if it is positive, otherwise, it will output zero and in this project it is used because pixels are always positive, Sigmoid activation function which is the logistic function used for binary classification problems.

- **Dropout Layer:** This layer is used to simulate having a large number of different network architectures by randomly dropping out nodes during training. This is remarkably effective regularization method to reduce overfitting and improve generalization error in deep neural networks of all kinds.
- **Flatten Layer:** This layer converts the multi-dimensional output to 1D vector of pixels.
- **Dense Layer:** A dense layer is a layer that is deeply connected with its preceding layer which means the neurons of the layer are connected to every neuron of its preceding layer.

6. Optimization

Optimization is a technique that finds the values of parameter/weights such that it minimizes or maximizes the objective function of interest. In this project, the objective function is an error/cost function that is needed to be minimized.

Several optimization techniques that are used are:

➤ **Stochastic Gradient Descent:**

In this technique, weights are initialized with random values, a single sample is selected in random from the dataset, then for each sample the algorithm selects optimized weights such that it reduces error/cost function, at last weights are updated & the model is evaluated.

➤ **Adam (Adaptive Moment Estimation):**

This optimization technique basically uses momentum which pushes out our algorithm or model from locally optimal solution or local minima along with RMSProp which is an adaptive learning algorithm that tries to improve AdaGrad which is a gradient-based optimization technique that adapts the learning rate to the parameters(weights), performing smaller updates (i.e. low learning rates) for parameters(weights) associated with frequently occurring features, and larger updates (i.e. high learning rates) for parameters associated with infrequent features. Here, learning rate refers to how quickly an algorithm updates its parameters(weights) or size of the step taken by the algorithm.

7. Steps Performed