# LUNG DISEASE DETECTION Using VGG16 State-of-the-Art Algorithm

## Description

This project as the title suggests focuses on the detection of any kind of lung disease or no disease when certain image(s) is/are given as inputs. So, this is a very basic project which actually shows how deep learning especially convolution neural networks can be used in the medical and healthcare industry.

Here we have used the State-of-the-Art convolutional neural network model VGG16 using transfer learning.

## Advantages of using Vgg16

- This model is small in size (only 16 layers) but efficient in terms of accuracy for predictions with respect to its size.
- Its biggest advantage is that the layers present in the model have consistent parameters (i.e., for each layer of the same type we have same value of the parameters).
  So, it focuses more on the layers rather than on the other hyper parameters of these layers like filter-size, number of units, strides, etc.
- This helps is improving the consistency and efficiency of the model for predictions.
- E.g., is that for convolutional layers the stride of 1 or filter size of 3X3 is fixed for all conv. Layers.
- Present in KERAS applications and so required no separate downloading.
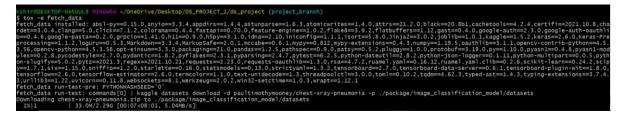
More on the layered architecture of VGG16 in the low-level document.

## Using the repository (Use the commands given in the Readme file)

### 1] For fetching of the dataset from Kaggle:

Active Kaggle account is required. The dataset used in this repository is Chest X-Ray Images (Pneumonia) available on the following link: https://www.kaggle.com/paultimothymooney/chest-xray-pneumonia



Fetching the data from Kaggle

After getting the zip file unzip it and get ready to use this dataset. The images in the dataset will look something like this:

The normal chest X-ray (left panel), Bacterial pneumonia (middle), whereas viral pneumonia (right) manifests in both lungs.

So, the live data for uploading needs to be similar to the form of above images.

## 2] For training the MODEL:

First of the images that we have fetched from Kaggle we need to tabulate them in the form of a table consisting of image path and the following class it belongs to namely (No Disease i.e., Normal, Pneumonia Bacteria, and Pneumonia Virus)



In the research environment we get this is the type of table. Similarly, we will get in our package training. Training the model is essentially training it on this table after converting the image names into vectors.

Create a pipeline with a customized transformer to handle big data and convert it into vector array and the VGG16 model. Then fit the pipeline on this training data and save it for further use. Set the epochs and batch size base upon computational strength of the system.

Model summary



Training of model with testing it on test images.

Following are the results:                    Epochs=2        Batch Size= 64

On the first epoch I got a training accuracy of 0.8495 and validation accuracy of 0.1509

On the second epoch I got training accuracy of 0.9025 and validation accuracy of 0.2808

This shows that if the epochs and batch size are increased, we can achieve better results.

# 3] For running the API:

First, we build the API using fastapi module leveraging the python async.io web framework. Then we run this API with the uvicorn ASGI web server.

Also, we perform logging using loguru module.

https://drive.google.com/file/d/1-vfWtJmDneP24zPrlfkjW1U5ARf9lQ7D/view?usp=sharing

Refer to the above video.

**FOR RUNNING THE REPOSITORY REFER TO THE COMMANDS IN README FILE.**