# Constituting an End-to-End Deep Neural Network to precisely distinguish distinct types of Pneumonia using Chest X-Ray Images.

Swastik Nath.
Udacity Machine Learning Engineer Final Capstone Project Report.

PROJECT OVERVIEW:

Pneumonia is still one of the most infected lungs disease whose modus operando is to infect the alveoli of the lungs and exhibit manifestations like breathing difficulties, fever, chest pain and if remain without proper diagnosis for a long time, it might end up being fatal. Every year, pneumonia concerns approximately 450 million people globally (7% of the population) and results in about 4 million deaths. With the influx of antibiotics and vaccines in the 20th century, survival has been considerably improved. Notwithstanding, in developing countries, and also among the very old, the very young and the chronically ill, pneumonia prevails as a leading cause of mortality. Pneumonia often shrinks the period of agony amidst those already approaching to death and has consequently been denominated "the old man's friend". Pneumonia generally commences as an upper respiratory tract infection that leads into the lower respiratory tract. It is a variety of pneumonitis (lung inflammation). The normal verdure of the upper airway gives shield by colliding with pathogens for nutrients. In the deeper airways, reflexes of the glottis, activities of correlative proteins and immunoglobulins are essential for safekeeping. Micro aspiration of contaminated excretions can taint the coarser airways and produce pneumonia.

The detection of Pneumatic inferection in the lungs can be determined from Chest Xray images by expert doctors. But Karmany, Goldbaum et. al in their research paper published in the Cell suggested that we can infuse the power of Artificial Intelligence with it to create a clinical decision support system to properly diagnose the disease just by looking at the X-Ray Images like the expert doctors do. According to the investigation done by Daniel S. Karmany, Michael Goldbaum et. al1 they introduced the architecture of CNN model already pre-trained on the ImageNet Dataset which carries almost 10 million images, where the model's feature extractor's parameters. The latest anxieties regarding viral epidemics are one of the most demanding research sectors and imbuing AI with that research would be great. In this clinical decision support diagnostic tool based on deep learning frameworks, we need to use a method called Transfer Learning which can train a neural network with just fractions of the data requirement.

The dataset Version 2.0 and Version 3.0 for this model is made accessible by Daniel Karmany, Kang Zhang, Micheal Goldbaum et.al1 from Mendeley Data by Elsevier2, the dataset is associated to the University of California, San Diego and Guangzhou Women and Children's Medical Centre. The data contains thousands of expert validated Chest X-Ray images and are split into Train and Test set of independent images, and each of these sets there is two

more distributions named Normal and Pneumonia. In the distribution Normal, there are thousands of Pneumonia Negative X-ray images which is named with the patient's id and type of the image class such as normal. In the Pneumonia distribution, we get other two types of images such as bacterial or viral contamination with their names as their representing labels appended with patient's id. The dataset is completely anonymized, so there is no possible way to trace back to the patient from the images. To feed these data to a machine learning model we need to redistribute the images into properly configured directories.

The whole dataset both Version 2.0 and Version 3.0 is available in the Mendeley Data website as zip files where version 2.0 is around 1.15 GB in size while the Version 3.0 is around 7.9 GB in size.

## PROBLEM STATEMENT:

To be particular in expressions of machine learning vocabularies, this predicament of classifying medical images based on specialist identified X-Ray images, is a type of Supervised Image-based Multiclass Classification problem, which is simply, every single image instance is a type of a single class. In this intricacy in hand, there is an entirety of 3 classes the images can be arranged to, the first one is Bacterial Pneumonia, the second one is Viral Pneumonia and the third one is the Pneumonia Negative. The photographs are X-Ray images upon which we need to train our model and render judgment. The X-Ray images being BW the images might have to be converted to RGB, notably when inferring as the ImageNet model is trained upon 3 channel RGB images. The images also need to be resized.

In contemporary advancements encompassing deep learning architectures, Convolutional Neural Networks (CNN) are established to be considered sound in terms of extracting and interpreting features of images uniquely. As the research paper by Daniel S. Karmany, Michael Goldbaum, et. al1 recommended that using Deep Convolutional Network which couples Convolutional Layers, RELU, Dropout Layers in different sequences and trained upon the dataset Imagenet comprising about 10 million images of different objects and freezes the feature extractor layers while modifying the classification dataset according to the number of associating classes. We then use a Cross-Entropy loss function to measure the deviation of the predictions of each class and use Stochastic Gradient Descent which behind the screen uses the Softmax Activation function, which transforms meanings of the output layers as the probabilities of the individual classes, which makes much more sense in the real world.

## METRICS:

To accurately evaluate how our models are realising we used the metrics namely training set batch loss, training set accuracy, test set batch loss and test set accuracy.

In our architectural decision scenario, the benchmark VGG19 model with the features extractor layer's gradient calculation set to false and we calculate the average test batch loss and the model's accuracy in the test set. In the case of SageMaker Deployment we did not set the gradient calculations to false and got better performance in terms of model's test set accuracy.

In our Resnet50 design decision scenario, we got to decide the number of epochs based on the test set loss and test set accuracy as the complexity of the model and the unvarying dataset makes the model prone to overfit to the train dataset. So, we decided on training for less number of epochs and got substantially worse train accuracy but even better test accuracy, which certainly means one thing, the model generalizes well across the dataset.

In case of SageMaker deployment we use both the Version 2.0 and the Version 3.0 dataset and that's the reason why we get even better test set accuracy with both the VGG19 model but without freeze parameters and the Resnet 50 model.

DATA EXPLORATION:

So, in our AWS SageMaker Notebook instance, we download the file by copying the https link of the ZIP file of the Version 2.0 and Version 3.0 and issue a linux wget command in the notebook instance which goes ahead and download the file into the AWS SageMaker Notebook Instance workspace. The dataset comes as a ZIP archive file, we extract the compressed files into a directory and see that the dataset contains thousands of BW .jpeg images with their filename acting as their respective class labels. Each of the images are collected from real patients and are labelled by experts of the field. The dataset is devoid of anomalies and corrupt files as the dataset is maintained by several reputed research institutes.

Next, we create a new directory in the notebook instance to hold the properly distributed directories. So, at first, we create two more directories inside the new directory to hold the train and test set. Now, at this stage, we use python libraries and control loops and distribute the images from Normal and Pneumonia directory into 3 separate directories which will work as the label of these images for both the training and test set. So, basically what we did was we changed the distribution of the images and put the images in their specific label directory like we put Bacteria infected Pneumonia X-Ray images into their directory.
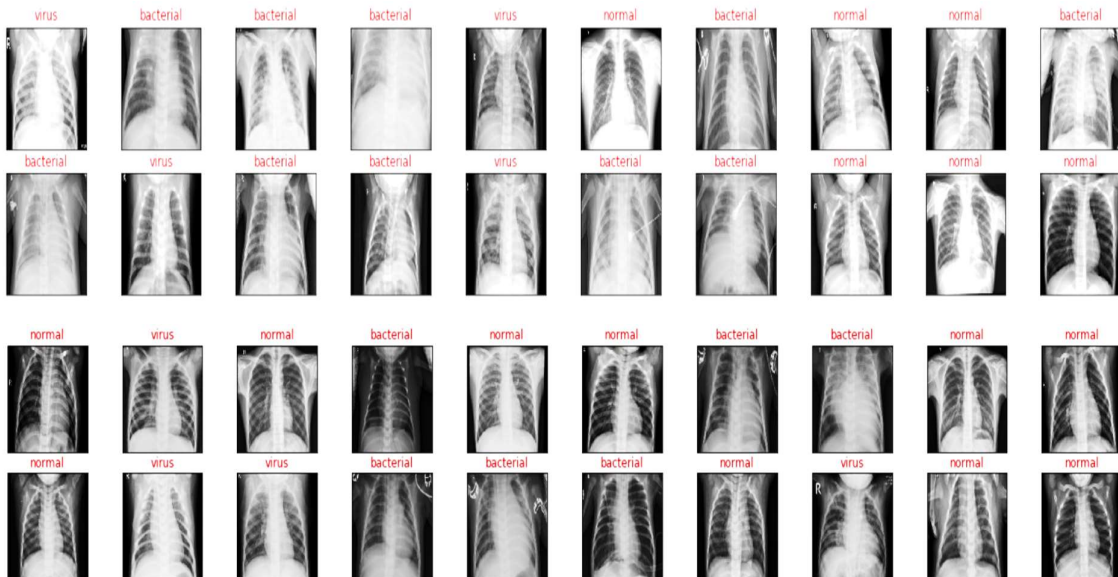
As we have chosen PyTorch as our preferred deep learning framework, the framework requires that the images are distributed into the same number of directories as there are target classes and treats the name of the directory as the label of the image. Next up, we shuffle up the dataset images and divide them into batches to feed our PyTorch based Machine Learning Model. In the process, we also resize them to properly fit the model's requirements and convert them into multidimensional tensors. We use

torchvision.transforms package to do so. PyTorch models perform the calculation on multidimensional tensors.

We do not get into Image Augmentation process here as the dataset we have doesn't vary much and the complex architectures of Deep CNNs, augmentation result into overfitting to the training data. We repeat the process for the test dataset too, and also devoid of any kind of image augmentation process. As we intend to deploy or model via the AWS SageMaker platform, to use the training and the test dataset with a SageMaker PyTorch Estimator, we need to upload the dataset already processed and properly distributed to an Amazon Simple Storage Service (S3) object in the same datacentre region as the model would be deployed to.  We use the .upload_data method of the SageMaker Session object in order to get hold of the S3 Bucket and upload all the .jpeg image files with redistributed directory structure.

EXPLORATORY VISUALIZATIONS:

The dataset actually consists of .jpeg images containing the BW Chest X-Ray images for normal patients, and also X-Ray images of the people infected with bacterial or viral pneumonia, as we prepare the data to feed our PyTorch based model, we must first shuffle the data, then batch the data and transform the images into tensors to perform model's internal matrix calculations. We check to see the batches of the images and observe that each of the batch has got proper randomness and have got each of the classes at least once and that is where lies the beauty of shuffling. We used the PyTorch's DataLoader function to shuffle, batch and prepare the dataset to be fed to the model. In the visualization below we check to see if all the images are properly loaded and prepared for feeding the model.



We check to see two batches and verify that the data is properly labelled and shuffled across the train and test dataset.

ALGORITHMS AND TECHNIQUES:

As the choice of deep neural network architecture, we use the VGG19 Batch Normalization and Resnet 50, both with pretrained weights and biases from the ImageNet classification challenge training. We change the underlying architecture of the model to some small extent, basically in the case of VGG19 Batch Normalization, I added a RELU activation layer and a Dropout Layer with probability of dropout being set to 50% of the total neurons and lastly I added a Linear Output Layer with the out_features set to the number of target classes in this case that is 3.

In the case of Resnet 50, I just made changes to the last layer and changed the last output linear layer's out_features to the number of target classes in this case which is again 3 and I at the first initialized the model's parameters with the parameters learnt from the training on the Imagenet Dataset.

As the loss function in both of the cases of VGG19 and Resnet 50 we have used the Cross Entropy, because it is one of the most efficient function there is to properly calculate the deviation of the model's prediction from it's label especially in cases like multi class or multi label classification problems. We also use Stochastic Gradient Descent as the optimizer here which goes ahead and updates the values of the parameters of each layer based upon the loss of the current epoch and current training datapoint. Another cause of using the SGD as the optimizer is because it uses the SoftMax function behind the screen which generates probabilities for the target class labels and from those probabilities the highest one is the class that's the model predicted. The probabilities are the measure how confident the model is in classifying that specific datapoint.

In case of training and testing of the model we feed the model shuffled batches of X-Ray images

To make inference over the images it first needs to be converted to Python Image Library (PIL) Image object and then it needs to be resized and after which it is transformed into a multidimensional tensor, which participate in the calculation at the CNN architecture. In case of inference, as the images are in BW format, we stack up three BW layers and creates a 3-channel image from the single channel BW Image and feed that image to the model endpoint to get the inference result.

We desired to create an end-to-end model so, we trained and deployed the model using the AWS SageMaker. At the very first of all we combined the Version 2.0 and Version 3.0 dataset together, transformed each image into a PIL Image, resized each image so that the internal matrix multiplications never run into error, and then transformed them into PyTorch Tensors.

We then used the PyTorch's ImageFolder function to load, label and transform the image data, we then used the DataLoader function of PyTorch to load the images into batches

after shuffling them, and next we upload the data to an Amazon S3 Bucket within the region of the notebook instance. We then go on to create SageMaker PyTorch Estimators and pass in the training script to follow the PyTorch Model training instructions, basically it contains all the details about the training procedures of the model after which we direct the S3 bucket location of both the training and test dataset as we uploaded them to the S3 bucket before and call the .fit method of the PyTorch Estimatior object. After the training is over, we compile a SageMaker PyTorchModel object to send inference request to it in real time and we pass in an inference script which includes the instructions to handle the input stream of bytearray as images and process them and then invokes the deployment endpoint. We have created the inference script in such a way that it provides the inferred class and the confidence of the model in the prediction.

The Lambda function can be invoked when an object is uploaded to the S3 and that object is then sent to the lambda function and it performs the pre-processing over the byte array object and then performs different transformations and send the result to a CloudWatch Log.

## SETTING THE BENCH MARK:

We have used only the Version 2.0 Dataset to participate in the architectural design decision procedure before we head over to training with the Amazon SageMaker and deploy the model to save resources and space.

According to the paper by **Karmany, Goldbaum et. al**[1], I downloaded the VGG19 Batch Normalization model with pretrained model parameters on the ImageNet Dataset, and made the gradient calculation for the features (Feature Extractor) layer to false leading to freezing of the layer parameters during the training process and trained it for up to 10 epochs. The accuracy of the model of the over the test dataset was quite impressive up to 78% with 490 correct predictions out of 624 total test datapoints.

## DATA PREPROCESSING:

in our AWS SageMaker Notebook instance, we download the file by copying the https link and issue a wget command which goes ahead and download the file into the AWS SageMaker Notebook Instance workspace. Next, we create a new directory in the notebook instance to hold the properly distributed directories. So, at first, we create two more directories inside the new directory to hold the train and test set. Now, at this stage, we use python libraries and control loops and distribute the images from Normal and Pneumonia directory into 3 separate directories which will work as the label of these images for both the training and test set. So, basically what we did was we changed the distribution of the images and put the images in their specific label directory like we put Bacteria infected Pneumonia X-Ray images into their directory. As we have chosen PyTorch as our preferred deep learning framework, the framework requires that the images are distributed into the same number of directories as there are target classes and treats the name of the directory as the label of the image. Next up, we shuffle up the dataset images and divide them into

batches to feed our PyTorch based Machine Learning Model. In the process, we also resize them to properly fit the model's requirements and convert them into multidimensional tensors. We use torchvision package to do so. PyTorch models perform the calculation on multidimensional tensors. We do not get into Image Augmentation process here as the dataset we have doesn't vary much and the complex architectures of Deep CNNs, augmentation result into overfitting to the training data. We repeat the process for the test dataset too, and also devoid of any kind of image augmentation process. As we intend to deploy or model via the AWS SageMaker platform, to use the training and the test dataset with a SageMaker PyTorch Estimator, we need to upload the dataset already processed and properly distributed to an Amazon Simple Storage Service (S3) object in the same datacentre region as the model would be deployed to.

IMPLEMENTATIONS:

To resolve this problem what we need to do is to use a Deep Convolutional Neural Network, with each of layer's parameters set to the values as it learned throughout its training on the ImageNet Dataset. The research paper by **Karmany et. al** suggested that we can freeze the feature extractor layers of the neural network architecture and update only the weights and biases of the last few layers, but in my case, I found that, doing so makes the model quite inaccurate in both training and test set. So, to remediate that, what I have done is initialized the model with the learned weights and biases for each of the layer but made the model update its full parameter range following the current gradient. It made the model to converge faster as it was able to get around 82% test set accuracy with only 4 epochs of training.

As the choice of deep neural network architecture, we use the VGG19 Batch Normalization and Resnet 50, both with pretrained weights and biases from the ImageNet classification challenge training. We change the underlying architecture of the model to some small extent, basically in the case of VGG19 Batch Normalization, I added a RELU activation layer and a Dropout Layer with probability of dropout being set to 50% of the total neurons and lastly I added a Linear Output Layer with the out_features set to the number of target classes in this case that is 3.

In the case of Resnet 50, I just made changes to the last layer and changed the last output linear layer's out_features to the number of target classes in this case which is again 3 and I at the first initialized the model's parameters with the parameters learnt from the training on the Imagenet Dataset.

As the loss function in both of the cases of VGG19 and Resnet 50 we have used the Cross Entropy, because it is one of the most efficient function there is to properly calculate the deviation of the model's prediction from it's label especially in cases like multi class or multi label classification problems. We also use Stochastic Gradient Descent as the optimizer here which goes ahead and updates the values of the parameters of each layer based upon the loss of the current epoch and current training datapoint. Another cause of using the SGD as

the optimizer is because it uses the SoftMax function behind the screen which generates probabilities for the target class labels and from those probabilities the highest one is the class that's the model predicted. The probabilities are the measure how confident the model is in classifying that specific datapoint.

In case of training and testing of the model we feed the model shuffled batches of X-Ray images

To make inference over the images it first needs to be converted to Python Image Library (PIL) Image object and then it needs to be resized and after which it is transformed into a multidimensional tensor, which participate in the calculation at the CNN architecture. In case of inference, as the images are in BW format, we stack up three BW layers and creates a 3-channel image from the single channel BW Image and feed that image to the model endpoint to get the inference result.

We faced quite a few troubles, while processing the input image over http as bytearray and send it over to the model for inference, so we had to device a new Lambda Function that does this for us.

REFINEMENT:

The benchmark model might not be performing as good as the Resnet model because the model was pretrained on RGB images with so much variation of images, but the X-Ray images doesn't variate quite a much as the Imagenet Dataset varies. So, the model cannot update parameters of the feature extractors based upon the loss values of predictions, and this might drive the model to a slightly worse performance than the Resnet model with full parameter range gradient calculation across the training set.

```python
import torch.nn as nn
classes = ['bacterial', 'normal', 'virus']
in_f = vgg19.classifier[6].out_features
vgg19.classifier.add_module(module = nn.ReLU(inplace=True), name = '7')
vgg19.classifier.add_module(module = nn.Dropout(p=0.5, inplace=False), name='8')
vgg19.classifier.add_module(module= nn.Linear(in_features = in_f, out_features = len(classes), bias=True), name='9')
print(vgg19.classifier)

if train_on_gpu:
    vgg19.cuda()
```

In the image above we are fine tuning the model to adapt to the classification problem we currently have for the VGG19 Batch Normalization model. In the image below we set the very last feature of the model to adapt to our training job in hand.

```python
def __init__(self, output_dim):

    super(Resnet, self).__init__()
    resnet = torchvision.models.resnet50(pretrained=True)

    classes = ['bacterial', 'normal', 'virus']
    in_f = resnet.fc.in_features
    resnet.fc = nn.Linear(in_features = in_f, out_features = len(classes), bias = True)

    self.model = resnet
```

So, to get even better test set accuracy we implement the model by finetuning the classification parameters such as for the case of VGG19 Batch Normalization, we add a few layers to the end of the last linear layer and for the case of the Resnet Model we change up the number of the output of the last layer to be equal to the number of target class labels. We use the full parameter range for gradient calculation at the time of training.

```python
# Load the training data.
train_loader, test_loader = _get_train_data_loader(args.batch_size, args.data_dir, args.valid_dir)

model = Resnet(args.output_dim).to(device)

  Defining an optimizer and loss function for training
criterion = nn.CrossEntropyLoss()

optimizer = optim.SGD(model.parameters(), lr=0.001)

# Trains the model (given line of code, which calls the above training function)
train(model, train_loader, test_loader, args.epochs, criterion, optimizer, device)
```

Using the Full parameter range of the Resnet model

MODEL EVALUATION AND VALIDATION:

The final model is a Resnet50 model which takes input of (1,3,224,224) multidimensional image tensor and results in a list of 3 probabilities of the model's confidence of the datapoint being one of the 3 target class.

While training the model there are 3 hyperparameters to consider, one is the train epochs, batch size of the train and the test dataset and the third one is the Learning Rate of the Stochastic Gradient Descent Algorithm. The default value of momentum of the SGD optimizer is well enough for this training task.

With trial and error, we set the value of the above hyperparameters for the Resnet50 model to the following:

| NAME OF THE HYPERPARAMETER | VALUE OF THE HYPERPARAMETER |
|---|---|
| TRAIN EPOCHS: | 3 |
| BATCH SIZE | 64 |
| LEARNING RATE FOR SGD | 0.01 |

To check for the robustness of the model we go ahead and perform inference over individual images in a test dataset using the endpoint of the SageMaker deployed model and we get the same result as we got from the model while validating the model with the test dataset. So, we get an idea that the model's results are accurate and is not getting influenced by any other factor.

We perform inference over negative out of context images and we can see that although the model has predicted a class but it does so with a very low confidence score. So, by this means we can conclude that the weights and parameters have been updated such that the model specializes in predicting only the Chest X-ray images.

JUSTIFICATION:

In this section we compare the final models with the benchmark architecture along with the evaluation metrics associated with each of the models.

| ARCHITECTURE | TEST LOSS | TEST ACCURACY | CORRECTLY PREDICTED LABELS | TOTAL TEST LABELS |
|---|---|---|---|---|
| VGG 19 BATCH NORMALIZATION WITH GRADIENT CALCULATION FOR THE CLASSIFIER LAYER SET TO FALSE (BENCHMARK MODEL) | Average test loss : 0.22 | 78 % | 490 | 624 |
| VGG19 BATCH NORMALIZATION WITH GRADIENT CALCULATION SET TO TRUE TRAINED WITH DATASET VERSION 2.0 AND 3.0(Deployed to SageMaker) | Average Loss in Test set: 0.223 | 79.19% | 989 | 1248 |
| RESNET 50 MODEL WITH GRADIENT CALCULATION SET TO TRUE TRAINED WITH DATASET VERSION 2.0 | Average Loss in Test set: 0.21 | 79% | 498 | 624 |
| RESNET 50 MODEL WITH GRADIENT CALCULATION SET TO TRUE TRAINED WITH DATASET VERSION 2.0 AND 3.0 (Deployed to SageMaker) | Average Loss in Test set: 0.0095 | 82.08% | 1024 | 1248 |

The model when deployed will act as a clinical decision-making algorithm while in use, so apart from the consultant expert will be used to strengthen the confidence of the consultant expert, the accuracy of 82.08% accuracy and the very minimum loss in the test set is quite good in terms of the complexity of the challenge. As the research paper suggested the model I tried and discussed here takes quite a different approach and extends the plethora of the model from academic only constraints to production environment and to the grasp of the daily users.

CITATIONS:

1. Kermany, Daniel; Zhang, Kang; Goldbaum, Michael (2018), "Large Dataset of Labeled Optical Coherence Tomogra Kermany, Daniel; Zhang, Kang; Goldbaum, Michael (2018), "Large Dataset of Labeled Optical Coherence Tomography (OCT) and Chest X-Ray Images", Mendeley Data, v3phy (OCT) and Chest X-Ray Images", Mendeley Data, v3

2. Kermany, Daniel; Zhang, Kang; Goldbaum, Michael (2018), "Labeled Optical Coherence Tomography (OCT) and Chest X-Ray Images for Classification", Mendeley Data, v2

3. Kermany, Daniel; Goldbaum, Michael; et. al "Identifying Medical Diagnosis and Treatable Diseases by Image-Based Deep Learning", Cell, Volume 172, Issue 5, P1122-1131.E9, February 22, 2018, CellPress : DOI: https://doi.org/10.1016/j.cell.2018.02.010

4. Yadav, S.S., Jadhav, S.M. Deep convolutional neural network based medical image classification for disease diagnosis. J Big Data 6, 113 (2019). https://doi.org/10.1186/s40537-019-0276-2