# Data Science Career Track

# Capstone 2 -

# Final Report

by Edward Franke
06/11/2019

Detect Pneumonia is chest x-rays Project – Final Report

EXECUTIVE SUMMARY:

The purpose for this project is to find a correlation connected to chest x-rays containing pneumonia that can separate them in real time compared to normal healthy chest x-rays.  The dataset is a cleaned dataset from Kaggle.com.  Tensorflow and Keras has discovered connections and methods to determine which x-rays have pneumonia with varying success.

IDEA:  A model to detect pneumonia is chest x-rays. (problem to solve)
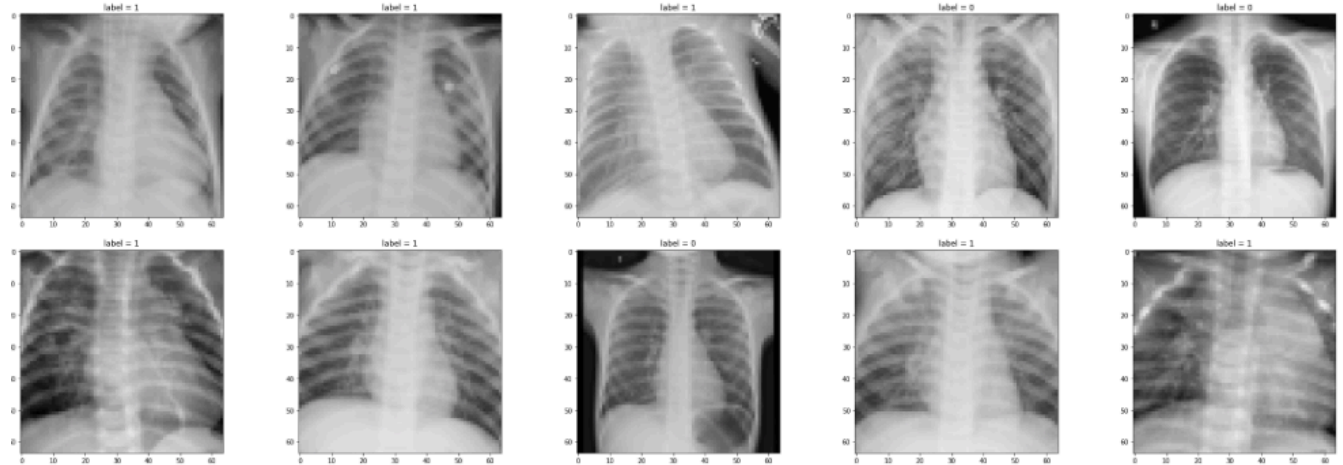
CLIENT:  Medical Industry

REASON:  Pneumonia is a very serious condition that has the potential for death.  I have personal knowledge of how serious it can be.  The sooner it can be detected, the better the chance for survival and less damage to the lungs.

DATA:  From Kaggle, 2 cleaned datasets with 5863 images and over 112,000 images.

SOLUTION:  Create a model or analysis to discover what makes an x-ray image of the chest to have pneumonia and to capture this automatically.

DETAILS:  See below.

DELIVERABLES:  Code and a presentation outlining the discoveries.

## Pneumonia Dataset

```python
1  # Metrics
2
3  # Getting predictions
4  predictions = model.predict(x=x_test)
5
6  acc = accuracy_score(y_test, np.round(predictions))*100
7  tn, fp, fn, tp = confusion_matrix(y_test, np.round(predictions)).ravel()
8
9  print('Accuracy: {}%'.format(acc))
10 print('Precision: {}%'.format(tp/(tp+fp)*100))
11 print('Recall: {}%'.format(tp/(tp+fn)*100))
```

```
Accuracy: 68.910256410256641%
Precision: 82.236842105263315%
Recall: 64.10256410256641%
```

Tensorflow and Keras has discovered some connections.

```
_____
Layer (type)                 Output Shape              Param #
================================================================
input_1 (InputLayer)         (None, 64, 64, 3)         0
_____
conv2d_1 (Conv2D)            (None, 64, 64, 16)        448
_____
conv2d_2 (Conv2D)            (None, 64, 64, 16)        2320
_____
max_pooling2d_1 (MaxPooling2 (None, 32, 32, 16)        0
_____
conv2d_3 (Conv2D)            (None, 32, 32, 32)        4640
_____
conv2d_4 (Conv2D)            (None, 32, 32, 32)        9248
_____
batch_normalization_1 (Batch (None, 32, 32, 32)        128
_____
max_pooling2d_2 (MaxPooling2 (None, 16, 16, 32)        0
_____
conv2d_5 (Conv2D)            (None, 16, 16, 64)        18496
_____
conv2d_6 (Conv2D)            (None, 16, 16, 64)        36928
_____
batch_normalization_2 (Batch (None, 16, 16, 64)        256
_____
max_pooling2d_3 (MaxPooling2 (None, 8, 8, 64)          0
_____
flatten_1 (Flatten)          (None, 4096)              0
_____
dense_1 (Dense)              (None, 256)               1048832
_____
dropout_1 (Dropout)          (None, 256)               0
_____
dense_2 (Dense)              (None, 64)                16448
_____
dropout_2 (Dropout)          (None, 64)                0
_____
dense_3 (Dense)              (None, 1)                 65
================================================================
Total params: 1,137,809
Trainable params: 1,137,617
Non-trainable params: 192
_____

None
```

```
In [69]:    1  # Metrics
            2
            3  # Getting predictions
            4  predictions = model.predict(x=x_test)
            5
            6  acc = accuracy_score(y_test, np.round(predictions))*100
            7  tn, fp, fn, tp = confusion_matrix(y_test, np.round(predictions)).ravel()
            8
            9  print('Accuracy: {}%'.format(acc))
           10  print('Precision: {}%'.format(tp/(tp+fp)*100))
           11  print('Recall: {}%'.format(tp/(tp+fn)*100))
```

Accuracy: 79.48717948717949%
Precision: 76.09561752988047%
Recall: 97.948717948717794%

```
In [70]:    1  print(confusion_matrix(y_test, np.round(predictions)).ravel())
```

[114 120    8 382]

## NIH Dataset

The NIH dataset contains 112,000 x-rays with varying ailments or no ailments.  The code used for the Pneumonia dataset is ineffective on this dataset.

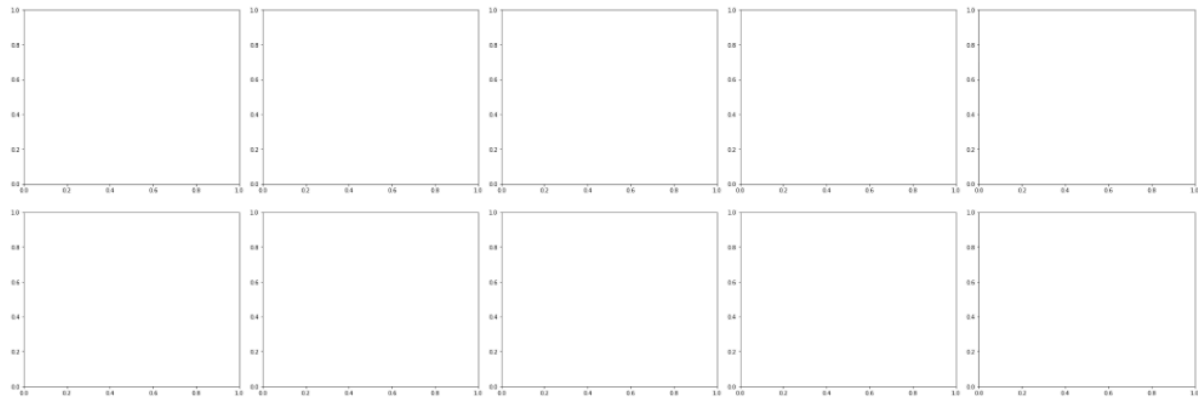Back to orginal code

```
In [ ]:     1  #nb_train_samples = 16188 #8094 #3036 #18046 #111589 #113243 #139987
            2  nb_train_samples = 88
            3  nb_validation_samples= 336
            4  epochs = int(nb_train_samples/batch_size)*3
            5  history = model.fit_generator(
            6      train_generator,
            7      steps_per_epoch=batch_size, #nb_train_samples/batch_size,
            8      epochs=epochs,
            9      validation_data=validation_generator,
           10      validation_steps=batch_size, #nb_validation_samples/batch_size, #val_batch_size,
           11      callbacks=callbacks_list,
           12      verbose=1)
```

```
---------------------------------------------------------------------------
IndexError                                Traceback (most recent call last)
<ipython-input-5-614047ec5bff> in <module>
      5
      6 for i in range(ax.shape[0]):
----> 7     ax[i].imshow(x_test[i], cmap='gray')
      8     ax[i].set_title('label = {}'.format(y_test[i]))

IndexError: index 0 is out of bounds for axis 0 with size 0
```
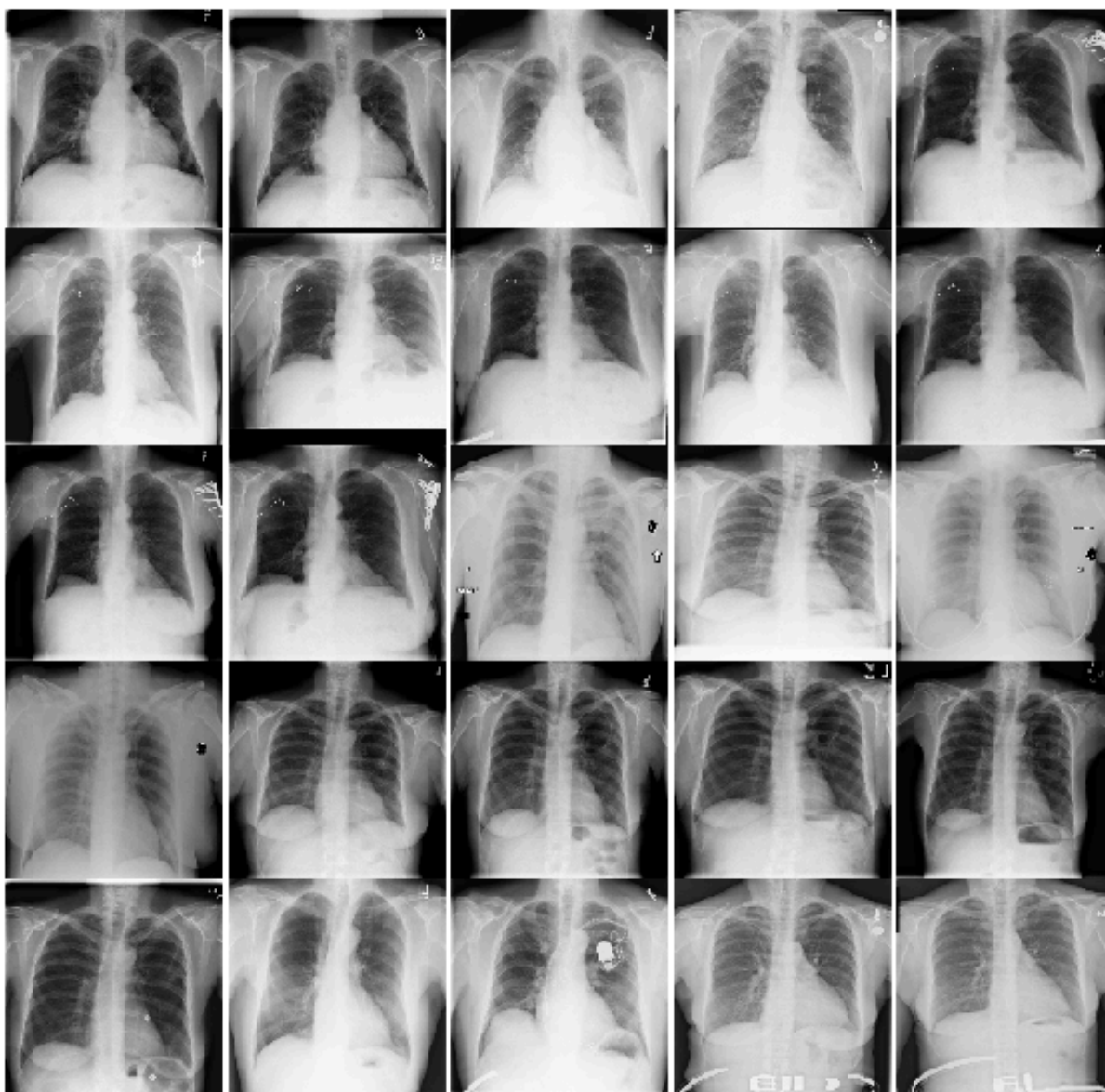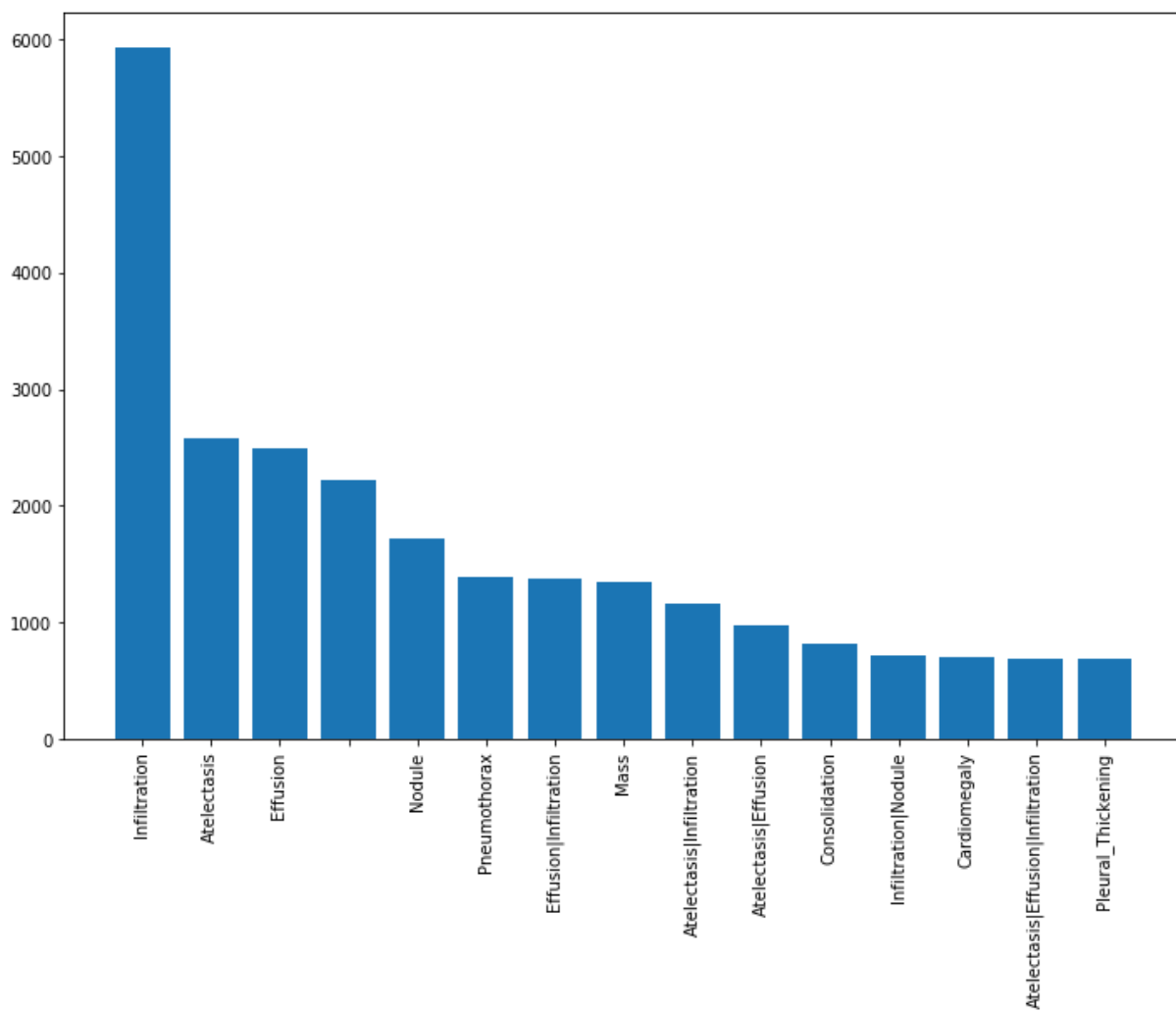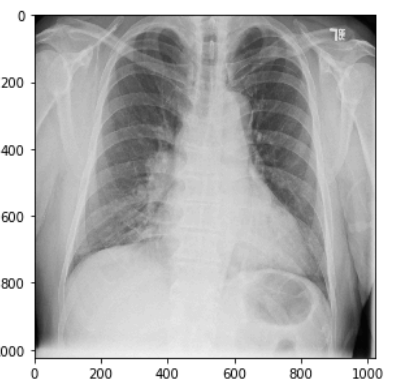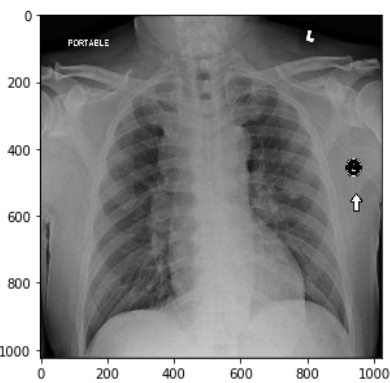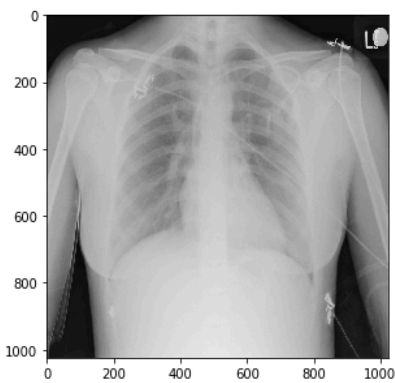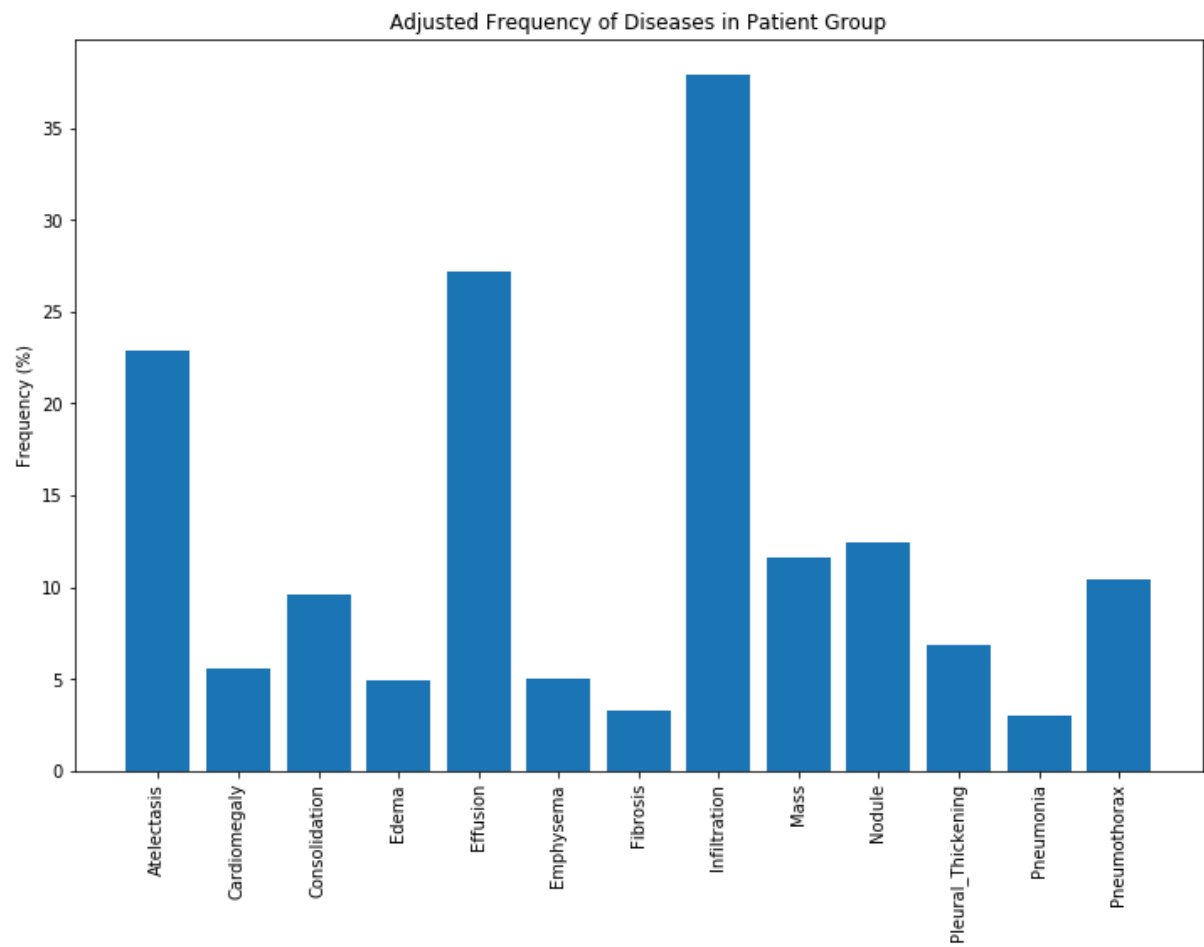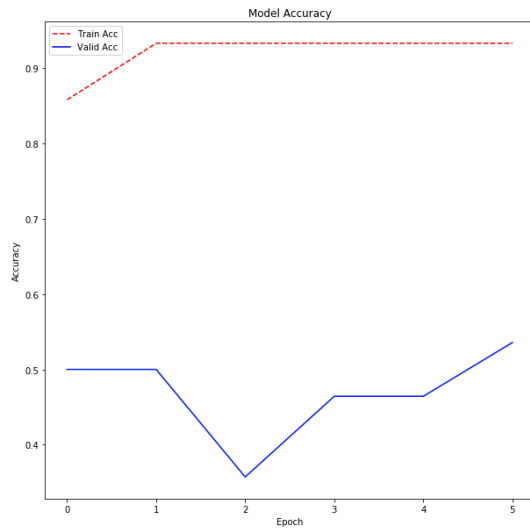


Code written specifically for this dataset gathers labels from a csv file and is also able to analyze the images to determine the ailment being searched for. This code was originally set to discover Fibromyalgia but was modified to detect Pneumonia.

Adjusted Frequency of Diseases in Patient Group

Model Accuracy

Tensorflow and Keras has discovered some connections.  With the number of Epochs increased, the accuracy will also increased but this will also require decent computer processing power due to the number of images being processed.



Confusion matrix, without normalization

TN = 115

FP = 66

FN = 101

TP = 70