

# Configuration Manual

MSc Research Project  
Data Analytics

**Prajwal S Mastamardi**  
Student ID: x19158718

School of Computing  
National College of Ireland

Supervisor: Dr. Barry Haycock

**National College of Ireland**  
**MSc Project Submission Sheet**  
**School of Computing**



**Student Name:** Prajwal S M  
**Student ID:** X19158718  
**Programme:** MSc in Data Analytics **Year:** 2021-2022  
**Module:** MSc Research Project  
**Lecturer:** Dr. Barry Haycock.  
**Submission Due Date:** 16/12/12  
**Project Title:** Anomaly identification in chest radiography  
**Word Count:** 600 **Page Count:** 5

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** .....

**Date:** .....

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
<b>Attach a Moodle submission receipt of the online project submission,</b> to each project (including multiple copies).	<input type="checkbox"/>
<b>You must ensure that you retain a HARD COPY of the project,</b> both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

<b>Office Use Only</b>	
Signature:	
Date:	
Penalty Applied (if applicable):	

# Configuration Manual

Prajwal S M  
Student ID: x19158718

## 1 Introduction

This setup guide will assist you in recreating the research "Lung Cancer Detection Using Classification Algorithms" from the ground up. This setup manual provides thorough information on the prerequisites for setting up, running, and testing this study using the specified framework.

## 2 Hardware Specification

- Operating System: Windows 10 Home Single Language (10.0, Build 18362)
- Processor: Intel(R) Core(TM) i7-8750H CPU @ 2.20GHz 2.21 GHz
- Installed RAM: 8.00 GB
- System Type: 64-bit OS, x64-based processor

## 3 Python Environment

Since the project was entirely written in Python, the Anaconda framework was used. Jupyter lab, Jupyter Notebook, R Studio, VS Code are among other preloaded environments in the Anaconda framework.

The Jupyter Notebook environment is launched as seen in Figure 1. Jupyter provides a user interface for writing code, creating models, and testing them

## 4 Libraries required

```
[1]: 1 import tensorflow as tf
2 from tensorflow import keras
3 from tensorflow.keras.layers import Dense, Activation, Dropout, Conv2D, MaxPooling2D, BatchNormalization
4 from tensorflow.keras.optimizers import Adam
5 from tensorflow.keras.metrics import categorical_crossentropy
6 from tensorflow.keras import regularizers
7 from tensorflow.keras.preprocessing.image import ImageDataGenerator
8 from tensorflow.keras.applications.vgg16 import VGG16
9 from tensorflow.keras.applications.vgg16 import preprocess_input
10 import tensorflow as tf
11 import tensorflow
12 from tensorflow import keras
13 from tensorflow.keras.applications.resnet50 import ResNet50
14 from tensorflow.keras.layers import Flatten, Dense, MaxPool2D, Dropout
15 from tensorflow.keras.models import Model
16 import numpy as np
17 import pandas as pd
18 import time
19 import os
20 import glob
21 from pathlib import Path
22 from tensorflow.keras import layers
23 from sklearn.model_selection import train_test_split
24 import matplotlib.pyplot as plt
25 from matplotlib.pyplot import imshow
26 import os
27 import cv2
28 import shutil
29 from pandas import DataFrame
30 from sklearn.metrics import confusion_matrix, classification_report
31 pd.set_option('display.width', 150)
32 import seaborn as sns
```

## 5 Dataset

There are 3616 COVID-19 positive instances in the database, as well as 10,192 Normal, 6012 Lung Opacity (Non-COVID lung infection), and 1345 Viral Pneumonia photos. As fresh x-ray images for COVID-19 pneumonia patients become available, author continue to update this database.

Database link: <https://www.kaggle.com/tawsifurrahman/covid19-radiography-database>

## 6 Loading the Dataset

**Case 1:** In this scenario the transfer learning algorithm such as VGG-16 and ResNet-50 were utilized to differentiate four classes namely Covid, Pneumonia, Non-covid infection (lung opacity) and Normal healthy lung images.

**Case 2:** In this experiment transfer learning algorithms were applied to classify three classes Covid, Pneumonia and normal. Dataset utilized to train model consist of Covid, pneumonia and normal.

**Case 3:** For this experiment, Anomaly dataset was created to train transfer learning algorithms, Anomaly dataset consist of normal images and anomaly set, anomaly set consist of images from Covid, pneumonia , and non-Covid infection. transfer learning algorithms were applied to differentiate between classes of normal and anomaly(infected), dataset consist of approx. 10,000 images of normal and anomaly images.

**Other cases:** dataset consisting of two classes namely normal and covid images were created for training the model.

Model training:

**Dataset loading**

```
In [17]: 1 train_images = train_generator.flow_from_dataframe(
2         dataframe=df_train,
3         x_col='path',
4         y_col='label',
5         color_mode='rgb',
6         class_mode='categorical',
7         target_size=(299, 299),
8         batch_size=BATCH_SIZE,
9         shuffle=True,
10        seed=seed,
11        subset='training'
12    )
13
```

Found 15239 validated image filenames belonging to 2 classes.

```
In [18]: 1 val_images = train_generator.flow_from_dataframe(
2         dataframe=df_train,
3         x_col='path',
4         y_col='label',
5         color_mode='rgb',
6         class_mode='categorical',
7         target_size=(299, 299),
8         batch_size=BATCH_SIZE,
9         shuffle=True,
10        seed=seed,
11        subset='validation'
12    )
```

Found 3809 validated image filenames belonging to 2 classes.

## VGG-16

```
VGG-16
In [20]: 1 vgg = VGG16(input_shape = [299,299,3], weights='imagenet', include_top=False)

In [21]: 1 # don't train existing weights
2 for layer in vgg.layers:
3     layer.trainable = False

In [22]: 1 # our Layers - you can add more if you want
2 x = Flatten()(vgg.output)
3
4 prediction1 = Dense(2, activation='sigmoid')(x)

In [23]: 1 # create a model object
2 model1 = Model(inputs=vgg.input, outputs=prediction1)

In [27]: 1 # tell the model what cost and optimization method to use
2 model1.compile(
3     loss='binary_crossentropy',
4     optimizer='adam',
5     metrics= "accuracy"
6 )

In [28]: 1 ACCURACY_THRESHOLD = 0.99
2 class myCallback(tf.keras.callbacks.Callback):
3     def on_epoch_end(self, epoch, logs={}):
4         if(logs.get('accuracy') > ACCURACY_THRESHOLD):
5             print("\nReached %2.2f%% accuracy, so stopping training!!" %(ACCURACY_THRESHOLD*100))
6             self.model.stop_training = True
7
8     es = myCallback()

In [29]: 1 history2 = model1.fit(
```

## Plotting

```
: 1 #plotting training values
2 sns.set()
3
4 acc = history2.history['accuracy']
5 val_acc = history2.history['val_accuracy']
6 loss = history2.history['loss']
7 val_loss = history2.history['val_loss']
8 epochs = range(1, len(loss) + 1)
9
10 #accuracy plot
11 plt.plot(epochs, acc, color='green', label='Training Accuracy')
12 plt.plot(epochs, val_acc, color='blue', label='Validation Accuracy')
13 plt.title('Training and Validation Accuracy')
14 plt.ylabel('Accuracy')
15 plt.xlabel('Epoch')
16 plt.legend()
17
18 plt.figure()
19 #Loss plot
20 plt.plot(epochs, loss, color='pink', label='Training Loss')
21 plt.plot(epochs, val_loss, color='red', label='Validation Loss')
22 plt.title('Training and Validation Loss')
23 plt.xlabel('Epoch')
24 plt.ylabel('Loss')
25 plt.legend()
26
27 plt.show()
```

## ResNet-50

```

1 res = ResNet50( input_shape=(299,299,3), include_top=False)

1 patience = 1
2 stop_patience = 5
3 factor = 0.5
4
5 callbacks = [
6     tf.keras.callbacks.EarlyStopping(patience=stop_patience, monitor='val_loss', verbose=1, restore_best_weights=True),
7     tf.keras.callbacks.ReduceLROnPlateau(monitor='val_loss', factor=factor, patience=patience, verbose=1)
8 ]

1 # We won't train all parameters again
2 for layer in res.layers:
3     layer.trainable = False

1 x2 = Flatten()(res.output)
2 out = Dense(units=2, activation='sigmoid', name='predictions_res')(x2)
3
4
5 # Creating our model
6 model2 = Model(inputs=res.input, outputs=out)

1 ACCURACY_THRESHOLD = 1.00
2 class myCallback(tf.keras.callbacks.Callback):
3     def on_epoch_end(self, epoch, logs={}):
4         if logs.get('accuracy') > ACCURACY_THRESHOLD:
5             print("\nReached %2.2f%% accuracy, so stopping training!!" %(ACCURACY_THRESHOLD*100))
6             self.model.stop_training = True
7
8 es = myCallback()

1 model2.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

```

## Common test data evaluation:

### Common test data

```

1 [48]: 1 predictions4 = model2.predict(x = test_data2)

1 [49]: 1 y_pred=np.argmax(predictions4,axis=-1)
2 y_pred
3 count =0
4 wrong=0
5 for i in y_pred:
6     if i ==1:
7         count=count+1
8     else: wrong = wrong+1
9 print(count)
10 print(wrong)
11

570
153

1 [50]: 1 test_data22 = tf.keras.preprocessing.image_dataset_from_directory(
2     "C:/Users/prjwl/Desktop/cov_test",
3
4     color_mode="rgb",
5
6     image_size=(299, 299),
7
8 )

Found 723 files belonging to 1 classes.

```

```

1 [51]: 1 predictions4 = model2.predict(x = test_data22)

```

```

1 [51]: 1 predictions4 = model2.predict(x = test_data22)

1 [52]: 1 y_pred=np.argmax(predictions4,axis=-1)
2 y_pred
3 count =0
4 wrong=0
5 for i in y_pred:
6     if i ==1:
7         count=count+1
8     else: wrong = wrong+1
9 print(count," accurately predicted covid images")
10 print(wrong," wrongly predicted")

570 accurately predicted covid images
153 wrongly predicted

```

```

1 [53]: 1 predictions_vgg = model1.predict(x = test_data22)

1 [54]: 1 y_pred=np.argmax(predictions_vgg,axis=-1)
2 y_pred
3 counts =0
4 wrongs=0
5 for i in y_pred:
6     if i ==1:
7         counts=counts+1
8     else: wrongs = wrongs+1
9 print(counts)
10 print(wrongs)

535
188

```

## References

1. Dataset Source :<https://www.kaggle.com/tawsifurrahman/covid19-radiography-database/code>
2. Code Reference for data loading:  
<https://www.kaggle.com/azharabdulaziz/worodproject95>
3. other net sources for implementation(ex:vizaulizations)
4. Code Reference for Tensorflow Learning: [https://www.tensorflow.org/api\\_docs/python/tf](https://www.tensorflow.org/api_docs/python/tf)
5. Code Reference for Learning Keras : [https://www.tensorflow.org/api\\_docs/python/tf/keras/](https://www.tensorflow.org/api_docs/python/tf/keras/)