

Configuration Manual

MSc Research Project

Data Analytics

Rishab Rao Gauravaram

Student ID: x19202504

School of Computing

National College of Ireland

Supervisor: Dr Catherine Mulwa

National College of Ireland

National College of Ireland

MSc Project Submission Sheet

School of Computing

Student Name: Rishab Rao Gauravaram

Student ID: X19202504

Programme: MSc Data Analytics

Module: MSc Data Analytics – Research Project

Lecturer:Dr. Catherine MulwaSubmission Due Date:16 August 2021

Word Count: 2193

Page Count: 39

Year: 2020-2021

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.

<u>ALL</u> 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: 23/09/2021 Date:

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

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

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

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Rishab Rao Gauravaram Student ID: x19202504

1 Introduction

This configuration manual will illustrate the environment setup and the configurations for the same to implement "Prediction and Classification of Electrocardiogram-Signals using Machine Learning using Apache Spark".

2 Environment Specification and Configuration

2.1 Hardware Configurations

The hardware configuration of the system used for this project is shown in the Figure

Device spec	ifications
Legion Y540 Device name	D-15IRH-PG0 LAPTOP-8UKBO5GP
Processor	Intel(R) Core(TM) i5-9300H CPU @ 2.40GHz 2.40 GHz
Installed RAM	16.0 GB (15.9 GB usable)
Device ID	76309F6C-1D3C-4316-B226-5633C4A2E562
Product ID	00327-35886-09734-AAOEM
System type	64-bit operating system, x64-based processor
Pen and touch	No pen or touch input is available for this display

2.2 Software Specifications and Requirements

1. To download the Apache ecosystem on a Virtual Machine, VirtualBox application was installed with Ubuntu Operating System.

The download link for VirtualBox is as follows – https://download.virtualbox.org/virtualbox/6.1.26/VirtualBox-6.1.26-145957-Win.exe

The download link for the Ubuntu OS execution file is as follows – <u>https://ubuntu.com/download/desktop/thank-you?version=20.04.2.0&architecture=amd64</u>

2. After downloading both the files, the following steps will provide a detailed description to setup VirtualBox and Ubuntu on the VirtualBox. After downloading the execution file for VirtualBox, run the execution file and then open it.

File Machine Help		
Tools	New Settings Discard Show	
DAP_Dec20	General Name: Ubuntu Operating System: Ubuntu (64-bit)	Preview
Hadoop Wered Off	System Base Memory: 11264 MB Processors: 4 Boot Order: Floow. Optical. Hard Disk	 2 I and an an
Vbuntu	Acceleration: VT-x/AMD-V, Nested Paging, KVM Paravirtualization Video Memory: 16 MB Graphics Controller: VMSVCA Remote Desktop Server: Disabled Decording: Disabled	
	Controller: IDE IDE Secondary Master: [Optical Drive] Empty Controller: SATA SATA Port 0: Ubuntu.vdi (Normal, 100.00 GB)	
	Audio Host Driver: Windows DirectSound Controller: ICH AC97	
	Retwork Adapter 1: Intel PRO/1000 MT Desktop (NAT) Adapter 2: Intel PRO/1000 MT Desktop (Host-only Adapter, "VirtualBox Host-Only Ethernet Adapter")	
	🖉 IISB	

3. VirtualBox Setup -

Click the "**New**" button. This is used to create a virtual machine as shown. Then press the "**Create**" button. Assign the Memory Size. In this case 11GB Ram was used.

		?	×
Create Virtual Machine			
Memory size			
Select the amount of memory (virtual machine.	RAM) in megabyt	es to be alloca	ted to the
The recommended memory size	e is 1024 MB.		
		11	264 🌻 ME
4 MB		16384 MB	
	N	lext	Cancel

The next step is to create the hard disk. Select the option as shown. Then press "Next".

	?	×
← Create Virtual Machine		
Hard disk		
If you wish you can add a virtual hard disk to the new either create a new hard disk file or select one from th another location using the folder icon.	m <mark>achine</mark> . Y ne list or fro	ou can om
If you need a more complex storage set-up you can sk make the changes to the machine settings once the m	tip this step achine is cr	and eated.
The recommended size of the hard disk is 10.00 GB .		
O Do not add a virtual hard disk		
Create a virtual hard disk now		
O Use an existing virtual hard disk file		
DAP_Dec20-disk001.vdi (Normal, 41.73 GB)		*
Create	Ca	ancel

Select VirtualBox Disk Image and press "Next"

		?	×
← Create Virtual Hard Dis	k		
Hard disk file type			
Please choose the type of file disk. If you do not need to us this setting unchanged.	that you would like to use e it with other virtualizatio	e for the new virtua n software you ca	al hard n leave
VDI (VirtualBox Disk Imag	e)		
○ VHD (Virtual Hard Disk)			
O VMDK (Virtual Machine Di	sk)		
	Expert Mode	Next	ancel

Choose the storage type. If the host device has less storage, then choose "dynamically allocated" option, else, choose "fixed size"

	?	×
 Create Virtual Hard Disk 		
Storage on physical hard disk		
Please choose whether the new virtual hard disk file should grov (dynamically allocated) or if it should be created at its maximum	v as it is us size (fixed	sed I size).
A dynamically allocated hard disk file will only use space on y disk as it fills up (up to a maximum fixed size), although it will r automatically when space on it is freed.	our physica not shrink a	al hard again
A fixed size hard disk file may take longer to create on some sy faster to use.	/stems but	is often
O Dynamically allocated		
Fixed size		
Next	Ca	incel

Select the file path for the Virtual Machine to be stored and storage size

		?	2
Create Virtual Hard Disk			
File location and size			
Please type the name of the new virtual hard di on the folder icon to select a different folder to	sk file into the box b create the file in.	elow or	click
C:\Users\risha\VirtualBox VMs\ubuntu\ubuntu.vo	di		
Select the size of the virtual hard disk in megab amount of file data that a virtual machine will be	ytes. This size is the a able to store on th	e limit o Ie hard	n the disk.
			100
4.00 MB	2.00 TB		100
4.00 MB	2.00 TB		100
4.00 MB	2.00 TB		100
4.00 MB	2.00 TB		100

After the virtual machine is created, open "Settings" and select the options as shown in the figure

-General

General	General				
System	Basic	Advanced	Description	Disk Encryption	
Display	Snapshot	Folder:	D:\Virtualbox V	/M\Ubuntu\Snapshots	Ŷ
Storage	Shared Clip	oboard: Bio	directional 👻		
Audio	Drag'ı	n'Drop: Bio	directional 🔻		
Network					
Serial Ports					
USB					
USB Shared Folders					
USB Shared Folders User Interface					
USB Shared Folders User Interface					
USB Shared Folders User Interface					

-System Storage



On the top right corner select the "empty disk" option and choose the Ubuntu ISO file previously downloaded.

Ubuntu - Settings				?	×
General	Storage				
 System Display Storage Audio Network Serial Ports USB Shared Folders User Interface 	Storage Devices Controller: IDE Controller: SATA Vubuntu.vdi	Attributes Optical Drive: Information Type: Size: Location: Attached to:	IDE Second	dary Mast	× @

-Network

General	Network		
System	Adapter 1 Ada	apter 2 Adapter 3 Adapter 4	
Display	Enable Networ	k Adapter	
Storage	Attached to:	Host-only Adapter 🔹	
Storage	Name:	VirtualBox Host-Only Ethernet Adapter	
USB			
Shared Folders			
Shared Folders User Interface			

4. Ubuntu Setup -

After VirtualBox settings configuration, close it and select the "Start" option. This will start the virtual machine.



Select the choice of language and click "Install Ubuntu"



Select the keyboard layout style

	Inst		
Keyboard layout			
Choose your keyboard layout: Dzongkina English (Australian) English (Cameroon) English (Ghana) English (Nigeria) English (South Africa) English (UK) English (US) Esperanto	k	English (US) English (US) - Cherokee English (US) - English (Colemak) English (US) - English (Dvorak) English (US) - English (Dvorak, alt. intl.) English (US) - English (Dvorak, intl., with dead keys) English (US) - English (Dvorak, Ieft-handed) English (US) - English (Dvorak, right-handed)	
Type here to test your keyboard Detect Keyboard Layout			
		Quit Back Continu	e

Select all the checkboxes except Minimal installation.

Ubuntu 20.04 [Running] - Oracle VM VirtualBox				-		×
Apr 24 2	:0:53				* •	1
Instal	u					(
Updates and other software						
What apps would you like to install to start with? Normal installation We have within a fire after a ferrer or and a dia data.						
Web browser, utilities, orrice sortware, games, and media players. Minimal installation Web browser and basic utilities						
Other options						
Download updates while installing Ubuntu This saves time after installation.	k					
Install third-party software for graphics and Wi-Fi hardw	vare and ad	ditional med	lia forma	its		
This software is subject to license terms included with its document	tation. Some	is proprietary.				
		Quit	Bac	k	Contir	ue
	000					
					STRG-R	ECHT

Choose the installation type.



The following screen will request for disk formatting This will **not** make any changes to the host system. The clean-up and formatting are **only** for the virtual machine.

💕 Ubuntu 20.04 [Running] - Oracle VM VirtualBox	- 🗆 ×
File Machine View Input Devices Help	
Apr 24 20:54	₽ • •
Install	
Installation type	
This computer currently has no detected operating systems. What w	would you like to do?
Write the changes to di	isks?
If you continue, the changes listed below yill be written to the disk changes manually.	ss. Otherwise, you will be able to make further
The partition tables of the following devices are changed: SCSI3 (0,0,0) (sda)	
The following partitions are going to be formatted: partition #1 of SCS13 (0,0,0) (sda) as partition #5 of SCS13 (0,0,0) (sda) as ext4	
	Go Back Continue
	Back Install Now
	0
Q	💿 💯 🚅 🧷 🔲 🔲 🚰 🔯 🐼 STRG-RECHTS

Enter the username and password. DO NOT FORGET THE LOGIN DETAILS. SAVE IT SOMEWHERE IF NECESSARY

		Apr 24 22:	6		÷ •	+
		Install				
Who are you	?					
		T				
	Your name:					
Y	our computer's name:	The serve it uses whe	a ik kallas karakkas asa			
	Pick a username.	The name it uses whe	n it taks to other cor	nputers.		
	Choose a password:					
	choose a password.					
	onrirm your password:	Oberlauter	et an Iba			
		 Log in automa Require my pa 	ssword to log in			
			-			
				Back	Contir	nue

The installation will commence. It will take some time for the OS to completely installed.



Post installation the virtual machine will request for a reboot. Press "Restart Now"



After the system restarts, enter the login details and open the terminal by pressing "Ctrl + Alt + t" and enter the following command as shown in the figure. This will install all essential Ubuntu packages.

```
rishab@rishab:~$ sudo apt install build-essential dkms linux-headers-$(uname -r)
```

2.3 Hadoop Installation

- Open the Terminal (Ctrl + Alt + t) and update Ubuntu OS using the commands sudo apt-get update sudo apt-get upgrade
- 2. Install Java-8 jdk using the following command sudo apt install openjdk-8-jdk
- 3. Check where Java is installed (Java Path). This is required to configure Hadoop.

```
rishab@rishab:~$ sudo update-alternatives --config java
[sudo] password for rishab:
There is only one alternative in link group java (providing /usr/bin/java):
    /usr/lib/jvm/java-8-openjdk-amd64/jre/bin/java
Nothing to configure.
```

4. Open the /etc/profile file using the command – sudo nano /etc/profile

Add the following lines of code to set the path for Java. And then save the file using the command -

source /etc/profile



 Open the sysctl.conf file using the command – sudo nano /etc/profile/sysctl.conf Add the following lines of code.

- 6. To make these changes to take immediate effect, reboot the system.
- 7. Configure the SSH for Hadoop. Create a user group for Hadoop using the following commands –

```
rishab@rishab:~$ sudo addgroup hadoopgroup
```

rishab@rishab:~\$ sudo adduser -ingroup hadoopgroup hduser

8. Install SSH

rishab@rishab:~\$ sudo apt-get install ssh

Enable SSH

rishab@rishab:~\$ sudo systemctl enable ssh

Start SSH

rishab@rishab:~\$ sudo systemctl start ssh

9. Switch user by typing the following command



10. Generating public/private keys folder

```
hduser@rishab:~$ ssh-keygen -t rsa -P ""
Generating public/private rsa key pair.
Enter file in which to save the key (/home/hduser/.ssh/id_rsa):
```

11. Copy the keys to the authorized keys folder

hduser@rishab:~\$ cat /home/hduser/.ssh/id_rsa.pub >> /home/hduser/.ssh/authorized_keys

12. Change the permissions for the keys

```
hduser@rishab:~$ cd .ssh/
hduser@rishab:~/.ssh$ chmod 666 ./authorized_keys
```

13. Download Hadoop 3.3.0 using the command

hduser@rishab:-\$ wget https://archive.apache.org/dist/hadoop/common/hadoop-3.3.0/hadoop-3.3.0.tar.gz

- 14. Extract the contents of the Hadoop zip file using the command tar -xvf hadoop 3.3.0.tar.gz
- 15. Move the extracted file into **cd/usr/local.** Change ownership using the following command -

sudo chown -R hduser:hadoopgroup /usr/local/Hadoop-3.3.0

16. Create a symbolic link with hadoop using the command

```
hduser@rishab:~$ sudo ln -sf /usr/local/hadoop-3.3.0/ /usr/local/hadoop/
```

17. Open the **bashrc file** using the command sudo nano ./.bashrc and configure it using the following lines of code. Then save the file using **source**



18. Open the Hadoop environment to set the java path



19. Navigate to the hadoop directory using **cd** /usr/local/hadoop/etc/hadoop and navigate the .xml files and add the following –

sudo nano core-site.xml



sudo nano mapred-site.xml

```
<configuration>
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>
```

sudo nano yarn-site.xml



sudo nano hdfs-site.xml



- 20. Navigate to the home directory using the command cd
- 21. Format the namenode

hduser@rishab:~\$ sudo hdfs namenode -format

22. Start Hadoop services by first navigating to the Hadoop directory **cd /usr/local/hadoop**

shab:/usr/local/hadoop\$ start-dfs.sh WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. Starting namenodes on [localhost] WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. Starting datanodes WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. Starting secondary namenodes [rishab] WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. hduser@rishab:/usr/local/hadoop\$ start-yarn.sh WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. Starting resourcemanager WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. Starting nodemanagers WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX. hduser@rishab:/usr/local/hadoop\$ jps 11136 DataNode 11345 SecondaryNameNode 11683 NodeManager 12024 Jps 11533 ResourceManager 10959 NameNode

2.4 Apache Spark Installation

- 1. Download and extract Apache Spark 3.1.2 in **cd /usr/local** using the link <u>https://www.apache.org/dyn/closer.lua/spark/spark-3.1.2/spark-3.1.2-bin-hadoop3.2.tgz</u>
- 2. Create a symbolic link of spark using the command sudo ln –s spark-3.1.2-bin-hadoop3.2 spark
- 3. Change the ownership of the Spark directory to Hadoop group and ownership to hduser as seen in section 2.3 (sudo chown -R hduser:hadoopgroup spark/)
- 4. Configure the .bashrc file using the following lines of code and save it using source



5. Navigate to cd/usr/local/spark and start the services using sbin/start-master.sh

2.5 Anaconda Installation

1. Download Anaconda using the following link: https://repo.anaconda.com/archive/Anaconda3-2021.05-Linux-x86_64.sh 2. Once the download is complete, extract it using the command:

hduser@rishab:~\$ bash ~/Downloads/Anaconda3-2021.05-Linux-x86_64.sh

Once the installation commences, read the terms and conditions, and answer "yes" to everything. Keep in mind to answer "yes" when the installer requests permission to prepend anaconda3 path to the .bashrc file.

3. Then source the .bashrc file for the changes to take immediate effect The terminal should look like the following:

(base) hduser@rishab:~\$

- 4. To implement Deep Neural Networks, create a separate environment on Anaconda and install TensorFlow and Keras.
- First step is to install conda forge. This is used to download and install TensorFlow. Use the following commands on the terminal after installing Anaconda: conda config –add channels conda-forge conda config –set channel_priority strict
- 6. Before installing Keras and TensorFlow a new environment must be created. Open anaconda navigator

(base) hduser@rishab:~\$ anaconda-navigator

Go to environments and create a new environment named "TensorFlow".

7. After the new environment is created, select the TensorFlow environment and install TensorFlow Package

🗹 keras	营 Deep learning library for theano and tensorflow	2.2.5
opt_einsum	O Optimizing einsum functions in numpy, tensorflow, dask, and more with contraction order optimization.	3.3.0
d tensorboard	Q Tensorflow's visualization toolkit	↗ 2.4.0
tensorboard-plugin- wit	0	1.6.0
densorflow	O Tensorflow is a machine learning library.	↗ 2.4.1
densorflow-base	O Tensorflow is a machine learning library, base package contains only tensorflow.	↗ 2.4.1
estimator	O Tensorflow estimator is a high-level tensorflow api that greatly simplifies machine learning programming.	2.5.0

3 Implementation of the research project

3.1 Download the dataset

- a) Download the wfdb package using the following line on the terminal: **pip install wfdb**
- b) Open a new python file using the terminal as shown

(base) hduser@rishab:~/Final\$ nano Data_Download.py

c) Type the following lines of code to download the data

F	hduser@rishab: ~/Final	Q	Ξ		×
GNU nano 4.8	Data_Download.py				
<pre>import wfdb # It will take some # The dataset will b def download(): wfdb.dl_database</pre>	time to download the dataset e downloaded inside the directory ('mitdb', dl_dir = 'data')				
ifname == "ma download()	in":				

d) Save the file and exit the text editor. Run the python script using the following line:

(base) hduser@rishab:~/Final\$ python Data_Download.py

e) The data will be downloaded under the folder name "Data". The contents of the folder will look like the following:



3.2 Data Preprocessing, Transformation and Exploratory Data Analysis

a) Open a new python file using a text editor (nano, vim, vi, gedit, etc) in the terminal for Data Processing, Data Transformation, Signal extraction and EDA.

nano preprocess.py

b) Import the following packages

```
1 import numpy as np
2 import sys, os
3 import wfdb
4 import matplotlib.pyplot as plt
5 import pywt
6 import pickle as pk
7 from collections import Counter
8 import time
9 import pandas as pd
10 start_time = time.time()
```

c) The names of the data files, five types of classes, different types of beats and different types of beats assigned to each class are stored in different lists and dictionaries.

```
12 #The data names, beats and sub beats are annotated as prescribed by the authors of the database
13 #The link for the annotations - <u>https://archive.physionet.org/physiobank/annotations.shtml</u>
14

      14

      15

      15

      16

      '108', '109', '111', '112', '113', '114', '115', '116',

      17

      '117', '118', '119', '121', '122', '123', '124', '200',

      18
      '201', '202', '203', '205', '207', '208', '209', '210',

      19
      '212', '213', '214', '215', '217', '219', '220', '221',

      20
      '222', '223', '228', '230', '231', '232', '233', '234']

21
22
'E':'E'
28
                       '/':'0', 'f':'0', '0':'0'}
29
30
31 #Empty lists to store features and labels
32 X = []
33 Y = []
34 num = 100
```

d) Data Preprocessing using Denoising¹ (DWT for decomposition and IDWT for reconstruction of the wavelets), normalising the wavelets and extracting the waveform for five different classes of heartbeats.



e) Exploratory Data Analysis on the Data

```
89 #EDA and Visualising the extracted signals
90 #Find the number of records in each class of the heartbeats
91 values, counts = np.unique(Y, return_counts=True)
91 Values, counts = np.unque(r, recurr_counts=rroc)

92 print(values, counts)

93 plt.figure(figstze=(20,18))

94 my_circle=plt.circle( (0,0), 0.7, color = 'white')

95 plt.pie(counts, labels=['W', 'S', 'V', 'F', 'Q'], colors=['#ff9999', '#66b3ff', '#99ff99', '#ffcc99', '#c2c2f0'],autopct='%2.2f%%', textprops={'fontsize': 18})

96 plt.title('Different classes of heartbeats', fontsize = 25)
97 p=plt.qcf()
98 #p.gca().add_artist(my_circle)
99
100 plt.show()
                                   104 #Plotting the graphs for each class of heartbeat
                                    105 values, indexes = np.unique(Y, return_index=True)
                                    106 for i in indexes:
                                    107
                                                plt.plot(X[i])
                                   108
                                                 plt.xlabel("beat sample")
                                    109
                                                 plt.ylabel("time")
                                                 if Y[i]==0:
                                    110
                                    111
                                                         plt.title("Class N")
                                               elif Y[i]==1:
                                    112
                                               plt.title("Class S")
elif Y[i]==2:
                                    113
                                    114
                                   115 plt.title("Class V")
116 elif Y[i]==3:
                                    117
                                                       plt.title("Class F")
                                    118
                                                 else:
                                    119
                                                       plt.title("Class Q")
```

plt.show()

120

¹ https://stackoverflow.com/questions/32184232/filter-ecg-signal-with-wavelet

f) Storing the transformed and cleaned data into a .csv that will be stored on the HDFS.

```
123 #Data to run ML algorithms on Spark
124 Features = pd.DataFrame(X)
125 #print(Features)
126
127 Labels = pd.DataFrame(Y)
128 #print(Labels)
129
130 ECG = pd.merge(Features, Labels, right_index = True, left_index = True)
131 #print(ECG)
132 ECD_Data = ECG.to_csv(r'ECG.csv', index=False)
133
```

g) Storing the transformed and cleaned data in a .pk file that can be used to implement models locally

```
136 #Data to run ML algorithms locally
137 fn = "ECG_Data"+".pk"
138 with open(fn, "wb") as fw:
139      pk.dump(X, fw, protocol=pk.HIGHEST_PROTOCOL)
140      pk.dump(Y, fw, protocol=pk.HIGHEST_PROTOCOL)
```

h) Save and run the python file using python preprocess.py

OUTPUTS -







3.3 Experiment 1 – Local implementation of the classification and prediction models

3.3.1 Implementation of Deep Neural Network

- a) On the terminal change the environment using **conda activate TensorFlow.** In this research project the name of the environment is TensorFlow.
- b) Open a python file on the terminal to using a text editor to implement DNN

```
(base) hduser@rishab:~$ conda activate TensorFlow
(TensorFlow) hduser@rishab:~$ nano DNN.py
```

c) Import the required libraries

```
1 from matplotlib import pyplot
2 import numpy as np
 3 import pickle as pk
 4 import os, sys
 5 from collections import Counter
 6 from sklearn.svm import SVC
 7 from sklearn.metrics import confusion matrix
 8 from sklearn.metrics import roc_auc_score
 9 from sklearn.model_selection import GridSearchCV
10 import time
11 import tensorflow as tf
12 from tensorflow.keras.models import Sequential
13 from tensorflow.keras.layers import *
14 from tensorflow.keras.utils import to_categorical
15 from tensorflow.keras import backend as K
16
17 from tensorflow.keras import callbacks
18 from tensorflow.keras.optimizers import Adam
19
20
21 start_time = time.time()
```

d) Load the dataset

```
23 #Load the pickled data
24 file = "ECG_Data"+".pk"
25 with open(file, "rb") as fn:
26   X = pk.load(fn)
27   Y = pk.load(fn)
```

e) Split the data into train data and test data in the ratio of 70:30

```
29 #Randomly split the Data into Train and Test in the ratio 70:30
30 r_seed = 1229
31 data_len = len(X)
32 np.random.seed(r_seed)
33 idx = list(range(data len))
34 np.random.shuffle(idx)
35
36 train_len = int(data_len*0.7)
37 test_len = int(data_len*0.3)
38
39 X_train = X[idx][:train_len]
40 X_test = X[idx] train_len:
41 Y_train = Y[idx][:train_len]
42 Y_test = Y[idx][train_len:]
43
44 X_train = np.expand_dims(X_train, axis=-1)
45 X_test = np.expand_dims(X_test, axis=-1)
46
47 print(X_train.shape)
48 print(X_test.shape)
49 print(Counter(Y_train))
50 print(Counter(Y_test))
```

f) Create the DNN model

```
52 #Input shape and number of classes for classification
53 f_size = X_train.shape[1]
54 class_num = 5
55
56 #Learning rate and batch size
57 lr = 0.005
58 batch size=128
59
60 #Convert the labels inot categorical values
61 Y_train = tf.keras.utils.to_categorical(Y_train, num_classes=class_num)
62
63 #Create the DNN model, using 4 layers and compile it
64 def make_model():
65
      model = Sequential()
      model.add(Conv1D(18, 7, activation='relu', input_shape=(f_size,1)))
66
67
      model.add(MaxPooling1D(2))
      model.add(Conv1D(18, 7, activation='relu'))
model.add(MaxPooling1D(2))
68
69
70
      model.add(Flatten())
      model.add(Dense(100, activation='relu'))
71
      model.add(Dense(class_num, activation='softmax'))
72
      model.compile(loss='categorical_crossentropy', optimizer=Adam(lr=lr))
73
74
      return model
75
76 model = make_model()
```

g) Train and test the model and then evaluate it using accuracy, specificity, sensitivity, and execution time.

```
78 bin label = lambda x: min(1,x)
 79
 80 #for e in range(1, 300+1):
81 for e in range(1, 50+1):
82 #Train the model for 50 epochs
       model.fit(X_train, Y_train, batch_size=batch_size, epochs=1, verbose=0)
 83
84
 85
       #Fit the model on the test dataset
 86
        y_pred = model.predict(X_test)
 87
       y_pred = np.argmax(y_pred, axis=1)
 88
       #Evaluation
 89
 90
       acc = np.sum(y_pred==Y_test)/len(Y_test)
 91
 92
       #Map the predicted labels and the test labels for evaluation
       y_true = list(map(bin_label, Y_test))
 93
       y_pred = list(map(bin_label, y_pred))
 94
95
96
       #Extracting true positive, true negative, false positive and false negative
97
       #from the confusion matrix
 98
       tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()
       SE = tp/(tp+fn)
99
100
       SP = tn/(fp+tn)
101
102
       print("Epoch: %d | SE: %.4f | ACC %.4f | AUC %.4f | SP: %.4f"%(e, SE, acc, auc, SP))
103
104 #Final Result
105 print("Final Result: \n")
106 print("SE: %.4f | ACC: %.4f | AUC: %.4f | SP: %.4f " %(SE, acc, auc, SP))
107
108 print("--- %s seconds ---" % (time.time() - start_time))
```

h) Close the python file and save it. Run the program using **python DNN.py**

OUTPUT –

Epoch: 1	SE: 0.8145	ACC 0.9590	SP: 0.9935
Epoch: 2	SE: 0.8917	ACC 0.9694	SP: 0.9906
Epoch: 3	SE: 0.8783	ACC 0.9737	SP: 0.9970
Epoch: 4	SE: 0.9068	ACC 0.9752	SP: 0.9925
Epoch: 5	SE: 0.9171	ACC 0.9809	SP: 0.9965
Epoch: 6	SE: 0.9119	ACC 0.9800	SP: 0.9969
Epoch: 7	SE: 0.9404	ACC 0.9842	SP: 0.9960
Epoch: 8	SE: 0.9298	ACC 0.9818	SP: 0.9958
Epoch: 9	SE: 0.9308	ACC 0.9840	SP: 0.9974
Epoch: 10	SE: 0.9362	ACC 0.9843	SP: 0.9969
Epoch: 11	SE: 0.9350	ACC 0.9833	SP: 0.9959
Epoch: 12	SE: 0.9395	ACC 0.9838	SP: 0.9965
Epoch: 13	SE: 0.9444	ACC 0.9851	SP: 0.9964
Epoch: 14	SE: 0.9473	ACC 0.9794	SP: 0.9889
Epoch: 15	SE: 0.9503	ACC 0.9854	SP: 0.9956
Epoch: 16	SE: 0.9503	ACC 0.9861	SP: 0.9964
Epoch: 17	SE: 0.9383	ACC 0.9851	SP: 0.9973
Epoch: 18	SE: 0.9332	ACC 0.9844	SP: 0.9973
Epoch: 19	SE: 0.9505	ACC 0.9849	SP: 0.9946
Epoch: 20	SE: 0.9181	ACC 0.9825	SP: 0.9987
Epoch: 21	SE: 0.9426	ACC 0.9854	SP: 0.9967
Epoch: 22	SE: 0.9392	ACC 0.9848	SP: 0.9968
Epoch: 23	SE: 0.9496	ACC 0.9861	SP: 0.9963
Epoch: 24	SE: 0.9470	ACC 0.9847	SP: 0.9956
Epoch: 25	SE: 0.9494	ACC 0.9865	SP: 0.9970
Epoch: 26	SE: 0.9546	ACC 0.9862	SP: 0.9955
Epoch: 27	SE: 0.9487	ACC 0.9846	SP: 0.9960
Epoch: 28	SE: 0.9564	ACC 0.9857	SP: 0.9946
Epoch: 29	SE: 0.9565	ACC 0.9873	SP: 0.9966
Epoch: 30	SE: 0.9604	ACC 0.9875	SP: 0.9955
Epoch: 31	SE: 0.9609	ACC 0.9867	SP: 0.9945
Epoch: 32	SE: 0.9445	ACC 0.9866	SP: 0.9979
Epoch: 33	SE: 0.9600	ACC 0.9866	SP: 0.9955
Epoch: 34	SE: 0.9574	ACC 0.9870	SP: 0.9959
Epoch: 35	SE: 0.9604	ACC 0.9873	SP: 0.9953
Epoch: 36	SE: 0.9463	ACC 0.9848	SP: 0.9959
Epoch: 37	SE: 0.9489	ACC 0.9873	SP: 0.9975
Epoch: 38	SE: 0.9482	ACC 0.9856	SP: 0.9960
Epoch: 39	SE: 0.9616	ACC 0.9878	SP: 0.9958
Epoch: 40	SE: 0.9541	ACC 0.9873	SP: 0.9968
Epoch: 41	SE: 0.9564	ACC 0.9874	SP: 0.9964
Epoch: 42	SE: 0.9489	ACC 0.9869	SP: 0.9971
Epoch: 43	SE: 0.9571	ACC 0.9868	SP: 0.9955
Epoch: 44	SE: 0.9656	ACC 0.9872	SP: 0.9945
Epoch: 45	SE: 0.9565	ACC 0.9850	SP: 0.9939
Epoch: 46	SE: 0.9628	ACC 0.9868	SP: 0.9946
Epoch: 47	SE: 0.9626	ACC 0.9879	SP: 0.9954
Epoch: 48	SE: 0.9611	ACC 0.9862	SP: 0.9941
Epoch: 49	SE: 0.9484	ACC 0.9873	SP: 0.9976
Epoch: 50	SE: 0.9701	T ACC 0.9869	SP: 0.9929
Final Rest			
SE. 0 0701		50 I SP. 0 000	20
JE. 0.9701	538135051727	seconds	
1057.1	330133031121	Seconds	

3.3.2 Implementation of Random Forest model

- a) On the terminal return to base environment and then open a python file to implement RF.
 - (TensorFlow) hduser@rishab:~\$ conda activate
 (base) hduser@rishab:~\$ nano RF.py
- b) Import the required packages

```
1 import pandas as pd
2 import numpy as np
3 import pickle as pk
4 import os, sys
5 import seaborn as sns
6 import time
7 from sklearn.model_selection import cross_val_score, KFold, StratifiedKFold
8 from sklearn.metrics import *
9 from collections import Counter
10 from sklearn.ensemble import RandomForestClassifier
11 import time
12
13 start_time = time.time()
```

c) Split the data into train and test in the ration of 80:20

```
15 #Load the pickled data
16 file = "ECG_Data"+".pk'
17 with open(file, "rb") as fn:
      X = pk.load(fn)
18
      Y = pk.load(fn)
19
20
21 #Randomly split the Data into Train and Test in the ratio 80:20
22 r_seed = 1229
23 data len = len(X)
24 np.random.seed(r_seed)
25 idx = list(range(data len))
26 np.random.shuffle(idx)
27
28 train_len = int(data_len*0.8)
29 test_len = int(data_len*0.2)
30
31 X_train = X[idx][:train_len]
32 X_test = X[idx][train_len:]
33 Y_train = Y[idx][:train_len]
34 Y_test = Y[idx][train_len:]
35
36 print(X_train.shape)
37 print(X_test.shape)
38 print(Counter(Y_train))
39 print(Counter(Y_test))
```

d) Train and test the model and, then evaluate it using accuracy, specificity, sensitivity, and execution time.

```
41 #Random Forest Classification Model. Fit the model on the training set
42 rf = RandomForestClassifier(n estimators=100)
43 rf_model = rf.fit(X_train, Y_train)
44
45 #Test the trained model on the training dataset
46 y_pred = rf_model.predict(X_test)
47
48 #Evaluation
49 acc = accuracy_score(Y_test, y_pred)
50 conf_mat = confusion_matrix(Y_test, y_pred)
51 print(type(conf_mat))
52 print(conf_mat)
53
54 #Extracting true positive, true negative, false positive and false negative
55 #from the confusion matrix
56 fp = conf_mat.sum(axis=0) - np.diag(conf_mat)
57 fn = conf_mat.sum(axis=1) - np.diag(conf_mat)
58 tp = np.diag(conf mat)
59 tn = conf_mat.sum() - (fp + fn + tp)
60
61 #Converting the values into float values
62 fp = fp.astype(float)
63 fn = fn.astype(float)
64 tp = tp.astype(float)
65 tn = tn.astype(float)
66
67 #Computing sensitivity and specificity from tn, tp, fn, fp
68 SE = tp/(tp+fn)
69 SP = tn/(fp+tn)
70 SE = np.average(SE)
71 SP = np.average(SP)
72
73 print(" SE: %.4f | ACC: %.4f | SP: %4f" %(SE, acc, SP))
74
75 print("--- %s seconds ---" % (time.time() - start time))
```

e) Save and close the python file. Run the program using python RF.py

OUTPUT -

SE: 0.8362 | ACC: 0.9809 | SP: 0.981380 --- 359.75102639198303 seconds ---

3.3.3 Implementation of Support Vector Machine model

a) Open a python file on the terminal to implement SVM.

```
(base) hduser@rishab:~$ nano SVM.py
```

b) Import the libraries

```
1 import numpy as np
 2 import pandas as pd
 3 import pickle as pk
 4 import os, sys
 5 import seaborn as sns
6 import time
 7 from sklearn.model_selection import cross_val_score, KFold, StratifiedKFold
8 from sklearn.model_selection import GridSearchCV
9 from sklearn.svm import LinearSVC
10 from sklearn.model_selection import RandomizedSearchCV
11 from scipy.stats import reciprocal, uniform
12 from sklearn.svm import SVC
13 from sklearn.metrics import *
14 from collections import Counter
15
16 start_time = time.time()
17
```

c) Split the data into train and test in the ration of 80:20

```
15 #Load the pickled data
16 file = "ECG Data"+".pk"
17 with open(file, "rb") as fn:
      X = pk.load(fn)
18
      Y = pk.load(fn)
19
20
21 #Randomly split the Data into Train and Test in the ratio 80:20
22 r_seed = 1229
23 data_len = len(X)
24 np.random.seed(r_seed)
25 idx = list(range(data_len))
26 np.random.shuffle(idx)
27
28 train_len = int(data_len*0.8)
29 test_len = int(data_len*0.2)
30
31 X_train = X[idx][:train_len]
32 X_test = X[idx][train_len:]
33 Y_train = Y[idx][:train_len]
34 Y_test = Y[idx][train_len:]
35
36 print(X_train.shape)
37 print(X_test.shape)
38 print(Counter(Y_train))
39 print(Counter(Y_test))
```

d) Train and test the model and, then evaluate it using accuracy, specificity, sensitivity, and execution time

```
44 #Support Vector Machine Classification Model. Fit the model on the training dataset
45 svm_clf = SVC(gamma="scale")
46 svm_clf.fit(X_train, Y_train)
47
48 #Hyperparameter tuning og C and gamma variable to prevent overfitting
49 param_distributions = {"gamma": reciprocal(0.001, 0.1), "C": uniform(1, 10)}
50
50
51 #3-fold cross-validation, iterated 10 times
52 rnd_search_cv = RandomizedSearchCV(svm_clf, param_distributions, n_iter=10, verbose=2, cv=3)
53 rnd_search_cv.fit(X_train, Y_train)
55 #Fitting the best estimator on the Training dataset
56 rnd_search_cv.best_estimator_.fit(X_train, Y_train)
57
58 #Fit the trained model on the test data
59 y_pred = rnd_search_cv.best_estimator_.predict(X_test)
60
61 #Evaluation
62 acc = accuracy_score(Y_test, y_pred)
63 conf_mat = confusion_matrix(Y_test, y_pred)
64
65 #Extracting true positive, true negative, false positive and false negative
66 fp = conf_mat.sum(axis=0) - np.diag(conf_mat)
67 fn = conf_mat.sum(axis=1) - np.diag(conf_mat)
68 tp = np.diag(conf_mat)
69 tn = conf_mat.sum() - (fp + fn + tp)
70
71 #Converting the values into float values
72 fp = fp.astype(float)
73 fn = fn.astype(float)
74 tp = tp.astype(float)
75 tn = tn.astype(float)
76
77 #Computing sensitivity and specificity from tn, tp, fn, fp
78 SE = tp/(tp+fn)
79 \text{ SP} = \text{tn}/(\text{fp+tn})
80 SE = np.average(SE)
81 SP = np.average(SP)
82
83 print(" SE: %.4f | ACC: %.4f | SP: %4f" %(SE, acc, SP))
84
85 print("--- %s seconds ---" % (time.time() - start_time))
```

e) Save the file and close it. Execute the file using python SVM.py

OUTPUT -

Fitting 3 folds for each of 10 candidates, totalling 30 fits	
<pre>[CV] ENDC=1.9493313853883334, gamma=0.00355224454904102;</pre>	total time= 5.1min
<pre>[CV] ENDC=1.9493313853883334, gamma=0.00355224454904102;</pre>	total time= 4.9min
<pre>[CV] ENDC=1.9493313853883334, gamma=0.00355224454904102;</pre>	total time= 5.1min
[CV] ENDC=3.970768395498694, gamma=0.07180971082644513;	total time= 2.4min
[CV] ENDC=3.970768395498694, gamma=0.07180971082644513;	total time= 2.4min
[CV] ENDC=3.970768395498694, gamma=0.07180971082644513;	total time= 2.3min
[CV] ENDC=1.9793791377756678, gamma=0.006850493921060555;	total time= 3.7min
[CV] ENDC=1.9793791377756678, gamma=0.006850493921060555;	total time= 3.6min
[CV] ENDC=1.9793791377756678, gamma=0.006850493921060555;	total time= 3.8min
[CV] ENDC=4.702307178313259, gamma=0.04682148275809945;	total time= 2.2min
[CV] ENDC=4.702307178313259, gamma=0.04682148275809945;	total time= 2.0min
[CV] ENDC=4.702307178313259, gamma=0.04682148275809945;	total time= 2.2min
<pre>[CV] ENDC=9.306955424014042, gamma=0.0026964216806020817;</pre>	total time= 3.5min
<pre>[CV] ENDC=9.306955424014042, gamma=0.0026964216806020817;</pre>	total time= 3.5min
<pre>[CV] ENDC=9.306955424014042, gamma=0.0026964216806020817;</pre>	total time= 3.6min
<pre>[CV] ENDC=8.805966025659725, gamma=0.024961736585479025;</pre>	total time= 2.2min
<pre>[CV] ENDC=8.805966025659725, gamma=0.024961736585479025;</pre>	total time= 2.1min
<pre>[CV] ENDC=8.805966025659725, gamma=0.024961736585479025;</pre>	total time= 2.2min
[CV] ENDC=8.761939120022886, gamma=0.028187311100850966;	total time= 2.2min
[CV] ENDC=8.761939120022886, gamma=0.028187311100850966;	total time= 2.2min
[CV] ENDC=8.761939120022886, gamma=0.028187311100850966;	total time= 2.2min
[CV] ENDC=4.274866401320841, gamma=0.07036430390020663;	total time= 1.9min
[CV] ENDC=4.274866401320841, gamma=0.07036430390020663;	total time= 1.9min
[CV] ENDC=4.274866401320841, gamma=0.07036430390020663;	total time= 1.9min
[CV] ENDC=10.258605798671827, gamma=0.0013549245402061753;	total time= 4.0min
[CV] ENDC=10.258605798671827, gamma=0.0013549245402061753;	total time= 3.9min
[CV] ENDC=10.258605798671827, gamma=0.0013549245402061753;	total time= 4.0min
[CV] ENDC=9.720765392469097, gamma=0.0010623844616929794;	total time= 4.2min
[CV] ENDC=9.720765392469097, gamma=0.0010623844616929794;	total time= 4.1min
[CV] END C=9.720765392469097, gamma=0.0010623844616929794;	total time= 4.2min
SE: 0.7441 ACC: 0.9637 SP: 0.965496	
6156.775834560394 seconds	

3.4 Experiment 2 – Implementation of the classification and prediction models on Apache Spark

The first step in this experiment is to load the data into Hadoop Distributed File System (HDFS).

(base) hduser@rishab:~/Final\$ hdfs dfs -copyFromLocal ECG.csv /user/hduser

To check whether the dataset has been successfully loaded or not type the following command:

(base) hduse	r@rishab:~/Final\$ hdfs dfs -ls							
WARNING: HADOOP_PREFIX has been replaced by HADOOP_HOME. Using value of HADOOP_PREFIX.								
Found 6 item	IS							
drwxr-xr-x	- hduser supergroup 0 2021-07-26 09:45 .sparkStaging							
- r w- rr	1 hduser supergroup 143741699 2021-08-02 16:01 ECG.csv 🛑							
drwxr-xr-x	- hduser supergroup 0 2021-07-27 18:35 ECG_Data							
- r w-rr	1 hduser supergroup							
drwxr-xr-x	- hduser supergroup 0 2021-07-26 09:18 Final_Data							
-rw-rr	1 hduser supergroup							

3.4.1 Implementation of Deep Neural Network model

a) Import libraries

```
1 from
         _future__ import print_function
 2 from sklearn.metrics import confusion_matrix
 3 import pyspark.sql.functions as F
4 from matplotlib import pyplot
5 from pyspark.mllib.evaluation import MulticlassMetrics
6 from collections import Counter
7 from pyspark.sql.functions import col
8 import pandas as pd
9 from pyspark.sql import SparkSession
10 from pyspark import SparkContext
11 from sklearn.metrics import *
12 import numpy as np
13 import tensorflow as tf
14 from tensorflow.keras.models import Sequential
15 from tensorflow.keras.layers import *
16 from tensorflow.keras.utils import to_categorical
17 from tensorflow.keras import backend as K
18 from tensorflow.keras import callbacks
19 from tensorflow.keras.optimizers import Adam
20 from pyspark.sql.types import FloatType
21 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
22 from pyspark.ml.feature import VectorAssembler
23 from pyspark.ml.tuning import ParamGridBuilder, TrainValidationSplit
24 from tensorflow.keras import optimizers, regularizers
25 from hyperas import optim
26 from hyperopt import Trials, STATUS_OK, tpe
27 from hyperas.distributions import choice, uniform
```

b) Create Spark Context and start Spark Session

```
29 #Create a Spark Context
30 sc = SparkContext()
31
32 #Create a function to read the data and output train dataset and test dataset
33 def data():
34   spark = SparkSession \
35    .builder \
36    .appName("CNN in Spark") \
37    .getOrCreate()
```

c) Import the dataset from hdfs.

38
39 Data = spark.read.format("csv").option("header", "true").load('hdfs:///user/hduser/ECG.csv')
40

d) Change the data type for the columns from object to float. Then rename the 0_y column to 'label'. Select the features and labels and convert them into NumPy arrays.

```
41 Data = Data.select(*(col(c).cast("float").alias(c) for c in Data.columns))
42 Data = Data.withColumnRenamed("0_y", "label")
43 X = np.array(Data.select(*(col(c) for c in Data.columns if c not in {'label'})).collect())
44 Y = np.array(Data.select('label').collect())
45 Y = Y.flatten()
```

e) Split the features and labels into training set and test set in the ratio 80:20.

```
47
      class num = 5
      r seed = 1229
48
49
      data_len = len(X)
50
      np.random.seed(r seed)
51
      idx = list(range(data len))
52
      np.random.shuffle(idx)
53
54
      train_len = int(data_len*0.8)
      test_len = data_len-train_len
55
56
57
      X_train = X[idx][:train_len]
      X_test = X[idx][train_len:]
58
59
      Y_train = Y[idx][:train_len]
60
      Y_test = Y[idx] train_len:
61
62
63
      Y_train = tf.keras.utils.to_categorical(Y_train, num_classes=class_num)
64
      Y_test = tf.keras.utils.to_categorical(Y_test, num_classes=class_num)
65
      return X_train, Y_train, X_test, Y_test
66
```

f) Function to create a Deep neural Network using 4 layers.

```
69 #Function to create a Deep Neural Network model using 4 layers
70 def model(X_train, Y_train, X_test, Y_test):
71
      model = Sequential()
72
      model.add(Dense(512, input_shape=(201,)))
      model.add(Activation('relu'))
73
      model.add(Dropout({{uniform(0, 1)}}))
74
      model.add(Dense({{choice([256, 512, 1024])}}))
model.add(Activation('relu'))
75
76
      model.add(Dropout({{uniform(0, 1)}}))
77
78
      model.add(Dense(5))
79
      model.add(Activation('softmax'))
80
      #Compile the model
      model.compile(loss='categorical crossentropy', metrics=['accuracy'], optimizer='adam')
81
82
      #Train the model
      result = model.fit(X_train, Y_train,
83
84
                 batch_size={{choice([64, 128])}},
85
                 epochs=5,
86
                 verbose=2
87
                 validation_split=0.2)
88
89
      #Find the best accuracy for each epoch
      validation_acc = np.amax(result.history['val_accuracy'])
90
91
      print('Best validation acc of epoch:', validation_acc)
92
93
      #Return the best accuracy and status of the model
94
       return {'loss': -validation_acc, 'status': STATUS_OK, 'model': model}
```

g) Use HyperOpt² to implement keras parallely. This function also helps select the best model (optimal model) during training that can be tested and validated.

```
96 #Using HyperOpt function, the DNN can be executed parallely
97 best_run, best_model = optim.minimize(model=model,
98
99
100
100
101
102 X_train, Y_train, X_test, Y_test = data()
103 print("Evalutation of best performing model:")
104 print(best_model.evaluate(X_test, Y_test))
```

h) Use the best configuration model for evaluation.

```
106 #Select the best configuration of the model and test it on the test data
107 y pred = best model.predict(X test)
108 y_pred = np.argmax(y_pred, axis=1)
109 y true = list(Y test)
110 y true = np.argmax(y true, axis=1)
111
112 #print the confusion matrix
113 conf = confusion_matrix(y_true, y_pred).ravel()
114
115 #Reshape the single array confusion matrix into nxn matrix
116 conf = np.reshape(conf,(5,5))
117
118 print(type(conf))
119 #Extract tn, tp, fn, fp from the confusion matrix
120 fp = conf.sum(axis=0) - np.diag(conf)
121 fn = conf.sum(axis=1) - np.diag(conf)
122 tp = np.diag(conf)
123 tn = conf.sum() - (fp + fn + tp)
124
125 fp = fp.astype(float)
126 fn = fn.astype(float)
127 tp = tp.astype(float)
128 tn = tn.astype(float)
129
130 #Calculate the sensitivity and specificity
131 SE = tp/(tp+fn)
132 SP = tn/(fp+tn)
133 SE = np.average(SE)
134 SP = np.average(SP)
135
136 print(" SE: %.4f | SP: %4f" %(SE, SP))
```

i) Save the file and clost it. Run the python file using **spark-submit.**

(TensorFlow) hduser@rishab:~/Final/Spark\$ spark-submit SparkCNN.py

² https://pythonrepo.com/repo/maxpumperla-hyperas-python-deep-learning

OUTPUT -

Epoch 1/5 548/548 - 19s - loss: 0.4416 - accuracy: 0.8875 - val_loss: 0.3745 - val_accuracy: 0.9167 Epoch 2/5 548/548 - 16s - loss: 0.3051 - accuracy: 0.9256 - val_loss: 0.2539 - val_accuracy: 0.9381 Epoch 3/5 548/548 - 12s - loss: 0.2791 - accuracy: 0.9308 - val_loss: 0.2307 - val_accuracy: 0.9384 Epoch 4/5 548/548 - 17s - loss: 0.2603 - accuracy: 0.9344 - val_loss: 0.2070 - val_accuracy: 0.9419 Epoch 5/5 548/548 - 14s - loss: 0.2551 - accuracy: 0.9351 - val_loss: 0.2027 - val_accuracy: 0.9436 Best validation acc of epoch: 0.9436450600624084 | 5/5 [05:58<00:00, 71.79s/trial, best loss: -0.9436450600624084] 100% Evalutation of best performing model: 685/685 [=============================] - 5s 7ms/step - loss: 0.2196 - accuracy: 0.9394 [0.2196389138698578, 0.9393870234489441]

3.4.2 Implementation of Random Forest model

a) Import the libraries

<class 'numpy.ndarray'> SE: 0.5314 | SP: 0.938900

```
1 from sklearn.metrics import classification_report, confusion_matrix
 2 from pyspark.sql.functions import rand
 3 import sklearn
4 from pyspark.mllib.evaluation import MulticlassMetrics
5 from pyspark.sql.functions import col
6 import pandas as pd
7 from pyspark.sql import SparkSession
8 from pyspark import SparkContext
9 import numpy as np
10 from pyspark.sql.types import FloatType
11 from pyspark.ml.classification import LinearSVC, OneVsRest
12 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
13 from pyspark.ml.feature import VectorAssembler
14 from pyspark.ml.tuning import ParamGridBuilder, TrainValidationSplit
15 from pyspark.ml import Pipeline
16 from pyspark.ml.classification import RandomForestClassifier
17 import pyspark.sql.functions as F
```

b) Create Spark Context and start Spark Session

```
19 #Create a Spark Context
20 sc = SparkContext()
21
22 #Enable a spark session
23 spark = SparkSession \
24   .builder \
25   .appName("RF in Spark") \
26   .getOrCreate()
```

c) Read the data from hdfs. Change the data type from object to float. Select all the feature columns and convert them into a vector. Select features and labels columns.

```
28 #Read thhe data from HDFS
29 Data = spark.read.format("csv").option("header", "true").load('hdfs:///user/hduser/ECG.csv')
30
31 #Rename the last coulm to labels
32 Data = Data.withColumnRenamed("0_y", "label")
33
34 #Convert all the column into floating type
35 Data = Data.select(*(col(c).cast("float").alias(c) for c in Data.columns))
36
37 #Select all the feature columns (201 columns) and
38 #store it a vector named features
39 Features assem = VectorAssembler(inputCols=[c for c in Data.columns if c not in {'label'}],
 outputCol="features")
40
41 #Apply the transformation to the data
42 Data = Features_assem.transform(Data)
```

d) Split the data in inito train and test ration in the ration 70:30. Create the random forest model with 200 trees with a tree depth of 5. Train the model and evaluate the performance of the model on the test data.

```
50 #Split the Data into train and test data in the ratio 70:30
51 Train Data, Test Data = Final Data.randomSplit([0.7,0.3])
52
53 #Create a random forest model
54 rf = RandomForestClassifier(labelCol="label"
                               featuresCol="features",
55
56
                               numTrees = 200,
57
                               maxDepth = 5,
58
                               maxBins = 30)
59
60 #Fit the model on the training data
61 rf_model = rf.fit(Train_Data)
62
63 #Test the trained model on the test dataset
64 predictions = rf_model.transform(Test_Data)
65
66 #Select the features and labels from the predicted data for evaluation
67 y_true = predictions.select(['label']).collect()
68 y_pred = predictions.select(['prediction']).collect()
69
70 #Print the classification report
71 print(classification_report(y_true, y_pred))
72
73 #Create a confusion matrix
74 conf_mat = confusion_matrix(y_true, y_pred)
75 print(conf_mat)
76
77 #Extract tn, tp, fn, fp from confusion matrix
78 fp = conf_mat.sum(axis=0) - np.diag(conf_mat)
79 fn = conf mat.sum(axis=1) - np.diag(conf mat)
80 tp = np.diag(conf_mat)
81 tn = conf_mat.sum() - (fp + fn + tp)
82
83 fp = fp.astype(float)
84 fn = fn.astype(float)
85 tp = tp.astype(float)
86 tn = tn.astype(float)
87
88 #Calculate the specificity and sensitivity
89 SE = tp/(tp+fn)
90 SP = tn/(fp+tn)
91 SE = np.average(SE)
92 SP = np.average(SP)
93
94 print(" SE: %.4f | SP: %4f" %(SE, SP))
```

e) Save the file and close it. Rub the program using the following command

(base) hduser@rishab:~/Final/Spark\$ spark-submit SparkRF.py

OUTPUT -

		precisio	n	recall	f1-score	support
	0.0	0.9	0	1.00	0.95	27222
	1.0	0.0	0	0.00	0.00	875
	2.0	0.9	6	0.55	0.70	2104
	3.0	0.0	0	0.00	0.00	259
	4.0	1.0	0	0.59	0.74	2358
accur	асу				0.91	32818
масго	avg	0.5	7	0.43	0.48	32818
weighted	avg	0.8	8	0.91	0.88	32818
[[27201	0	20	0	1]		
[874	0	1	Θ	0]		
[942	0	1161	0	1]		
[231	0	28	Θ	0]		
[958	0	5	0	1395]]		
SE: 0.42	285	SP: 0.89	2237			

3.4.3 Implementation of Support Vector Machine model

a) Import the required libraries.

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 import sklearn
3 from pyspark.mllib.evaluation import MulticlassMetrics
4 from collections import Counter
5 from pyspark.sql.functions import col
6 import pandas as pd
7 from pyspark.sql import SparkSession
8 from pyspark import SparkContext
9 import numpy as np
10 from pyspark.sql.types import FloatType
11 from pyspark.ml.classification import LinearSVC, OneVsRest
12 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
13 from pyspark.ml.feature import ParamGridBuilder, TrainValidationSplit
```

b) Create Spark Context and start Spark Session

```
16 #Create a Spark Context
17 sc = SparkContext()
18
19 #Enable a spark session
20 spark = SparkSession \
21 .builder \
22 .appName("SVM in Spark") \
23 .getOrCreate()
```

c) Read the data from hdfs. Change the data type from object to float. Select all the feature columns and convert them into a vector. Select features and labels columns.

```
25 #Read thhe data from HDFS
26 Data = spark.read.format("csv").option("header", "true").load('hdfs:///user/hduser/ECG.csv')
27
28 #Convert all the column into floating type
29 Data = Data.select(*(col(c).cast("float").alias(c) for c in Data.columns))
30
31 #Rename the last coulm to labels
32 Data = Data.withColumnRenamed("0_y", "label")
33
34 #Select all the feature columns (201 columns) and
35 #store it a vector named features
36 Features_assem = VectorAssembler(inputCols=[c for c in Data.columns if c not in {'label'}],
  outputCol="features")
37
38 #Apply the transformation to the data
39 Data = Features_assem.transform(Data)
40
41 #Select the features and labels columns for classification
42 Final_Data = Data.select(col("label"), col("features"))
```

d) Split the data into train and test. Create the support vector machine using Linear Support Vector Classifier. Since there are five classes, one vs all classifier is attached to the model to handle multi-classification. The model is trained for 25 iterations and evaluated.

```
44 #Split the Data into train and test data in the ratio 70:30
45 Train Data, Test Data = Final Data.randomSplit([0.8,0.2])
46
47 #Create a linear model with 25 iterations
48 #Fit one vs all classifier-to enable multiclass classification
49 lsvc = LinearSVC(maxIter=25, regParam=0.001)
50 ovr = OneVsRest(classifier=lsvc)
51
52 #Fit the model on the training data
53 model = ovr.fit(Train_Data)
54
55 #Test the model on test data
56 predictions = model.transform(Test_Data)
58 #Select the features and labels from the predicted data for evaluation
59 y_true = predictions.select(['label']).collect()
60 y_pred = predictions.select(['prediction']).collect()
61
62 #Print the classification report
63 print(classification_report(y_true, y_pred))
64
65 #Create a confusion matrix
66 conf_mat = confusion_matrix(y_true, y_pred)
67 print(conf_mat)
68
69 #Extract tn, tp, fn, fp from confusion matrix
70 fp = conf_mat.sum(axis=0) - np.diag(conf_mat)
71 fn = conf_mat.sum(axis=1) - np.diag(conf_mat)
72 tp = np.diag(conf_mat)
73 tn = conf_mat.sum() - (fp + fn + tp)
74
75 fp = fp.astype(float)
76 fn = fn.astype(float)
77 tp = tp.astype(float)
78 tn = tn.astype(float)
80 #Calculate the accuracy specificity and sensitivity
81 \text{ acc} = (tp+tn)/(tp+tn+fp+fn)
82 acc = np.average(acc)
83 SE = tp/(tp+fn)
84 SP = tn/(fp+tn)
85 SE = np.average(SE)
86 SP = np.average(SP)
87
88 print(" SE: %.4f | SP: %4f | acc: %.4f" %(SE, SP, acc))
```

OUTPUT -

		precision		recall	f1-score	support
	0.0	0.83		1.00	0.91	18131
	1.0	0.00		0.00	0.00	577
	2.0	1.00		0.00	0.00	1467
	3.0	0.00		0.00	0.00	159
Ì	4.0	0.00		0.00	0.00	1592
accur	acv				0.83	21926
масго	avg	0.37		0.20	0.18	21926
weighted	avg	0.75		0.83	0.75	21926
[[18131	0	Θ	0	0]		
[577	0	Θ	0	0]		
[1466	O	1	0	0]		
[159	0	0	0	0]		
[1592	0	0	0	0]]		
<class 'r<="" td=""><td>umpy</td><td>.ndarray'></td><td></td><td></td><td></td><td></td></class>	umpy	.ndarray'>				
SE: 0.20	001	SP: 0.800	053	acc:	0.9308	

4 Conclusion

All the above steps were implemented to create Electrocardiogram classification and prediction using the MIT-BIH dataset. All the tools used were mentioned throughout the implementation stages. Two experiments – Local implementation and Apache Spark were implemented to answer the research questions and research objectives.

It is important to use the specified versions of Hadoop and Apache Spark, as errors or glitches might arise if any other version is used.

REFERENCES

[1] <u>https://pywavelets.readthedocs.io/en/latest/</u>

- [2] <u>https://github.com/MIT-LCP/wfdb-python</u>
- [3] https://keras.io/guides/writing_a_training_loop_from_scratch/