

Configuration Manual

MSc Research Project
Data Analytics

Shveta Srivastava
Student ID: 18194851

School of Computing
National College of Ireland

Supervisor: Dr. Rashmi Gupta

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Shveta Srivastava
Student ID:	18194851
Programme:	Data Analytics
Year:	2020-2021
Module:	MSc Research Project
Supervisor:	Dr. Rashmi Gupta
Submission Due Date:	16-08-2021
Project Title:	Configuration Manual
Word Count:	1189
Page Count:	20

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:	16th August 2021

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
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.	<input type="checkbox"/>

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):	

Genetic Algorithm Optimized Deep Learning Model for Parkinson's Disease Severity Detection

Shveta Srivastava
18194851

1 Introduction

This configuration manual contains the information for complete implementation of the research project including hardware and software used, so as to replicate the work done at any time. This step-by-step manual will aid the users to understand the code and its implementation.

2 Hardware Specification

The project was carried out on a Dell G3 15 laptop using the configuration shown in Figure 1 and Figure 2.

Device specifications	
HP Pavilion Laptop 14-ce0xxx	
Device name	LAPTOP-J34R64GG
Processor	Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz 1.80 GHz
Installed RAM	8.00 GB (7.88 GB usable)
Device ID	AB5EA88C-8281-4E73-B739-7FB851606739
Product ID	00325-96500-78779-AAOEM
System type	64-bit operating system, x64-based processor
Pen and touch	No pen or touch input is available for this display

Figure 1: Device Specifications

3 Software Requirement and Environment Setup

The research work was implemented using R and Python. The process of setting up R environment and Python environment are shown below.

OS Name	Microsoft Windows 10 Home
Version	10.0.19043 Build 19043
Other OS Description	Not Available
OS Manufacturer	Microsoft Corporation
System Name	LAPTOP-J34R64GG

Figure 2: System Specifications

3.1 R

R is a software environment that is free for use for building machine learning tools, graphical tools and statistical computing (Matloff; 2011). R (Version 4.1.1) can be installed using Cran R website,¹.

3.2 Anaconda for Python

Anaconda is the data science toolkit mainly for python programming (Van Rossum et al.; 2000). For python, Anaconda was installed from the website for Windows². Once Anaconda is installed successfully, Jupyter Notebook can be launched and used for Python Programming as shown in Figure 3.

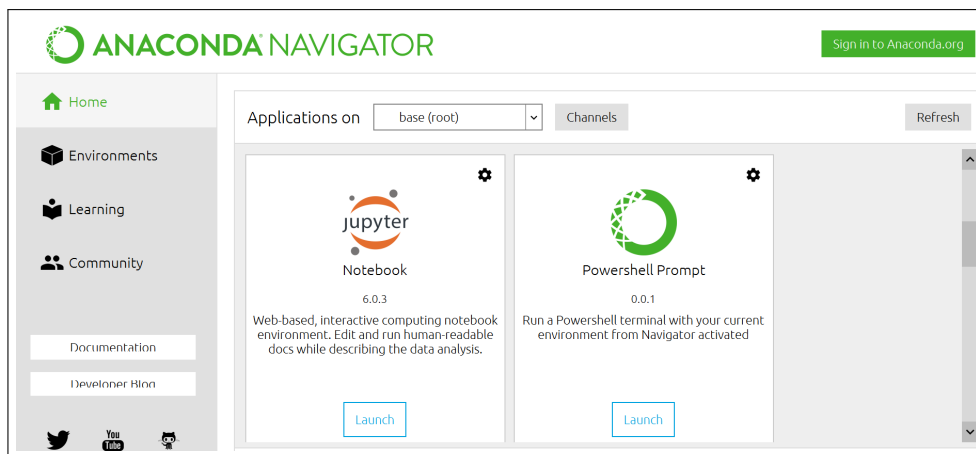


Figure 3: Launch Jupyter Notebook from Anaconda Navigator

3.3 Tensor Flow

Tensorflow environment can be installed by using the CMD.EXE prompt following the instructions given at Anaconda online user guide ³. Abadi et al. (2016) describes several tensorflow applications for machine learning algorithm in detail.

¹<https://cran.r-project.org/bin/windows/base/>

²<https://www.anaconda.com/products/individual>

³<https://docs.anaconda.com/anaconda/user-guide/tasks/tensorflow/>

3.4 Google Colaboratory

Google Colaboratory, also called 'Google Colab' is a free cloud service which offers Jupyter Notebook for running python by selecting GPU and TPU setup within Colab (Carneiro et al.; 2018). Dataset can be retrieved from Google Drive by mounting it in Colab Notebook.

4 Data Collection and Merging

Parkinson's Disease patient's motor and non-motor assessment data has been collected from PPMI (Parkinson's Progressive Markers Initiative) on approval of request from the PPMI website ⁴.

The merging and data transformation has been done using R programming. Downloaded all the motor and non-motor assessments data and uploaded on R environment.

4.1 Read and Understand Datasets

The following datasets have been downloaded from the website. The structure of the datasets have been visualized and understood the important features from each of the dataset, using code dictionary given at PPMI.

```
1
2 #read the datasets
3
4 UPDRS1 <- read.csv("MDS_UPDRS_Part_I.csv", header = T)
5 head(UPDRS1)
6 str(UPDRS1)
7
8 UPDRSII <- read.csv("MDS_UPDRS_Part_II.csv", header = T)
9 str(UPDRSII)
10
11 UPDRSIII <- read.csv("MDS_UPDRS_Part_III.csv", header = T)
12 str(UPDRSIII)
13 head(UPDRSIII)
14
15 SCOPA <- read.csv("SCOPA-AUT.csv", header = T)
16 str(SCOPA)
17
18 MOCA <- read.csv("Montreal_Cognitive_Assessment__MoCA.csv", header = T)
19 str(MOCA)
20
21 BOSTON <- read.csv("Modified_Boston_Naming_Test.csv", header = T)
22 str(BOSTON)
23
24 STAI <- read.csv("State-Trait_Anxiety_Inventory.csv", header= T)
25 str(STAI)
26
27 QUIP <- read.csv("QUIP_Impulsive Compulsive Disorder.csv", header= T)
28 str(QUIP)
29
30 GDS <- read.csv("Geriatric_Depression_Scale.csv", header= T)
31 str(GDS)
32
```

⁴<https://www.ppmi-info.org/>

```

33 HVLTL <- read.csv("Hopkins_Verbal_Learning_Test.csv", header = T)
34 str(HVLTL)
35
36 LNS <- read.csv("Letter-Number_Sequencing.csv", header = T)
37 str(LNS)
38
39 SFT <- read.csv("Semantic_Fluency.csv", header = T)
40 str(SFT)
41
42 SDM <- read.csv("Symbol_Digit_Modalities.csv", header = T)
43 str(SDM)
44
45 CLOCK <- read.csv("Clock_Drawing.csv", header = T)
46 str(CLOCK)
47
48 COGNITIVE <- read.csv("Cognitive_Categorization.csv", header = T)
49 str(COGNITIVE)
50
51 LFT <- read.csv("Lexical_Fluency.csv", header = T)
52 str(LFT)
53
54 TRAIL <- read.csv("Trail_Making_A_and_B.csv", header = T)
55 str(TRAIL)
56
57 EPWORTH <- read.csv("Epworth_Sleepiness_Scale.csv", header = T)
58 str(EPWORTH)
59
60 REM <- read.csv("REM_Sleep_Disorder_Questionnaire.csv", header = T)
61 str(REM)

```

Listing 1: Reading and Understanding Datasets

The list of several motor and non-motor assessment datasets used in this study can be seen in R environment as shown in Figure 4 with the number of variables and the number of patient records in the PPMI database. It should be noted that LFT, BOSTON, CLOCK and TRAIL have been excluded from study due to very few number of patient records which will not add any value to the analysis and might hinder the model building process.

The target variable 'HY' severity scale scores are contained in 'UPDRS III' dataset.

4.2 Select important variables

Next, the important variables from all the datasets have been selected using 'Dplyr' function.

```

1
2 library(dplyr)
3 N_UPDRS1 <- UPDRS1 %>%
4   select(PATNO, EVENT_ID, NP1COG, NP1HALL, NP1DPRS, NP1ANXS, NP1APAT,
5     NP1DDS)
6 head(N_UPDRS1)
7
8 N_UPDRSII <- UPDRSII %>%
9   select(PATNO, EVENT_ID, NP2SPCH, NP2SALV, NP2SWAL, NP2EAT, NP2DRES,
10     NP2HYGN, NP2HWRT, NP2HOBB,
11     NP2TURN, NP2TRMR, NP2RISE, NP2WALK, NP2FREZ)

```

▶ BOSTON	438 obs. of 13 variables
▶ CLOCK	332 obs. of 18 variables
▶ COGNITIVE	7030 obs. of 18 variables
▶ EPWORTH	8227 obs. of 19 variables
▶ GDS	8546 obs. of 25 variables
▶ HVL T	6999 obs. of 23 variables
▶ LFT	443 obs. of 13 variables
▶ LNS	6969 obs. of 34 variables
▶ MOCA	7860 obs. of 38 variables
▶ QUIP	8229 obs. of 24 variables
▶ REM	8836 obs. of 32 variables
▶ SCOPA	8229 obs. of 46 variables
▶ SDM	6986 obs. of 15 variables
▶ SFT	6995 obs. of 16 variables
▶ STAI	8535 obs. of 50 variables
▶ TRAIL	436 obs. of 16 variables
▶ UPDRS1	13904 obs. of 17 variables
▶ UPDRSII	13908 obs. of 24 variables
▶ UPDRSIII	15843 obs. of 54 variables

Figure 4: List of Datasets from PPMI

```

11 N_UPDRSIII <- UPDRSIII %>%
12   select(PATNO, EVENT_ID, NP3SPCH, NP3FACXP, NP3RIGN, NP3RIGRU,
13     NP3RIGLU, PN3RIGRL, NP3RIGLL, NP3FTAPR,
14     NP3FTAPL, NP3HMOVR, NP3HMOVL, NP3PRSPR, NP3PRSPL, NP3TTAPR,
15     NP3TTAPL, NP3LGAGL,
16     NP3LGAGR, NP3RISNG,
17     NP3GAIT, NP3FRZGT, NP3PSTBL, NP3POSTR, NP3BRADY, NP3PTRMR,
18     NP3PTRML, NP3KTRMR, NP3KTRML, NP3RTARU,
19     NP3RTALU, NP3RTARL, NP3RTALL, NP3RTALJ, NP3RTCON, DYSKPRES, NHY
20   )
21
22 N_SCOPA <- SCOPA %>%
23   select(PATNO, EVENT_ID, SCAU1, SCAU2, SCAU3, SCAU4, SCAU5, SCAU6,
24     SCAU7, SCAU8,
25     SCAU9, SCAU10, SCAU11, SCAU12, SCAU13, SCAU14, SCAU15, SCAU16,
26     SCAU17,
27     SCAU18, SCAU19, SCAU20, SCAU21, SCAU24, SCAU26A, SCAU26B,
28     SCAU26C, SCAU26D)
29 str(N_SCOPA)
30
31 N_MOCA <- MOCA %>%
32   select(PATNO, EVENT_ID, MCAALTTM, MCACUBE, MCACLCKC, MCACLCKN,
33     MCACLCKH, MCALION, MCA RHINO,
34     MCACAMEL, MCAFDS, MCABDS, MCAVIGIL, MCASER7, MCASNTNC, MCAVFNUM, MCAVF,
35     MCAABSTR, MCAREC1,

```

```

27 MCAREC2, MCAREC3, MCAREC4, MCAREC5, MCADATE, MCAMONTH, MCAYR, MCADAY,
    MCAPLACE, MCACITY, MCATOT)
28
29 N_REM <- REM %>%
30   select(PATNO, EVENT_ID, DRMVIVID, DRMAGRAC, DRMNOCTB, SLPLMBMV,
    SLPINJUR, DRMVERBL, DRMFIGHT,
31     DRMUMV, DRMOBJFL, MVAWAKEN, DRMREMEM, SLPDSTRB, STROKE, HETRA,
    RLS, NARCLPSY, DEPRS,
32     EPILEPSY, BRNINFM)
33
34 head(N_REM)
35
36 N_GDS <- GDS %>%
37   select(PATNO, EVENT_ID, GDSSATIS, GDSDROPD, GDSEEMPTY, GDSBORED,
    GDSGSPR, GDSAFRAD, GDSHAPPY,
38   GDSHPLS, GDSHOME, GDSMEMRY, GDSALIVE, GDSWRTLS, GDSENRGY, GDSHOPLS,
    GDSBETER)
39
40 N_STAI <- STAI %>%
41   select(PATNO, EVENT_ID, STAIAD1, STAIAD2, STAIAD3, STAIAD4, STAIAD5,
    STAIAD6, STAIAD7, STAIAD8,
42     STAIAD9, STAIAD10, STAIAD11, STAIAD12, STAIAD13, STAIAD14,
    STAIAD15, STAIAD16, STAIAD17,
43     STAIAD18, STAIAD19, STAIAD20, STAIAD21, STAIAD22, STAIAD23,
    STAIAD24, STAIAD25, STAIAD26,
44     STAIAD27, STAIAD28, STAIAD29, STAIAD30, STAIAD31, STAIAD32,
    STAIAD33, STAIAD34, STAIAD35,
45     STAIAD36, STAIAD37, STAIAD38, STAIAD39, STAIAD40)
46
47 N_QUIP <- QUIP %>%
48   select(PATNO, EVENT_ID, TMGAMBLE, CNTRLGMB, TMSEX, CNTRLSEX, TMBUY,
    CNTRLBUY, TMEAT, CNTRLEAT,
49     TMTORACT, TMTMTACT, TMTRWD)
50
51 N_EPWORTH <- EPWORTH %>%
52   select(PATNO, EVENT_ID, ESS1, ESS2, ESS3, ESS4, ESS5, ESS6, ESS7,
    ESS8)
53
54 N_SFT <- SFT %>%
55   select(PATNO, EVENT_ID, VLTANIM, VLTVEG, VLTFRUIT)
56 str(N_SFT)
57 library(Amelia)
58 missmap(N_SFT)
59
60 N_HVLT <- HVLT %>%
61   select(PATNO, EVENT_ID, HVLTRT1, HVLTRT2, HVLTRT3, HVLTRDLY, HVLTRC,
    HVLTFPRL, HVLTFPUN)
62 missmap(N_HVLT)
63
64 N_LNS <- LNS %>%
65   select(PATNO, EVENT_ID, LNS_TOTRAW)
66
67 N_SDM <- SDM %>%
68   select(PATNO, EVENT_ID, SDMTOTAL)

```

Listing 2: Selecting Variables from each dataset

4.3 Merge datasets

The datasets have been merged together by unique patient number 'PATNO' and their visit ID 'EVENT_ID'. The Datasets with very few records have been avoided. The final merged dataset has been saved in the google drive for further retrieval in Colab notebook for further analysis.

```
1
2 #merge
3 D1 <- merge(N_UPDRS1,N_UPDRSIII, by= c("PATNO", "EVENT_ID"))
4 str(D1)
5
6 D2 <- merge(D1,N_UPDRSII, by= c("PATNO", "EVENT_ID"))
7 str(D2)
8
9 D3 <- merge(D2,N_REM, by= c("PATNO", "EVENT_ID"))
10 str(D3)
11
12 D4 <- merge(D3,N_GDS, by= c("PATNO", "EVENT_ID"))
13 str(D4)
14
15 D5 <- merge(D4,N_STAI, by= c("PATNO", "EVENT_ID"))
16 str(D5)
17
18 D7 <- merge(D5,N_QUIP, by= c("PATNO", "EVENT_ID"))
19
20 D8 <- merge(D7,N_EPWORTH, by= c("PATNO", "EVENT_ID"))
21 str(D8)
22
23 D9 <- merge(D8,N_SCOPA, by= c("PATNO", "EVENT_ID"))
24
25
26 #exclude MOCA, LNS, HVLIT, SDM, SFT
27 #as the number of rows reduced drastically
28
29 #D9 <- merge(D8,N_MOCA, by= c("PATNO", "EVENT_ID"))
30 #str(D9)
31
32 #D10 <- merge(D9,N_LNS, by= c("PATNO", "EVENT_ID"))
```

Listing 3: Merging datasets

4.4 Coding the Events

The table of the number of total visits from participants can be visualized as shown in figure 5. 'Event_ID' of 'PW', 'ST', 'U01', 'V01', 'V03', 'V05' have been excluded from study due to insufficient data. The rest are coded numerical values for analysis.

BL	PW	ST	U01	V01	V02	V03	V04	V05	V06	V08	V10	V12	V13	V14	V15	V16
1558	10	201	2	12	1029	2	1535	2	1399	1135	985	818	489	524	262	42

Figure 5: List of Datasets from PPMI

```
1 Final_Data <- D9
2 str(Final_Data)
3
```

```

4 #FIND DEPENDENT VARIABLE 'NHY'
5
6 table(Final_Data$NHY)
7 table(Final_Data$EVENT_ID)
8
9 # select visits that are of interest for modeling
10
11 #encoding of visits from baseline till visit 14
12 Final_Data$EVENT_ID <- ifelse(Final_Data$EVENT_ID == 'BL',1,
13                               ifelse(Final_Data$EVENT_ID == 'V02', 2,
14                                       ifelse(Final_Data$EVENT_ID == 'V04
15                                               ',3,
16                                               ifelse(Final_Data$EVENT_ID
17                                                       == 'V06',4,
18                                                       ifelse(Final_Data$
19                                                           EVENT_ID == 'V08',5,
20                                                           ifelse(Final_
21                                                               Data$EVENT_ID == 'V10',6,
22                                                               ifelse
23                                                                   (Final_Data$EVENT_ID == 'V12',7,
24                                                                       ifelse(Final_Data$EVENT_ID == 'V13',8, 9)
25                                                                                       ))))))))
26
27 write.csv(Final_Data, file = "Final_Data.csv")

```

Listing 4: Coding Visits of Interest

5 EDA and Pre-processing

5.1 Load the data in Google Colab

Rest of the analysis is done using Python on Jupyter Notebook in Google Colab. The data is retrieved from Google Drive which has to be mounted on Colab using authentication code. Then the libraries are imported followed by uploading the 'Final Data' from the data path.

```

1
2 from google.colab import drive
3 drive.mount('/content/drive/')
4
5 import tensorflow as tf
6 import pandas as pd
7 import numpy as np
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10 from sklearn.preprocessing import MinMaxScaler
11 from sklearn.model_selection import train_test_split, cross_val_score
12
13 data_path = '/content/drive/MyDrive/Final_Data.csv'
14 data = pd.read_csv(data_path)

```

Listing 5: Read the Data

5.2 Exploratory Data Analysis

The dataset is visualized initially using the following steps.

```
1 data.head()
2 data.shape # Dataset has 9472 rows and 175 features
3 data.dtypes # all the variables have either integer or float values
4 data.describe() # describes the count of variables, mean, std.
5 deviation, range of values with min, max values.
6 sns.catplot(x='NHY',kind='count',data=data) # The severity scale based
7 on 'HY' index is from '1' to '5' indicating various stages of
8 Parkinson's disease.
9 sns.catplot(x='EVENT_ID',kind='count',data=data) # Number of patients
10 participating by visits. number '9' is the final visit. It can be
11 observed that there is a decline in visits.
12 print(data['PATNO'].value_counts())
13 n = len(pd.unique(data['PATNO']))
14 print("Number of Patients participating in assessments :",
15       n) #Number of Patients participating in assessments : 1478
```

Listing 6: Exploratory Data Analysis

The following figure 6 and figure 7 shows the plot of the target variable 'NHY' and 'EVENTS' respectively.

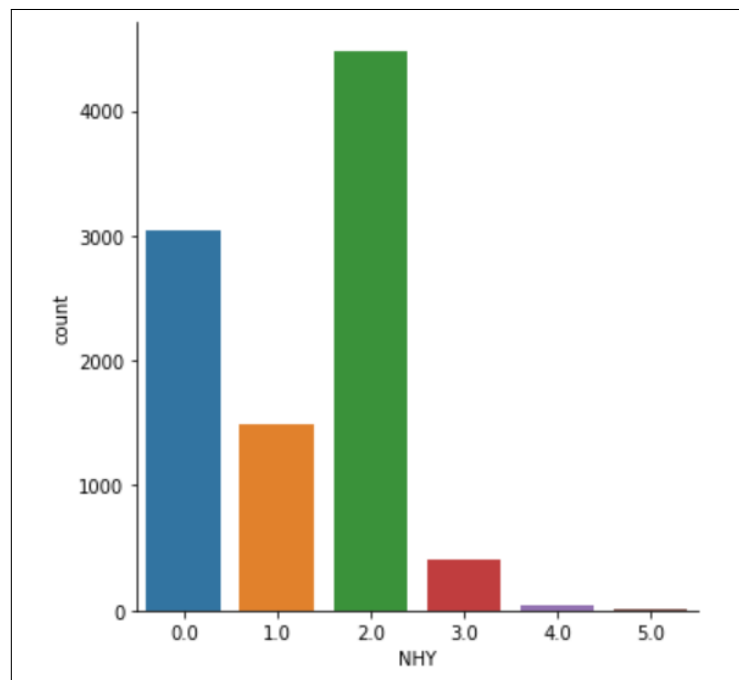


Figure 6: Different stages of Disease Severity

5.3 Handling Missing Records

Missing variables can be visualized using heatmap. Missing records has been replaced with previous assessment score for that patient's visits.

```
1 # plot the missmap to visualize the columns with missing records
2 sns.heatmap(data.isnull(), cbar=False, yticklabels=False, cmap='viridis
3             ')
```

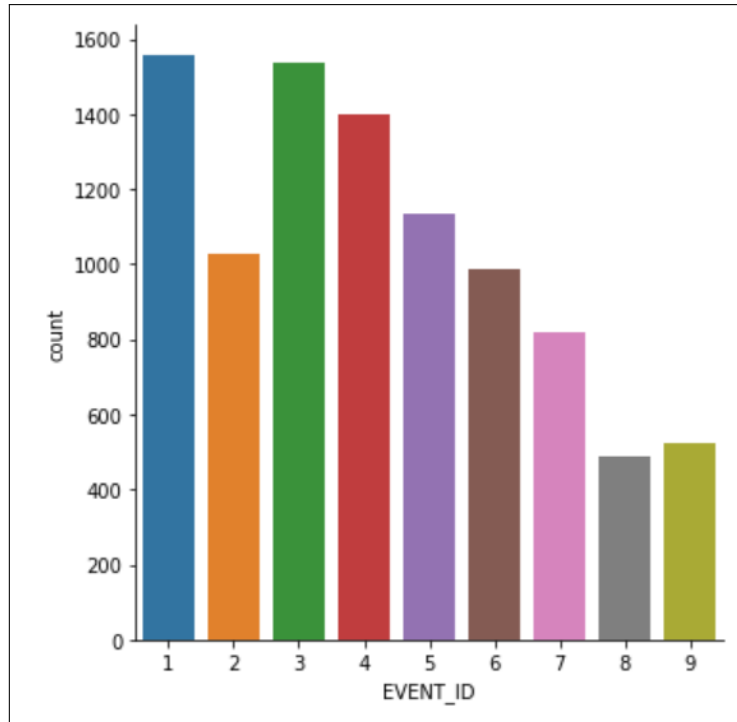


Figure 7: Different visits for PD Assessments

```

3 data= data.drop(['SCAU20','SCAU21','SCAU24'],axis=1) # These column has
  mostly missing values, so its better to remove them
4
5 data.isna().sum().sum() # 827 missing records in dataset
6 data['NHY'].isna().sum()
7
8 # fill NA with interpolate function after sorting each patient with
  their visits.
9 data = data.set_index(['PATNO']).sort_index(level='PATNO') #set index
  and sort data
10 data = data.interpolate(method='linear', axis=1)
11
12 # replace the late stages into moderate stages due to insufficient data
  , and as we want to predict if the patients are 'severe' and need
  medical attention or 'not-severe'
13 data['NHY'] = data['NHY'].replace([5.0, 4.0], 3.0)
14 PD_data = data
15 print(PD_data) # data is now ready for analysis

```

Listing 7: Handling Missing Records

The heat map for missing variables can be seen in figure 8.

6 Modelling

The data is split into features and labels where 'NHY' is label and rest all are features. The features are further standardized using min-max normalisation. The data is then split into train and test in the ratio of 80-20. The steps are shown below.

```

1 #The dataset is split into features(the independent variables) and
  labels(the dependent or target variable):

```

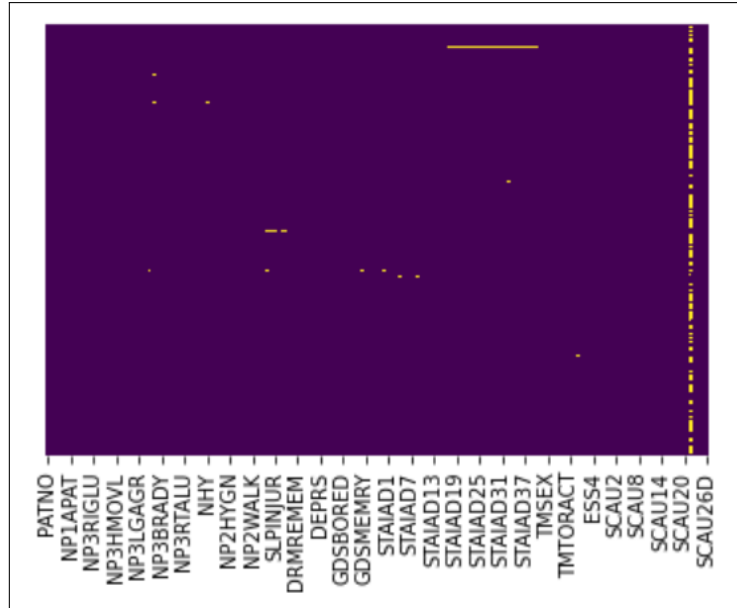


Figure 8: Heat Map for Missing Variables

```

2 features=PD_data.drop(['NHY'],axis=1)
3 labels=PD_data['NHY']
4 #normalize the data using the minmax scaler to bring the feature
  variables within the range -1 to 1:
5 scaler=MinMaxScaler((-1,1))
6 X=scaler.fit_transform(features)
7 y=labels
8
9 #split for training data is 80% and testing is 20%
10 X_train,X_test,y_train,y_test=train_test_split(X, y, test_size=0.2)
11 print('X_train num',X_train.shape[0])
12 print('X_test num',X_test.shape[0])

```

Listing 8: Standardization and Train Test split

6.1 LSTM

Long Short Term Memory Model is built for predicting Parkinson's Disease Severity. LSTM model deals well with the complexities of longitudinal data of PD patients, learning from long-term dependencies. The code is given below. Keras has been utilized to import models and layers.

```

1 from tensorflow import keras
2 from tensorflow.keras import layers
3 import matplotlib.pyplot as plt
4 from keras.models import Sequential
5 from keras.layers import Dense, Dropout, LSTM
6 from sklearn.model_selection import cross_val_score
7 from sklearn.model_selection import KFold
8 from keras.wrappers.scikit_learn import KerasClassifier
9 from keras.callbacks import EarlyStopping
10 from keras.callbacks import ModelCheckpoint
11 from sklearn import metrics
12 from sklearn.metrics import accuracy_score

```

```

13 from sklearn.metrics import confusion_matrix, classification_report,
    roc_auc_score, roc_curve
14 from keras.layers.core import Activation

```

Listing 9: import required libraries

The standardized data is reshaped into Time-series with time-steps equal to number of average visits for each patient, i.e. 5.

```

1 # reshape input to be 3D [samples, timesteps, features]
2 timesteps = 5
3 X_train = X_train.reshape((X_train.shape[0], timesteps, int(X_train.
    shape[1]/timesteps)))
4 X_test = X_test.reshape((X_test.shape[0], timesteps, int(X_test.shape
    [1]/timesteps)))
5 print(X_train.shape, y_train.shape, X_test.shape, y_test.shape)
6 #(7577, 5, 34) (7577,) (1895, 5, 34) (1895,)

```

Listing 10: Reshape according to time-steps

The LSTM model is built with 10% dropouts to avoid over-fitting after each dense layer. Softmax activation is used for outer layer to make sure the probability outcomes totals to 1. Optimizer used in compiling the LSTM architecture is 'adam' optimizer and metrics is 'sparse categorical accuracy'.

```

1 model = Sequential()
2 model.add(LSTM(44, input_shape=(X_train.shape[1], X_train.shape[2])))
3 model.add(Dense(20))
4 model.add(Activation('relu'))
5 model.add(Dropout(0.1))
6 model.add(Dense(10))
7 model.add(Activation('relu'))
8 model.add(Dropout(0.1))
9 model.add(Dense(4))
10 model.add(Activation('softmax'))
11
12 model.compile(loss='sparse_categorical_crossentropy', optimizer='adam',
    metrics= ['sparse_categorical_accuracy'])
13
14 model.summary()

```

Listing 11: LSTM Model Architecture

Model is implemented with early stopping for model loss and model checkpoint for validation accuracy.

```

1 early_stopping = keras.callbacks.EarlyStopping(monitor= 'val_loss',
    mode= 'min', patience= 10)
2
3 model_checkpoint = ModelCheckpoint ('model.h5', monitor= '
    val_sparse_categorical_accuracy', mode = 'max', verbose= 1,
    save_best_only= True)
4
5 history = model.fit(X_train, y_train, validation_data=(X_test, y_test),
    batch_size=32, epochs=100, shuffle= False,
6                     callbacks=[early_stopping, model_checkpoint],
    verbose=2)

```

Listing 12: Fit the Model

Plot the loss and accuracy for comparing validation sample results.

```

1 # plot the loss
2 history_df = pd.DataFrame(history.history)
3 history_df.loc[:, ['loss', 'val_loss']].plot(title="LSTM Model Loss")
4 plt.ylabel('loss')
5 plt.xlabel('epoch')
6
7 #plot model accuracy
8 plt.plot(history.history['sparse_categorical_accuracy'])
9 plt.plot(history.history['val_sparse_categorical_accuracy'])
10 plt.title('LSTM model accuracy')
11 plt.ylabel('accuracy')
12 plt.xlabel('epoch')
13 plt.legend(['train', 'validation'], loc='upper left')
14 plt.show()
15
16 # Make Predictions
17 y_score = model.predict(X_test)
18 y_pred = y_score
19
20 y_pred = (y_pred>0.5)
21 y_pred = y_pred.astype(int)
22 y_pred[0] #array([0, 1, 0, 0])

```

Listing 13: Plots for Model Loss and Model Accuracy

Accuracy for LSTM model is 88%. Model doesn't seem to overfit as the curves of train and validation data are not spaced apart. The LSTM model plots for loss function and model accuracy is shown in figure 9 and figure 10 respectively.

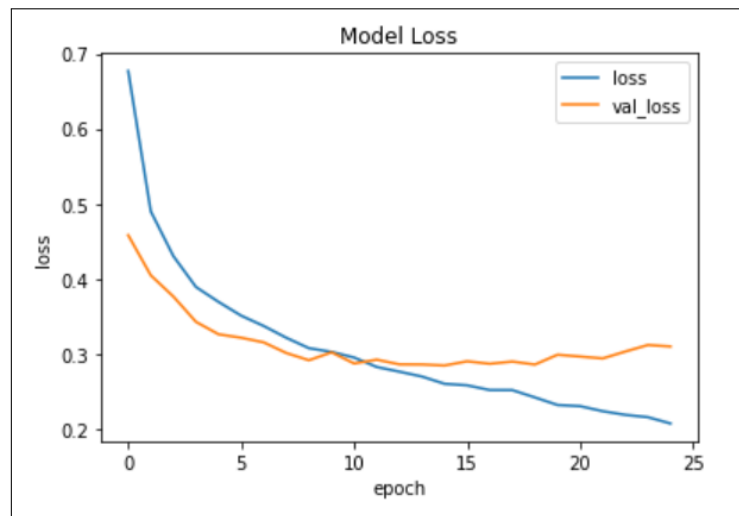


Figure 9: LSTM Model Loss Plot

6.2 DEEP NEURAL NETWORK- Multi-Layer Perceptron

The feed forward network has been built using similar steps as before for binary classification.

```

1 #Patients with 0.0 and 1.0 score on HY scale are considered 'unilateral
   normal' patients,
2 #whereas patients with score of 2 or more needs medical attention and
   are considered 'cognitively impaired' in this project.

```

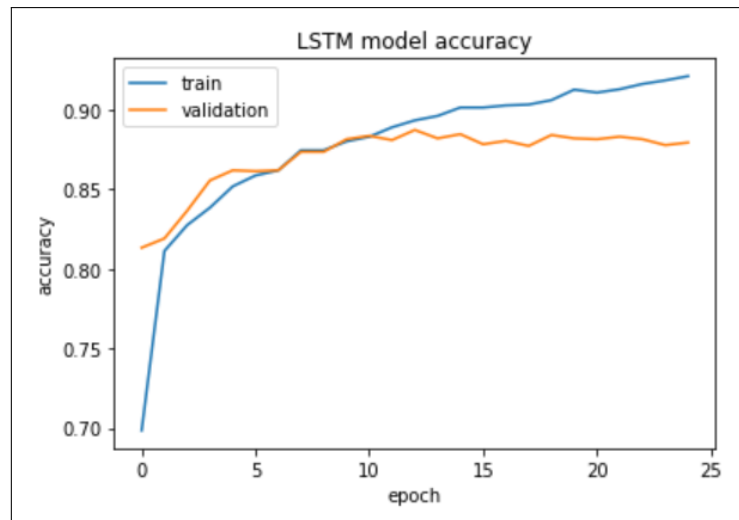


Figure 10: LSTM Model Accuracy Plot

```

3 PD_data['NHY'] = PD_data['NHY'].replace([0.0],1.0)
4 PD_data['NHY'] = PD_data['NHY'].replace([3.0], 2.0)
5
6 #split for training data is 80% and testing is 20%
7 X_train,X_test,y_train,y_test=train_test_split(X, y, test_size=0.2)
8 print(X_train.shape)
9 input_shape = [X_train.shape[1]] #(7577, 170)
10
11 #Model architecture
12 mlp_model = Sequential()
13 mlp_model.add(layers.Dense(170, activation='relu', input_dim=
    input_shape))
14 mlp_model.add(Dropout(0.1))
15 mlp_model.add(layers.Dense(70, activation='relu'))
16 mlp_model.add(Dropout(0.1))
17 mlp_model.add(layers.Dense(16, activation='relu'))
18 mlp_model.add(Dropout(0.1))
19 mlp_model.add(layers.Dense(1, activation='sigmoid'))
20
21 # compile model with optimizing algorithm, loss function and accuracy
    function
22 mlp_model.compile(optimizer="adam", loss="binary_crossentropy", metrics
    =["binary_accuracy"])
23
24 early_stopping = keras.callbacks.EarlyStopping(monitor= 'val_loss',
    mode= 'min',  patience= 10)
25
26 history = mlp_model.fit(X_train, y_train, validation_data=(X_test,
    y_test),batch_size=32, epochs=100, shuffle= False,
    callbacks=[early_stopping],verbose=2)
27
28 #Model Plots
29 history_df = pd.DataFrame(history.history)
30 history_df.loc[:, ['loss', 'val_loss']].plot(title="Cross-entropy")
31 history_df.loc[:, ['binary_accuracy', 'val_binary_accuracy']].plot(
    title="Accuracy")
32 #Make predictions
33 mlp_preds = mlp_model.predict(X_test)
34 #Evaluate

```



```

35 accuracy = mlp_model.evaluate (X_test, y_test)[1]
36 print ('Accuracy:', accuracy)
37 #60/60 [=====] - 0s 1ms/step - loss:
    -2475099947008.0000 - binary_accuracy: 0.4739
38 #Accuracy: 0.4738786220550537
39 print('Mean Absolute Error:', metrics.mean_absolute_error(y_test,
    mlp_preds))
40 print('Mean Squared Error:', metrics.mean_squared_error(y_test,
    mlp_preds))
41 print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(
    y_test, mlp_preds)))
42 #Mean Absolute Error: 0.5261213720316623
43 #Mean Squared Error: 0.5261213720316623
44 #Root Mean Squared Error: 0.72534224475875

```

Listing 14: Multi Layer Perceptron Model

7 Hybrid Model- GA-LSTM

Next, built LSTM model within genetic algorithm heuristic search engine to find the best hyperparameters of the model.

```

1 # import the required libraries. The model was built using bitstring
  and Deap framework.
2 !pip install bitstring
3 !pip install deap
4 from keras.models import Model
5 from deap import base, creator, tools, algorithms
6 from scipy.stats import bernoulli
7 from bitstring import BitArray
8 np.random.seed(1120)

```

Listing 15: install dependencies and libraries

```

Collecting bitstring
  Downloading bitstring-3.1.9-py3-none-any.whl (38 kB)
Installing collected packages: bitstring
Successfully installed bitstring-3.1.9

```

Figure 11: BitString Installation

```

Collecting deap
  Downloading deap-1.3.1-cp37-cp37m-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_12_x86_64.manylinux2010_x86_64.whl (160 kB)
    | 160 kB 12.1 MB/s
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from deap) (1.19.5)
Installing collected packages: deap
Successfully installed deap-1.3.1

```

Figure 12: DEAP Installation

The successful installation message will appear as shown in figure 11 and 12.

```

1 # reshape data into 1-dimensional array
2 data = np.reshape(np.array(PD_data['NHY']), (len(PD_data['NHY']), 1))
3 data.shape
4

```

```

5 # Use first 8525 points as training/validation and rest of the points
  as test set.
6 train_data = data[0:8525]
7 test_data = data[8525:]
8
9 # Prepare the dataset with X,Y pair where for a given time t, X is the
  previous values and Y is the future values
10 def prep_dataset(data, window_size):
11     X, Y = np.empty((0,window_size)), np.empty((0))
12     for i in range(len(data)-window_size-1):
13         X = np.vstack([X,data[i:(i + window_size),0]])
14         Y = np.append([Y,data[i + window_size,0]])
15     X = np.reshape(X,(len(X),window_size,1))
16     Y = np.reshape(Y,(len(Y),1))
17     return X, Y

```

Listing 16: Prepare dataset for genetic algorithm search

Next, Define a function to get window size and number of units by decoding GA solution. Then,prepare the dataset using GA searched window size, further divide the dataset into train and validation set. Finally, train the LSTM model, calculate fitness score of the current solution by Genetic Algorithm.

```

1
2 def train_evaluate(ga_individual_solution):
3     # Decode GA solution to integer for window_size and num_units
4     window_size_bits = BitArray(ga_individual_solution[0:6])
5     num_units_bits = BitArray(ga_individual_solution[6:])
6     window_size = window_size_bits.uint
7     num_units = num_units_bits.uint
8     print('\nWindow Size: ', window_size, ', Num of Units: ', num_units
9     )
10    # Return fitness score of 100 if window_size or num_unit is
  zero
11    if window_size == 0 or num_units == 0:
12        return 100,
13
14    # Segment the train_data based on new window_size; split into train
  and test (80/20)
15    X,Y = prep_dataset(train_data,window_size)
16    X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size
  = 0.20, random_state = 1120)
17
18    # Train LSTM model and predict on validation set
19    inputs = Input(shape=(window_size,1))
20    lstm = LSTM(num_units, input_shape=(window_size,1))(inputs)
21    predictions = Dense(1, activation='linear')(lstm)
22    gamodel = Model(inputs=inputs, outputs=predictions)
23    gamodel.compile(optimizer='adam',loss='binary_crossentropy')
24    gamodel.fit(X_train, y_train, epochs=10, batch_size=10,shuffle=True
  )
25    y_pred = gamodel.predict(X_test)
26    y_pred = y_pred.astype(int)
27
28    # Calculate the RMSE score as fitness score for GA
29    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
30    print('Validation RMSE: ', rmse,'\n')

```

```
31 return rmse,
```

Listing 17: Define Functions for Generating Genetic Algorithm Search

Using DEAP package, defined parameters for implementing GA. Here, the solution is represented by a binary representation '0's and '1's. Bernoulli distribution will randomly initialize the generation. Crossover and mutation are randomly selected by roulette wheel. Finally, the GA parameter values are arbitrarily initialized.

```
1 population_size = 4
2 num_generations = 4
3 gene_length = 10
4
5 # to minimize the RMSE score, weights are defined as -1.0.
6 creator.create('FitnessMax', base.Fitness, weights = (-1.0,))
7 creator.create('Individual', list, fitness = creator.FitnessMax)
8
9 toolbox = base.Toolbox()
10 toolbox.register('binary', bernoulli.rvs, 0.5)
11 toolbox.register('individual', tools.initRepeat, creator.Individual,
12     toolbox.binary,
13     n = gene_length)
14 toolbox.register('population', tools.initRepeat, list, toolbox.individual)
15
16 toolbox.register('mate', tools.cxOrdered)
17 toolbox.register('mutate', tools.mutShuffleIndexes, indpb = 0.6)
18 toolbox.register('select', tools.selRoulette)
19 toolbox.register('evaluate', train_evaluate)
20
21 population = toolbox.population(n = population_size)
22 r = algorithms.eaSimple(population, toolbox, cxpb = 0.4, mutpb = 0.1,
23     ngen = num_generations, verbose = False)
```

Listing 18: Define Genetic Algorithm parameters

```
1 # Print top N solutions - (1st only)
2 best_individuals = tools.selBest(population,k = 1)
3 best_window_size = None
4 best_num_units = None
5
6 for bi in best_individuals:
7     window_size_bits = BitArray(bi[0:6])
8     num_units_bits = BitArray(bi[6:])
9     best_window_size = window_size_bits.uint
10    best_num_units = num_units_bits.uint
11    print('\nTop Solution for Window Size: ', best_window_size, ', Num
12    of Units: ', best_num_units)
13
14 ##### Result for top solution
15 #Top Solution for Window Size: 12 , Num of Units: 8
```

Listing 19: Print the top solutions

```
1 # Train the model using best configuration on complete training set
2 #and make predictions on the test set
3 X_train,y_train = prep_dataset(train_data,best_window_size)
4 X_test, y_test = prep_dataset(test_data,best_window_size)
5
```

```

6 inputs = Input(shape=(best_window_size,1))
7 GA_LSTM = LSTM(best_num_units, input_shape=(best_window_size,1))(inputs
8 )
9 predictions = Dense(1, activation='linear')(GA_LSTM)
10 gamodel = Model(inputs = inputs, outputs = predictions)
11 gamodel.compile(optimizer='adam',loss='binary_crossentropy', metrics=
12 ['mean_squared_error'])
13 gamodel.fit(X_train, y_train, epochs=50, batch_size=32,shuffle=True)
14 y_pred = gamodel.predict(X_test)
15 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
16 print('Test RMSE: ', rmse)
17 #Epoch 50/50
18 #266/266 [=====] - 1s 4ms/step - loss: 0.1349
19 - mean_squared_error: 0.1349
20 #Test RMSE: 0.32939433126774104
21 #y_pred = gamodel.predict(X_test)
22 #print ('Accuracy:', metrics.accuracy_score(y_test, y_pred)) ##
23 ##### Not working on GA-LSTM model as the classification metrics
24 can't handle a mix of binary and continuous targets.

```

Listing 20: Train the LSTM model using the top solutions

Hence, Test Root Mean Square Error (RMSE) was significantly reduced from 0.72 of baseline approach to 0.33 in GA-LSTM model

8 Multiple Baseline Machine Learning Algorithms

Multiple baseline classification models have been trained and evaluated to monitor the classification performance on Parkinson's disease data.

```

1 #Importing various classification algorithm to find which algorithm
2 suits the best for the dataset:
3 from sklearn.neural_network import MLPClassifier
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.svm import SVC
6 from sklearn.gaussian_process import GaussianProcessClassifier
7 from sklearn.ensemble import GradientBoostingClassifier
8 from sklearn.gaussian_process.kernels import RBF
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.ensemble import ExtraTreesClassifier
11 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
12 from sklearn.naive_bayes import GaussianNB
13 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
14 from sklearn.linear_model import SGDClassifier

```

Listing 21: Importing Classifiers

```

1 names = ["Nearest_Neighbors", "Linear_SVM", "Polynomial_SVM", "RBF_SVM",
2          "Gaussian_Process",
3          "Gradient_Boosting", "Decision_Tree", "Extra_Trees", "
4          Random_Forest", "Neural_Net", "AdaBoost",
5          "Naive_Bayes", "QDA", "SGD"]
6 classifiers = [

```

```

6     KNeighborsClassifier(3),
7     SVC(kernel="linear", C=0.025),
8     SVC(kernel="poly", degree=3, C=0.025),
9     SVC(kernel="rbf", C=1, gamma=2),
10    GaussianProcessClassifier(1.0 * RBF(1.0)),
11    GradientBoostingClassifier(n_estimators=100, learning_rate=1.0),
12    DecisionTreeClassifier(max_depth=5),
13    ExtraTreesClassifier(n_estimators=10, min_samples_split=2),
14    RandomForestClassifier(max_depth=5, n_estimators=100),
15    MLPClassifier(alpha=1, max_iter=1000),
16    AdaBoostClassifier(n_estimators=100),
17    GaussianNB(),
18    QuadraticDiscriminantAnalysis(),
19    SGDClassifier(loss="hinge", penalty="l2")]

```

Listing 22: Fitting the classification models

```

1 scores = []
2 for name, clf in zip(names, classifiers):
3     clf.fit(X_train, y_train)
4     score = clf.score(X_test, y_test)
5     scores.append(score)
6
7 df = pd.DataFrame()
8 df['name'] = names
9 df['score'] = scores
10 df
11
12 #Plot the results for comparison
13 sns.set(style="whitegrid")
14 ax = sns.barplot(y="name", x="score", data=df)

```

Listing 23: Evaluate and compare the classifiers

The plot has been shown in figure 13. It can be seen that XG Boost model shows the best performance for classification.

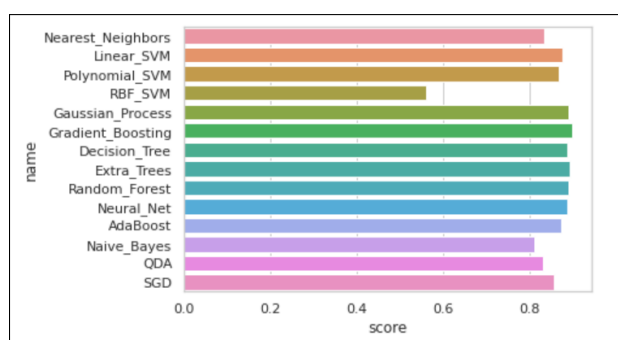


Figure 13: Comparison Plot

Further the classification scores can be seen in figure 14 which shows XG Boost algorithm with accuracy of 89.76% whereas deep learning neural network shows the accuracy of 88.75%. Even though the score for most of the machine learning algorithms are high, they do not consider the multiple time steps for multiple visits of patients.

	name	score
0	Nearest_Neighbors	0.833245
1	Linear_SVM	0.874934
2	Polynomial_SVM	0.868074
3	RBF_SVM	0.559894
4	Gaussian_Process	0.890237
5	Gradient_Boosting	0.897625
6	Decision_Tree	0.887071
7	Extra_Trees	0.891821
8	Random_Forest	0.889182
9	Neural_Net	0.887599
10	AdaBoost	0.874406
11	Naive_Bayes	0.810026
12	QDA	0.829551
13	SGD	0.857520

Figure 14: Classification Scores for different Models

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M. et al. (2016). Tensorflow: A system for large-scale machine learning, *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, pp. 265–283.
- Carneiro, T., Da Nóbrega, R. V. M., Nepomuceno, T., Bian, G.-B., De Albuquerque, V. H. C. and Reboucas Filho, P. P. (2018). Performance analysis of google colaboratory as a tool for accelerating deep learning applications, *IEEE Access* **6**: 61677–61685.
- Matloff, N. (2011). *The art of R programming: A tour of statistical software design*, No Starch Press.
- Van Rossum, G., Drake, F. L. et al. (2000). *Python reference manual*, iUniverse Indiana.