

Configuration Manual

Analysis and predictions of CO2 emissions using Neural
Networks

Msc Data Analytics

Jeet Jaikishan Vyas

Student ID: x19197161

School of Computing

National College of Ireland

Supervisor: Prof. Jorge Basilio

National College of Ireland
MSc Project Submission Sheet



School of Computing

Student Name: Jeet Jaikishan Vyas

Student ID: X19197161

Programme: Msc Data Analytics

Year: 2021

Module: Research Project

Lecturer:

Submission

Due Date: 16/12/2021

Project Title: Analysis and prediction of CO2 emissions using deep learning

Word Count: 1736 **Page Count:** 11

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: Jeet Jaikishan Vyas

Date: 16/12/2021

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):	

Configuration Manual

Jeet Jaikishan Vyas
Student ID: x19197161

1 Introduction

The configuration manual document will state the required software and hardware required for the research project “Analysis and predictions of CO2 emission using Neural Networks”. The manual contains the code for the Data collection/cleaning, clustering and models.

2 System and Software Requirements

Below are the required system configurations required for the research to be conducted.

2.1 Hardware Requirements:

Processor: Intel Core i5 – 7300U CPU @ 2.30 GHz to 2.40 GHz
Storage Capacity: 1 TB (Terabyte) HDD Hard Disk
System Type: 64-bit processor, x64
GPU: AMD Radeon / NVIDIA GeForce
Operating System: Windows 10 or 11 (64-bit operating system)
Ram: 8 GB

2.2 Software Requirements:

Programming and Data loading / Processing / Cleaning / Modelling: Jupyter Notebook by Anaconda.
Clustering: RStudio
Visualizations: Python / PowerBI.
Other tools: Microsoft word, Snipping tool, Microsoft excel.
Microsoft word are used for creation of tables and showcasing the figures. Snipping tool is used for getting the screenshot. Microsoft excel contains the data used for the research.

3 Data collection and cleaning

The flow of the research project is showcased below:

3.1 Collection of data and cleaning

The data used for the research in the format of CSV file has been presented. This showcases the raw data which has been used for the research implementation.

The data about the countries has been shown in Fig 2. The data contains 247 rows and 31 columns. The information about all the countries is distinguished in the data. The country names in three different columns, currency units for the countries, encoded country names, the population census and the industrial, water withdrawal, agricultural data. The total data shows the information that are relevant for a country.

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt

1 df_itr = pd.read_csv('E:\\NCI\\Semester 3\\Thesis\\Dataset\\Indicators.csv')
2 df_ctrs = pd.read_csv("E:\\NCI\\Semester 3\\Thesis\\Dataset\\Country.csv")
3 df_series = pd.read_csv("E:\\NCI\\Semester 3\\Thesis\\Dataset\\Series.csv")

```

Figure 3.3: Loading Libraries and Data Frame Creation

The raw data was loaded on the python environment by anaconda. The tables indicators, country and series were selected for the research. For the research, indicators and country data is being considered. The series data has been just used for information purposes. The pandas library in python used for data manipulation and analysis. The data should be loaded into data frames.

```

1 df_ctrs.isna().sum()

```

CountryCode	0
ShortName	0
TableName	0
LongName	0
Alpha2Code	3
CurrencyUnit	33
SpecialNotes	83
Region	33
IncomeGroup	33
Wb2Code	1
NationalAccountsBaseYear	42
NationalAccountsReferenceYear	193
SnaPriceValuation	49
LendingCategory	103
OtherGroups	188
SystemOfNationalAccounts	33
AlternativeConversionFactor	200
PppSurveyYear	56
BalanceOfPaymentsManualInUse	66
ExternalDebtReportingStatus	123
SystemOfTrade	47
GovernmentAccountingConcept	86
ImfDataDisseminationStandard	64
LatestPopulationCensus	34
LatestHouseholdSurvey	100
SourceOfMostRecentIncomeAndExpenditureData	89
VitalRegistrationComplete	135
LatestAgriculturalCensus	105
LatestIndustrialData	134
LatestTradeData	61
LatestWaterWithdrawalData	67
dtype: int64	

Figure 3.4: Data cleaning for country data

```

1 df_itr.isna().sum()
CountryName      0
CountryCode      0
IndicatorName     0
IndicatorCode     0
Year             0
Value            0

```

Figure 3.5: Sum of null values in indicators data.

3.1.1 Data Cleaning

To ensure there are no null values in the data, data cleaning techniques were implemented. The pandas library allows calculating the null values and filling the null values. Figure 3.3 shows the representation. The data frame was stored in a new data frame named countries which will be used for further research steps. On the other hand, after implementation there were no null values found in the indicators data.

```

1 df_itr.duplicated()
0      False
1      False
2      False
3      False
4      False
...
5656453  False
5656454  False
5656455  False
5656456  False
5656457  False
Length: 5656458, dtype: bool

1 df_countries.duplicated()
0      False
1      False
2      False
3      False
4      False
...
242     False
243     False
244     False
245     False
246     False
Length: 247, dtype: bool

1 df_series.duplicated()
0      False
1      False
2      False
3      False
4      False

```

Figure 3.5: Duplicated values evaluation on the data

In order to find any duplicates in the data, the duplicated function shows if there are any duplicates in the data.

4 Data Visualisations

For the research, there were bar graphs visualized using the matplotlib library in python. Three indicators were chosen for this step. The indicators were CO2 emissions per metric tons per capita, Age dependency ratio, CO2 emissions from liquid fuel consumptions. Visualisations for 6 countries were implemented for the research. The countries for this are United states of America, Australia, Ireland, India, China and Ireland. The graphs showed the

trend of the indicators from the year 1960 – 2010. The behaviour of the indicators according to the countries shows the flow. The indicators have been named and the plot function for bar graphs should be implemented. The world data shows the overall values of the specified indicator.

The following images show the graphical representations of overall world for the three indicators, the representations for the 6 countries including the world are visualized in python.

```

1 indicator1 = 'CO2 emissions \ (metric'
2 world0 = 'World'
3
4 x = df_process['IndicatorName'].str.contains(indicator1)
5 y = df_process['CountryName'].str.contains(world0)
6
7 total = df_process[x & y]
8
9 x = total['Year'].values
10 y = total['Value'].values
11 plt.title("CO2 Emissions metric ton per capita - World")
12 plt.xlabel("Years")
13 plt.ylabel("Emission range")
14 plt.bar(x,y)

```

Figure 4.1: CO2 emissions metric tons per capita – World

The trends of CO2 emissions metric tons per capita had seen a fluctuating one. From 1960 there has been a rise in the emissions up to 1980. There was slight fall after which the trend was quite stable and a rise in the emissions were seen after 2000 which kept rising until 2010.

```

1 indicator2 = 'Age dependency ratio'
2 world0 = 'World'
3
4 x1 = df_process['IndicatorName'].str.contains(indicator2)
5 y1 = df_process['CountryName'].str.contains(world0)
6
7 total = df_process[x1 & y1]
8
9 x1 = total['Year'].values
10 y1 = total['Value'].values
11 plt.title("Age dependency ratio - World")
12 plt.xlabel("Years")
13 plt.ylabel("Total %")
14 plt.bar(x1,y1)
15 plt.show()

```

Figure 4.2: Age dependency Ratio – World

The age dependency ratio follows the formula for the number of people engaged in the working class compared to the people not engaged in the working class. The people distinguished in the working population are between the age group of 18 – 59 comparing the people in age group of 0 – 17 and 59 and above. The classes are defined as minors and senior citizens or retired individuals. The indicator should be considered as an effect on the economy of an individual country.

```

1 indicator3 = 'CO2 emissions from liquid fuel consumption'
2 world0 = 'World'
3
4 x2 = df_process['IndicatorName'].str.contains(indicator3)
5 y2 = df_process['CountryName'].str.contains(world0)
6
7 total = df_process[x2 & y2]
8
9 x2 = total['Year'].values
10 y2 = total['Value'].values
11 plt.title("CO2 Emissions liquid fuel consumption - World")
12 plt.xlabel("Years")
13 plt.ylabel("Emission range")
14 plt.bar(x2,y2)
15 plt.show()

```

Figure 4.3: CO2 emissions from liquid fuel consumptions – world

The code showcases the indicator liquid fuel consumptions for the world. The overall trends represent the indicator name, country name includes the world and the value for the same.

5 Clustering

```
1 library(readr)
2 library(reshape)
3 library(cluster)
4
5 Indicators <- read.csv("E:\\NCI\\Semester 3\\Thesis\\Dataset\\Indicators.csv")
6
7 code <- c('SE.XPD.CTOT.ZS', 'GB.XPD.RSDV.GD.ZS', 'SE.ADT.LITR.ZS', 'EG.ELC.ACCS.ZS',
8           'SL.UEM.LTRM.ZS', 'SL.UEM.TOTL.ZS', 'FB.CBK.BRCH.P5', 'GC.DOD.TOTL.GD.ZS', 'EN.ATM.CO2E.PC',
9           'BX.KLT.DINV.WD.GD.ZS', 'NY.GDP.PCAP.CD', 'NY.GNP.MKTP.CD', 'FP.CPI.TOTL.ZG', 'NY.ADJ.NNTY.PC.CD',
10          'PA.NUS.FCRF', 'DT.DOD.PVLX.GN.ZS', 'IC.REG.DURS', 'FI.RES.TOTL.DT.ZS', 'SH.MED.BEDS.ZS',
11          'SP.POP.1564.TO.ZS', 'TX.VAL.TECH.CD', 'SP.DYN.CDRT.IN', 'SP.DYN.CBRT.IN', 'SI.POV.GINI',
12          'IC.BUS.EASE.XQ', 'IT.NET.USER.P2', 'IT.CEL.SETS.P2', 'IQ.CPA.TRAD.XQ', 'TM.VAL.FUEL.ZS.UN',
13          'TX.VAL.MANF.ZS.UN', 'NE.TRD.GNFS.ZS')
14
15 indic <- cast(Indicators, CountryName ~ IndicatorCode, mean)
16 indic <- indic[,c('CountryName', code)]
17 row.names(indic) <- indic[,1]
18 indic <- indic[,c(-1,-2,-4,-6,-9,-14,-17,-19,-25,-29)]
19 indic <- na.omit(indic)
20 Cnames <- rownames(indic)
21
22 indic.norm <- sapply(indic, scale)
23 row.names(indic.norm) <- Cnames
24
25 set.seed(42)
26 indic_cluster <- kmeans(indic.norm, 3)
27 indic_cluster$size
28
29 par(mfrow=c(1,1))
30 clusplot(indic.norm, indic_cluster$cluster, color = TRUE, shade = TRUE, labels = 2, lines = 0)
31
32 indic.norm <- data.frame(cbind(indic_cluster$cluster, rownames(indic), indic.norm))
33 colnames(indic.norm) <- c("Cluster_No", "CountryName", colnames(indic.norm)[c(-1,-2)])
34 rownames(indic.norm) <- NULL
35
36 developed_countries <- subset(indic.norm, Cluster_No == 1, select = CountryName)
37 developing_countries <- subset(indic.norm, Cluster_No == 2, select = CountryName)
38 underdeveloped_countries <- subset(indic.norm, Cluster_No == 3, select = CountryName)
39
40 print(developed_countries)
41 print(developing_countries)
42 print(underdeveloped_countries)
```

Figure 5: Clustering of Countries

The aim for clustering was to develop clusters of the countries in the indicators data. For this necessary step were carried in RStudio. The libraries for conducting are readr, reshape and cluster. The readr library helps in reading the CSV format files, reshape library helps in reshapes a data frame between wide and long formats in order to achieve repeated measurements in separate columns. The cluster library helps in getting the clusters using the algorithms on a specific data. The indicator data loaded in the environment as the first initial step. The indicator codes should be taken into consideration for creation of clusters. 31 indicators were chosen for the same. The transformation of data carried out by changing the formation of rows to columns. The cast function used to get the aggregate data with the aggregate function and formula. In this case, the relationship between the columns using the tilde function. The transformation makes the data ready to standardise and apply the algorithm for clustering. The sapply function takes the data frame as an input and the potential outcome will be a vector or matrix, therefore the data (indic) has been input as an object along with the function of scale. for the same to get a vector or matrix.

The K means clustering algorithm has been used for the cluster implementation. As KMeans being the most simple and quick clustering algorithm as it creates clusters with the mean

value close to the one feature to another. This helps in making a sophisticated form of clustering. Three clusters were formed as the KMean value was set at 3. Plotting the cluster created three clusters and the information about the clusters with the countries in the appropriate countries will be achieved.

6 Time Series and Models

```

1 from tqdm import tqdm
2 from time import time

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.nn.functional as f

1 from torch.utils.data import Dataset, DataLoader, TensorDataset

```

Figure 6: Libraries and utils

The libraries required to get the models using Pytorch. nn is the module for neural network models. Optim is the module which helps in optimizing the models deployed. Nn functional involves all functions which are applied on models. Tqdm library is used to show the progress bar.

```

1 def get_valid_range(country_code, ECON_COLS, CO2_COL, time_step=5, verbose=False):
2
3     try:
4         x = COUNTRIES_DATA[country_code].loc[ECON_COLS + [CO2_COL], :]
5     except KeyError:
6         if verbose:
7             print('{} - {} miss the following indicators:\n'.format(country_code, CODE_TO_NAMES[country_code]),
8                   ', '.join(set(ECON_COLS + [CO2_COL]) - set(COUNTRIES_DATA[country_code].index)))
9         return [], None
10
11     y_econ = (x.loc[ECON_cols, :].isna().sum(axis=0) == 0)['Value']
12     y_co2 = ~x.loc[CO2_COL, :].isna()['Value']
13
14     has_co2 = [year for year in y_co2[y_co2].index]
15     for interval in continuous_subarrays(has_co2, step=1, min_len=time_step):
16         y_co2.iloc[max(0, interval[0] - time_step):interval[0]] = True
17
18     y = y_econ & y_co2
19
20     is_valid = [year if y[year] else 0 for year in y.index]
21     res = continuous_subarrays(is_valid, step=1, min_len=time_step)
22     if verbose and len(res) == 0:
23         print('{} - {} has no valid interval.'.format(country_code, CODE_TO_NAMES[country_code]))
24
25     return res, x

```

Figure 6: Time series Data processing

In order to get the data ready for implementation of time series, a valid range of data need to be prepared which will withstand the factors while implementing the time series. Here, the country code (countries data), ECON_COLS are the economic indicator columns, CO2_COL are the CO2 emission indicator columns which are the target features for implementation. There has to be a set time step that will consider the time interval for the time series. Any null

values have been ignored. A continuous subarray created for validating the time step with the length of the same. The countries data and time step were returned.

```

1 def CO2TimeSeries(ECON_COLS, CO2_COL, time_step, valid_pct=0.2):
2
3     train_data, val_data = None, None
4     train_targets, val_targets = np.array([]), np.array([])
5
6     for country in df_process['CountryCode'].unique():
7
8         valid_ranges, df_countries = get_valid_range(country, ECON_COLS, CO2_COL, time_step)
9         if len(valid_ranges) == 0:
10             continue
11
12         for indicator in df_countries.index:
13             if indicator != CO2_COL:
14                 df_countries.loc[indicator, :] = \
15                     (df_countries.loc[indicator, :] - indicators_mean[indicator]) / indicators_std[indicator]
16             else:
17                 df_countries.loc[indicator, :] = [np.log(num) if not np.isnan(num) else np.nan \
18                     for num in df_countries.loc[indicator, :]]
19
20                 country_data, country_targets = None, np.array([])
21             econ_df, co2_df = df_countries.loc[ECON_COLS, :], df_countries.loc[CO2_COL, :]
22
23         for interval in valid_ranges:
24             for i in range(interval[0], interval[1] - time_step + 1):
25                 obs = econ_df.iloc[:, i:i+time_step].to_numpy()
26                 obs = obs[np.newaxis, :, :]
27                 if country_data is not None:
28                     country_data = np.concatenate((country_data, obs), axis=0)
29                 else:
30                     country_data = obs
31                 country_targets = np.append(country_targets, co2_df.iloc[i + time_step - 1])
32
33         num_train = round((1 - valid_pct) * country_data.shape[0])
34         country_train_data = country_data[:num_train, :, :]
35         country_train_targets = country_targets[:num_train]
36         country_val_data = country_data[num_train:, :, :]
37         country_val_targets = country_targets[num_train:]
38
39         assert country_data.shape[0] == len(country_targets)
40         assert len(country_val_targets) + len(country_train_targets) == country_data.shape[0]
41         assert len(country_val_targets) == country_val_data.shape[0]
42         assert len(country_train_targets) == country_train_data.shape[0]
43
44         if train_data is not None:
45             train_data = np.concatenate((train_data, country_train_data), axis=0)
46             val_data = np.concatenate((val_data, country_val_data), axis=0)
47         else:
48             train_data = country_train_data
49             val_data = country_val_data
50         train_targets = np.append(train_targets, country_train_targets)
51         val_targets = np.append(val_targets, country_val_targets)
52
53     assert train_data is not None,
54     return TensorDataset(torch.tensor(train_data), torch.tensor(train_targets)), \
55         TensorDataset(torch.tensor(val_data), torch.tensor(val_targets))

```

Figure 6.1: Time Series Implementation

Implementing the time series with the indicators and countries data. The train and validation were created using arrays in numpy library. The range of indicators if not null for economic and CO2 emissions will be considered. The second range of countries in valid range are not null then the data will concatenate else it will append. The same process was carried out to prepare the training and validation implementation.

```

1 train_ds, val_ds = CO2TimeSeries(ECON_cols, CO2_COL='EN.ATM.CO2E.KT', time_step=TIME_STEP, valid_pct=VALID_PCT)
2 print("Number of Training observations:", len(train_ds))
3 print("Number of Validation observations:", len(val_ds))
4
5 BATCH_SIZE = 32
6
7 train_dl = DataLoader(train_ds, batch_size=BATCH_SIZE, shuffle=True)
8 val_dl = DataLoader(val_ds, batch_size=BATCH_SIZE, shuffle=True)

```

Figure 6.2: Time series training and validation

The training and validation for the time series with the economic indicators and CO2 emissions metric ton per capita for CO2_COL with the valid time step declared in the implementation phase. The batch size set to 32 shows the number of samples to be shown to the networks. The batch size used to fit the model and control the number of predictions made at a certain time.

```

1 class LSTM(nn.Module):
2
3     def __init__(self, input_dims, hidden_dims, output_dims,):
4         super().__init__()
5         self.rec_layer = nn.LSTM(input_dims, hidden_dims, batch_first=True)
6         self.fc = nn.Linear(hidden_dims, output_dims)
7
8     def forward(self, x):
9
10        x = self.reshape(x)
11        outputs, (hidden_states, cell_states) = self.rec_layer(x)
12        outputs = outputs[:, -1, :]
13        outputs = outputs.squeeze(1)
14        outputs = self.fc(outputs)
15
16        return outputs.squeeze(1)
17
18    def reshape(self, x):
19        return torch.reshape(x, shape=(x.shape[0], x.shape[2], x.shape[1]))

```

Figure 6.3: 1Layer LSTM Model

```

1 class GRU(nn.Module):
2
3     def __init__(self, input_dims, hidden_dims, output_dims,):
4         super().__init__()
5         self.rec_layer = nn.GRU(input_dims, hidden_dims, batch_first=True)
6         self.fc = nn.Linear(hidden_dims, output_dims)
7
8     def forward(self, x):
9
10        x = self.reshape(x)
11        outputs, hidden_states = self.rec_layer(x)
12        outputs = outputs[:, -1, :]
13        outputs = outputs.squeeze(1)
14        outputs = self.fc(outputs)
15
16        return outputs.squeeze(1)
17
18    def reshape(self, x):
19        return torch.reshape(x, shape=(x.shape[0], x.shape[2], x.shape[1]))

```

Figure 6.3: Gated Recurring Unit (GRU)

```

1 class FeedForward(nn.Module):
2
3     def __init__(self, input_dims, time_step):
4         super().__init__()
5
6         # 5 Layers
7         self.flatten = nn.Flatten()
8         self.fc1 = nn.Linear(input_dims * time_step, 512)
9         self.fc2 = nn.Linear(512, 256)
10        self.fc3 = nn.Linear(256, 256)
11        self.fc4 = nn.Linear(256, 128)
12        self.fc5 = nn.Linear(128, 1)
13
14        self.activation = nn.ReLU()
15
16        def forward(self, x):
17
18            x = self.flatten(x)
19            x = self.activation(self.fc1(x))
20            x = self.activation(self.fc2(x))
21            x = self.activation(self.fc3(x))
22            x = self.activation(self.fc4(x))
23            x = self.fc5(x)
24
25        return

```

Figure 6.4: Feed Forward Model creation

Three models were created 1 layer LSTM, GRU and Feed Forward. They were obtained from pytorch in the neural network. All models are recurrent models only LSTM and GRU were created with the required dimensions for the models. The last hidden layer is grabbed and getting a reshaped form of hidden and output batch. The feed forward model was created with the linear functions consisting 4 layers and a dimension with inputs along with the time step. The observations are flattened.

```

1 import torch.backends.cudnn as cudnn
2 torch.cuda.empty_cache()
3 cudnn.benchmark = True
4
5 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

```

```

1 lstm_model = LSTM(input_dims=len(ECON_cols), hidden_dims=32, output_dims=1).to(device)
2 lstm_optimizer = optim.Adam(
3     lstm_model.parameters(),
4     lr=0.001
5 )
6 lstm_loss_fn = torch.nn.MSELoss()

```

```

1 train_dl = train_dl
2 train_model = lstm_model
3 train_optimizer = lstm_optimizer
4 train_loss_fn = lstm_loss_fn

```

```

1

```

```

1 gru_model = GRU(input_dims=len(ECON_cols), hidden_dims=32, output_dims=1).to(device)
2 gru_optimizer = optim.Adam(
3     gru_model.parameters(),
4     lr=0.001
5 )
6 gru_loss_fn = torch.nn.MSELoss()

```

```

1 train_dl = train_dl
2 train_model = gru_model
3 train_optimizer = gru_optimizer
4 train_loss_fn = gru_loss_fn

```

The cudnn benchmark is set true for hardware validation. CPU has been set for carrying out the models to run the dimensions.

The LSTM and GRU models are trained and optimized with calculating the loss using the loss function. The loss function will calculate the MSE loss.


```

1 train_np_x = np.stack([obs for obs, label in train_ds], axis=0)
2 train_np_y = np.stack([label for obs, label in train_ds])
3 val_np_x = np.stack([obs for obs, label in val_ds], axis=0)
4 val_np_y = np.stack([label for obs, label in val_ds])
5
6 train_np_x = train_np_x.reshape((train_np_x.shape[0], -1))
7 val_np_x = val_np_x.reshape((val_np_x.shape[0], -1))

```

```

1 xgb_model = XGBRegressor()
2 xgb_model.fit(train_np_x, train_np_y)

```

Figure 6.5: Random Forest with Gradient Boosting

The training and validation dimensions have been converted to stacks using numpy. They are reshaped and fitted in the model. The results forecasted using matplotlib and box plots should give the validation loss for the number of observations.

References

- [1] Adi Bronstein (2017) *A quick introduction to the “Pandas” Python Library* Available at: <https://towardsdatascience.com/a-quick-introduction-to-the-pandas-python-library-f1b678f34673> [25 October 2021].
- [2] Jason Brownlee (2017) *Stacked Long Short-term Memory Networks* Available at: <https://machinelearningmastery.com/stacked-long-short-term-memory-networks/> [10 November 2021]
- [3] Tasnuva Zaman (2019) *Building Neural Network Using Pytorch* Available at: <https://towardsdatascience.com/building-neural-network-using-pytorch-84f6e75f9a> [12 November 2021]
- [4] GeeksforGeeks (2020) *K-Means Clustering in R programming* Available at: <https://www.geeksforgeeks.org/k-means-clustering-in-r-programming/> [23 October 2021]