

Configuration Manual

MSc Research Project
Data Analytics

Ratna Pillai
Student ID: x18134297

School of Computing
National College of Ireland

Supervisor: Dr. Muhammad Iqbal

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Ratna Pillai
Student ID:	x18134297
Programme:	Data Analytics
Year:	2019
Module:	MSc Research Project
Supervisor:	Dr. Muhammad Iqbal
Submission Due Date:	12/12/2019
Project Title:	Configuration Manual
Word Count:	1836
Page Count:	26

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:	12 th December 2019

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

Ratna Pillai

x18134297

MSc Research Project in Data Analytics

12th December 2019

1 Introduction

This configuration manual specifies the hardware, software requirements and the programming phases of the implementation of the below research project in detail:

“Optimized Predictive Modelling to Unfold the Links of Crime with Education, Safety and Climate in Chicago”

2 System Configurations

2.1 Hardware

- **Processor:** Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz 1.80 GHz
- **RAM:** 8 GB
- **System Type:** Windows OS, 64-bit
- **GPU:** Intel(R) UHD Graphics Family, 4GB
- **Storage:** 1 TB HDD

2.2 Software

- **Microsoft Excel 2016:** This spreadsheet tool offered by Microsoft has been used for storing the downloaded datasets in flat files as csv (comma separated values).
- **Anaconda Distribution-Jupyter Notebook:** This is an open source software which has been downloaded from the anaconda distribution website.¹ This distribution support platforms like jupyter notebooks, spyder, R studio to run machine learning models using R or python. Exploratory data analysis, manipulation of data, pre-processing, transformation and visualizations in this study were done using Python (version 3.6.5) on Jupyter notebooks using this distribution.

¹<https://www.anaconda.com/distribution/>

- **Google Colaboratory:** Also called Colab, this is a free cloud service that aids the users with free GPU services to run machine learning models on an environment similar to Jupyter notebooks. For this study, colab is used for modelling and hyper parameter optimization. For enabling GPU settings, from the Runtime menu on the notebook screen, select change runtime type to GPU as shown in Figure 1. There is also an additional option to change it to TPU, which is efficient to use in cases of high data volume.

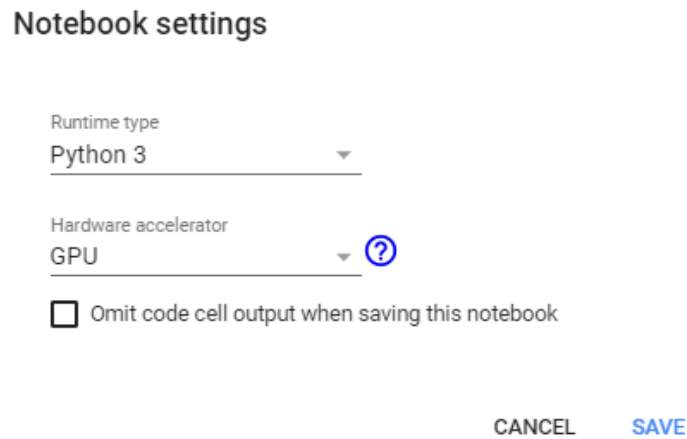


Figure 1: GPU Settings: Google Colaboratory

- **Power BI Desktop:** Part of Visualizations were done using Microsoft Power BI desktop app. This software is downloaded from Microsoft store website.²

3 Project Development

Implementation of this research work is entirely done using Python programming. Initial stages of the research project development includes data pre-processing, merging all the datasets in consideration, normalization and one hot encoding. Followed by the data preparation activities, predictive modelling is done using machine learning techniques in python using the sk-learn (scikit-learn) and keras libraries.

3.1 Data Preparation

Importing of the datasets and data manipulation is done using the pandas (dataframe) and numpy (arrays) libraries. Following sections detail the pre-processing steps carried out in each dataset.

3.1.1 Crime Dataset

Pre-processing of crime data primarily involved fixing the missing latitudes and longitudes as shown in Figure 2. Block addresses (which was already clean when the data was downloaded from the source) of the missing location co-ordinates were passed to the

²<https://www.microsoft.com/en-in/p/power-bi-desktop/9ntxr16hnw1t?activetab=pivot:overviewtab>

wrapper `censusbatchgeocoder`. This function contacts the geocodes API (Application Programming Interface) in the US (United States) and returns a dictionary of address detail for a particular address. This dictionary is converted in to a new dataframe and the missing values are filled using `lookup` and `fill_na` (Not Applicable) functions.

```

186 import censusbatchgeocoder
187 #Extract only null rows for location & zip codes
188 locnull = narcotics[narcotics['Latitude'].isnull()]
189 locnull.shape
190 city = 'Chicago'
191 state = 'IL'
192 address = locnull['Block']
193 address = address.drop_duplicates()
194 addressdata = pd.DataFrame(address)
195 addressdata['city'] = city
196 addressdata['state'] = state
197 addressdata.columns = ['address','city','state']
198 addressdata = addressdata.reset_index()
199 addressdata['id'] = np.arange(1,len(addressdata)+1)
200 addressdata = addressdata.drop(columns=['index'])
201 #print (addressdata)
202 #Rearrange columns - addressdata
203 cols = list(addressdata.columns)
204 a, b = cols.index('address'), cols.index('id')
205 cols[b], cols[a] = cols[a], cols[b]
206 addressdata = addressdata[cols]
207 addressdata.head()
208 fetchaddress = addressdata.to_dict("records")
209 results = censusbatchgeocoder.geocode(fetchaddress.to_dict("records"), zipcode=None)
210 #Lookup for zipcode
211 narcotics.Zip_Codes.replace('NaN', np.NaN, inplace=True)
212 narcotics.loc[narcotics['Zip_Codes'].isnull(), 'Zip_Codes'] = narcotics['Block'].map(locdata.ZipCode)
213 #Lookup for latitude
214 narcotics.loc[narcotics['Latitude'].isna(), 'Latitude'] = narcotics['Block'].map(locdata.LAT)
215 #Lookup for longitude
216 narcotics.loc[narcotics['Longitude'].isna(), 'Longitude'] = narcotics['Block'].map(locdata.LON)
217 #Fill missing Zipcodes
218 s = locdata.set_index('Block')['ZipCode']
219 narcotics['Zip_Codes'] = narcotics['Zip_Codes'].fillna(narcotics['Block'].map(s))

```

Figure 2: Handling of missing latitudes and longitudes in Crime dataset

3.1.2 High Schools Dataset

Four datasets for high school report (academic years 2015-2018) were downloaded from Chicago data portal and saved as csv (comma separated values) files. Each year school data was complex enough with more than 150 attributes and selection of relevant attributes was a challenge. Also, there were missing values which had to be handled as shown in Figure 3. Since the data was maintained school wise, the median value was used to replace the missing values in school data.

```

24 #2018 edu data
25 edurefined = edueight[['Long_Name', 'Community_Area', 'School_Latitude', 'School_Longitude', 'Primary_Category',
26 'Teacher_Attendance_Year_1_Pct', 'Teacher_Attendance_Year_2_Pct', 'Student_Attendance_Year_1_Pct',
28 'Suspensions_Per_100_Students_Year_1_Pct', 'Suspensions_Per_100_Students_Year_2_Pct',
29 'Average_Length_Suspension_Year_1_Pct', 'Average_Length_Suspension_Year_2_Pct', 'Mobility_Rate_Pct',
30 'Chronic_Trucancy_Pct', 'One_Year_Dropout_Rate_Year_1_Pct'],
34 #Compute average misconduct to suspensions rate
35 edurefined['Avg_Misconduct_Rate'] = (edurefined['Misconducts_To_Suspensions_Year_1_Pct']
36 + edurefined['Misconducts_To_Suspensions_Year_2_Pct'])/2
37 #Compute average suspension rate
38 edurefined['Avg_Suspension_Rate'] = (edurefined['Suspensions_Per_100_Students_Year_1_Pct']
40 #Compute average attendance rate
41 edurefined['Avg_Student_Attendance_Rate'] = (edurefined['Student_Attendance_Year_1_Pct']
43 #Compute average teacher attendance rate
44 edurefined['Avg_Teacher_Attendance_Rate'] = (edurefined['Teacher_Attendance_Year_1_Pct']
45 + edurefined['Teacher_Attendance_Year_2_Pct'])/2
46 edurefined['Average_Length_Suspension_Year_1_Pct'] = edurefined['Average_Length_Suspension_Year_1_Pct'].astype(str)
47 edurefined['Average_Length_Suspension_Year_2_Pct'] = edurefined['Average_Length_Suspension_Year_2_Pct'].astype(str)
48 edurefined['Average_Length_Suspension_Year_1_Pct'] = edurefined['Average_Length_Suspension_Year_1_Pct'].str.replace(' days', '')
49 edurefined['Average_Length_Suspension_Year_2_Pct'] = edurefined['Average_Length_Suspension_Year_2_Pct'].str.replace(' days', '')
50 edurefined['Average_Length_Suspension_Year_1_Pct'] = edurefined['Average_Length_Suspension_Year_1_Pct'].astype(float)
51 edurefined['Average_Length_Suspension_Year_2_Pct'] = edurefined['Average_Length_Suspension_Year_2_Pct'].astype(float)
52 edurefined['Avg_Suspension_Days'] = (edurefined['Average_Length_Suspension_Year_2_Pct']
53 + edurefined['Average_Length_Suspension_Year_1_Pct'])/2
54 edurefined.head()
55 #Replace missing values with median
56 edurefined['Avg_Misconduct_Rate'].fillna((edurefined['Avg_Misconduct_Rate'].median()), inplace=True)
80 #Convert location c-ordinates to geohash
81 import pygeohash as pgh
82 edurefined['geohash'] = edurefined.apply(lambda x: pgh.encode(x.School_Latitude, x.School_Longitude, precision=5), axis=1)
83 #Group by geohash
84 edu_df = edurefined.groupby('geohash').mean().reset_index()

```

Figure 3: High school education data pre-processing

- Firstly, only numerical attributes were considered from each education dataset. Suspension days column consisted of "days" keyword appended to number of days which was cleaned to retain only the number and "days" keyword was removed. Then, a check for missing value was done and filled using median value for that column.
- Since the misconduct rate, suspension rate, enrollment rate, freshman track rate were all maintained academic year wise (For eg. academic year 2015-16 data was stored in the file as suspension days year 1 and suspension days year 2), these columns were combined together using mean of both the columns and stored as Average Suspension Days.
- A similar approach was followed for the remaining three years data and the latitude and longitude of each school was converted to geohash using pygeohash library in python.
- Lastly, all the pre-processed data are merged and saved in a flat file in a csv format.

3.1.3 Locations Dataset

- Locations dataset identified were three namely, police station, speed camera locations and red light camera locations in Chicago. As these datasets comprised of latitude and longitude co-ordinates which were primarily required for this research, presence of missing values were checked as highlighted in Figure 4.

```

21 #Extract only required columns from police station, red-light camera and speed camera datasets
22 police = police[['LATITUDE','LONGITUDE']]
23 red = red[['LATITUDE','LONGITUDE']]
24 speed = speed[['LATITUDE','LONGITUDE']]
25 #Check for NA values in each dataset
26 red.isna().sum()
27 police.isna().sum()
28 speed.isna().sum()
29 #Verify the shape of each dataframe
30 print(red.shape)
31 print(speed.shape)
32 print(police.shape)

```

Figure 4: Locations dataset - Missing values check

- Nearest distance of police stations, red light cameras and speed cameras was calculated using the user defined distance function as denoted in Figure 5 for each of the crime instance.³

```

34 #Calculate distance between two latitudes and longitudes
35 import math
36 def distance(origin, destination):
37     lat1, lon1 = origin
38     lat2, lon2 = destination
39     radius = 6371 # km
40     dlat = math.radians(lat2-lat1)
41     dlon = math.radians(lon2-lon1)
42     a = math.sin(dlat/2) * math.sin(dlat/2) + math.cos(math.radians(lat1)) \
43         * math.cos(math.radians(lat2)) * math.sin(dlon/2) * math.sin(dlon/2)
44     c = 2 * math.atan2(math.sqrt(a), math.sqrt(1-a))
45     d = round(radius * c,2)
46     return d
47

```

Figure 5: Function to compute distance between two latitude and longitude co-ordinates

- Compute nearest speed camera distance from crime location For each speed camera location as highlighted in Figure 6, the latitudes and longitudes were used in the distance computation with the crime location (latitude and longitude). Then minimum distance for that crime instance was computed using `min()` function in python. Similarly, the code was run for the red-light camera locations and police stations as shown in Figure 7 and Figure 8.

³<https://gist.github.com/rochacbruno/2883505>

```

50 import datetime
51 #Calculate distance between two points and return minimum for speed camera locations
52 distanceall = pd.DataFrame(columns=['id','dist','latitude','longitude'])
53 distancemin = pd.DataFrame(columns=['id','dist','latitude','longitude'])
54 distanceall['dist'] = None
55 distancemin['dist'] = None
56 print(datetime.datetime.now())
57 for i in range(0,len(crime),1):
58     print(i)
59     for j in range(0,len(speed),1):
60         dist = distance((crime.iloc[i]['Latitude'].astype(float),crime.iloc[i]['Longitude'].astype(float)),
61                         (speed.iloc[j]['LATITUDE'].astype(float),speed.iloc[j]['LONGITUDE'].astype(float)))
62         lat = crime.iloc[i]['Latitude']
63         lon = crime.iloc[i]['Longitude']
64         distanceall.loc[j,'dist'] = dist
65         distanceall.loc[j,'latitude'] = lat
66         distanceall.loc[j,'longitude'] = lon
67         distanceall.loc[j,'id'] = i
68         d = len(distanceall)
69         if j == (len(speed)-1):
70             mindist = distanceall['dist'].min()
71         distancemin.loc[i,'dist'] = mindist
72         distancemin.loc[i,'latitude'] = lat
73         distancemin.loc[i,'longitude'] = lon
74         distancemin.loc[i,'id'] = i
75     print(datetime.datetime.now())

```

Figure 6: Computation of nearest speed camera distance

- Compute nearest Red-light camera distance from crime location

```

78 #Calculate distance between two points and return minimum for red-light camera locations
79 distanceall = pd.DataFrame(columns=['id','dist','latitude','longitude'])
80 distancemin = pd.DataFrame(columns=['id','dist','latitude','longitude'])
81 distanceall['dist'] = None
82 distancemin['dist'] = None
83 print(datetime.datetime.now())
84 for i in range(0,len(crime),1):
85     print(i)
86     for j in range(0,len(speed),1):
87         dist = distance((crime.iloc[i]['Latitude'].astype(float),crime.iloc[i]['Longitude'].astype(float)),
88                         (red.iloc[j]['LATITUDE'].astype(float),red.iloc[j]['LONGITUDE'].astype(float)))
89         lat = crime.iloc[i]['Latitude']
90         lon = crime.iloc[i]['Longitude']
91         distanceall.loc[j,'dist'] = dist
92         distanceall.loc[j,'latitude'] = lat
93         distanceall.loc[j,'longitude'] = lon
94         distanceall.loc[j,'id'] = i
95         d = len(distanceall)
96         if j == (len(speed)-1):
97             mindist = distanceall['dist'].min()
98         distancemin.loc[i,'dist'] = mindist
99         distancemin.loc[i,'latitude'] = lat
100        distancemin.loc[i,'longitude'] = lon
101        distancemin.loc[i,'id'] = i
102    print(datetime.datetime.now())
103

```

Figure 7: Computation of nearest red light camera distance

- Compute nearest police station distance from crime location


```

104 #Calculate distance between two points and return minimum for police stations
105 distanceall = pd.DataFrame(columns=['id','dist','latitude','longitude'])
106 distancemin = pd.DataFrame(columns=['id','dist','latitude','longitude'])
107 distanceall['dist'] = None
108 distancemin['dist'] = None
109 print(datetime.datetime.now())
110 for i in range(0,len(crime),1):
111     print(i)
112     for j in range(0,len(speed),1):
113         dist = distance((crime.iloc[i]['Latitude'].astype(float),crime.iloc[i]['Longitude'].astype(float)),
114                        (police.iloc[j]['LATITUDE'].astype(float),police.iloc[j]['LONGITUDE'].astype(float)))
115         lat = crime.iloc[i]['Latitude']
116         lon = crime.iloc[i]['Longitude']
117         distanceall.loc[j,'dist'] = dist
118         distanceall.loc[j,'latitude'] = lat
119         distanceall.loc[j,'longitude'] = lon
120         distanceall.loc[j,'id'] = i
121         d = len(distanceall)
122         if j == (len(speed)-1):
123             mindist = distanceall['dist'].min()
124         distancemin.loc[i,'dist'] = mindist
125         distancemin.loc[i,'latitude'] = lat
126         distancemin.loc[i,'longitude'] = lon
127         distancemin.loc[i,'id'] = i
128     print(datetime.datetime.now())

```

Figure 8: Computation of nearest speed camera distance

3.1.4 Weather Dataset

Daily weather data for the period 2015-2018 was scraped using (API) Application Programming Interface access from NCEI (formerly known as NOAA) portal. The data was scraped in 3 steps as followed by many data scientists as outlined in the medium portal⁴:

- Gained access to the API using token for authorization from NCEI web token site and accessed in python using requests library as demonstrated in Figure 11.⁵



Figure 9: Request the API access token

⁴<https://towardsdatascience.com/getting-weather-data-in-3-easy-steps-8dc10cc5c859>

⁵<https://www.ncdc.noaa.gov/cdo-web/token>

```

10 #needed to make web requests
11 import requests
12 #store the data we get as a dataframe
13 import pandas as pd
14 #convert the response as a structured json
15 import json
16 #mathematical operations on lists
17 import numpy as np
18 #parse the datetimes we get from NOAA
19 from datetime import datetime
20 #add the access token you got from NOAA
21 Token = 'XnvxAnqwtzCGLaepxKsuVihUwMhZkCti'
22 #Long Beach Airport station
23 station_id = 'GHCND:USW00094846'

```

Figure 10: API settings to access weather data

- Identified the relevant weather station for required data collection from the site.⁶

■ Data Tools: Find a Station

Retrieve weather records from observing stations by entering the desired location, data set, data range, and data can be specified as city, county, state, country, or ZIP code.

Figure 11: Find the relevant station

- With the help of python requests object, windspeed (wind), average temperature(avg_temp) and precipitation (prcp) values for three years on a daily basis was scraped. For each of the weather attribute, the datatype id was required to be modified in the URL of the request command.
- Average temperature data was fetched and stored in a dataframe. Refer Figure 12 for the python code.

⁶<https://www.ncdc.noaa.gov/cdo-web/datatools/findstation>

```

25 #initialize lists to store data - Average temperature (TAVG)
26 dates_temp = []
27 dates_prcp = []
28 temps = []
29 prcp = []
30 #for each year from 2015-2018 ...
31 for year in range(2015, 2019):
32     year = str(year)
33     print('working on year '+year)
34     #make the api call
35     r = requests.get('https://www.ncdc.noaa.gov/cdo-web/api/v2/data?datasetid=GHCND
36     &datatypeid=TAVG&limit=1000
37     &stationid=GHCND:USW00023129&startdate='+year+'-01-01
38     &enddate='+year+'-12-31', headers={'token':Token})
39     #load the api response as a json
40     d = json.loads(r.text)
41     #get all items in the response which are average temperature readings
42     avg_temps = [item for item in d['results'] if item['datatype']=='TAVG']
43     #get the date field from all average temperature readings
44     dates_temp += [item['date'] for item in avg_temps]
45     #get the actual average temperature from all average temperature readings
46     temps += [item['value'] for item in avg_temps]
47 #initialize dataframe
48 df_temp = pd.DataFrame()
49 #populate date and average temperature fields (cast string date to datetime and convert temperature from tenths of Celsius to Fahrenheit)
50 df_temp['date'] = [datetime.strptime(d, "%Y-%m-%dT%H:%M:%S") for d in dates_temp]
51 df_temp['avgTemp'] = [float(v)/10.0*1.8 + 32 for v in temps]

```

Figure 12: Average temperature data scraping

- Average windspeed was fetched and stored in a dataframe. Refer Figure 13 for the python code.

```

82 #initialize lists to store data - Average Wind Speed (AWND)
83 dates_temp = []
84 dates_prcp = []
85 temps = []
86 prcp = []
87 #for each year from 2015-2018 ...
88 for year in range(2015, 2019):
89     year = str(year)
90     print('working on year '+year)
91     #make the api call
92     r = requests.get('https://www.ncdc.noaa.gov/cdo-web/api/v2/data?datasetid=GHCND
93     &datatypeid=AWND&limit=1000&stationid=GHCND:USW00023129
94     &startdate='+year+'-01-01&enddate='+year+'-12-31', headers={'token':Token})
95     #load the api response as a json
96     d = json.loads(r.text)
97     #get all items in the response which are average temperature readings
98     avg_temps = [item for item in d['results'] if item['datatype']=='AWND']
99     #get the date field from all average temperature readings
100    dates_temp += [item['date'] for item in avg_temps]
101    #get the actual average temperature from all average temperature readings
102    temps += [item['value'] for item in avg_temps]
103 #initialize dataframe
104 df_wind = pd.DataFrame()
105 #populate date and average temperature fields (cast string date to datetime and convert temperature from tenths of Celsius to Fahrenheit)
106 df_wind['date'] = [datetime.strptime(d, "%Y-%m-%dT%H:%M:%S") for d in dates_temp]
107 df_wind['wind'] = [float(v)/10.0*1.8 + 32 for v in temps]

```

Figure 13: Average windspeed data scraping

- Precipitation was fetched and stored in a dataframe. Refer Figure 14 for the python code.

```

55 #initialize lists to store data - Precipitation (PRCP)
56 dates_temp = []
57 dates_prcp = []
58 temps = []
59 prcp = []
60 #for each year from 2015-2018 ...
61 for year in range(2015, 2019):
62     year = str(year)
63     print('working on year '+year)
64     #make the api call
65     r = requests.get('https://www.ncdc.noaa.gov/odo-web/api/v2/data?datasetid=GHCND
66     &datatypeid=PRCP&limit=1000&stationid=GHCND:USW00023129
67     &startdate='+year+'-01-01&enddate='+year+'-12-31', headers={'token':Token})
68     #load the api response as a json
69     d = json.loads(r.text)
70     #get all items in the response which are average temperature readings
71     avg_temps = [item for item in d['results'] if item['datatype']=='PRCP']
72     #get the date field from all average temperature readings
73     dates_temp += [item['date'] for item in avg_temps]
74     #get the actual average temperature from all average temperature readings
75     temps += [item['value'] for item in avg_temps]
76 #initialize dataframe
77 df_prcp = pd.DataFrame()
78 #populate date and average temperature fields (cast string date to datetime and convert temperature from tenths of Celsius to Fahrenheit)
79 df_prcp['date'] = [datetime.strptime(d, "%Y-%m-%dT%H:%M:%S") for d in dates_temp]
80 df_prcp['prcp'] = [float(v)/10.0*1.8 + 32 for v in temps]

```

Figure 14: Average precipitation data scraping

- After scraping the required weather data in pandas dataframe, these dataframes were merged to form one dataframe and finally saved to csv. Also, the date parameter was split as Year, Month and Day for merging with crime data.

```

109 #Merge all weather scrapes
110 print(df_temp.shape)
111 print(df_prcp.shape)
112 print(df_wind.shape)
113 df_prcp.head()
114 weather = df_temp.merge(df_prcp,on=['date'])
115 weather = weather.merge(df_wind,on=['date'],how='left')
116 weather.head()
117 #Split date column to year, month and day
118 weather['date'] = pd.to_datetime(weather['date'])
119 weather['Month'] = weather['date'].dt.month
120 weather['Day'] = weather['date'].dt.day
121 weather['Year'] = weather['date'].dt.year
122 weather = weather.drop(columns=['date'])
123 weather.to_csv('C:\Data Analytics\Sem 3\Dataset\LocationData\Distances\weatherapi.csv',index=False)

```

Figure 15: Merge average temperature, precipitation and average wind speed data

3.1.5 Data Merging

All the location co-ordinates present in each dataset in the form of latitude and longitude were converted to geohash. Also, the date parameter was split in to Year, Month and Day in each dataset. Based on the relevant attribute, each dataset was grouped by geohash and merged with crime.

- **Merge Distances**

The nearest distance computed for each crime instance was merged with crime dataset as shown in Figure 16.

```

122 #Convert latitude & longitude to geohash
123 import pygeohash as pgh
124 final['geohash'] = final.apply(lambda x: pgh.encode(x.Latitude, x.Longitude, precision=5), axis=1)
125 dist = final.copy()
126 final.shape
127 #Extract distances
128 distance = dist[['geohash','NearestPoliceDist','NearestSpeedCamDist','NearestRedCamDist']]
129 distance = distance.groupby(['geohash']).mean().reset_index()
130 distance.head()
131 distance = distance.round(2)
132 #Select required columns from crime dataset
133 hs = final[['geohash','Primary_Type','Year','Month','Day','WEEKDAY','Holiday','Time']]
134 hs.shape
135 #compute the crime rate for a geohash a latitude and longitude belongs to on a monthly basis
136 school = hs.groupby(['geohash','Primary_Type','Year','Month','WEEKDAY','Holiday','Time']).size().reset_index()
137 school.head()
138 schoolrate_df = school.rename(columns={0:'crimescount'})
139 schoolrate_df.head()
140 schoolrate_df = schoolrate_df.round(2) #Rounding to nearest place
141 #Merge distance data with crime
142 schoolrate_df = schoolrate_df.merge(distance,on=['geohash'],how='left')
143 schoolrate_df.shape #34440,11
144 schoolrate_df.columns
145 schoolrate_df.isna().sum()
146 schoolrate_df

```

Figure 16: Merge nearest distances with crime data

- **Merge Red-light and speed camera count**

In addition to nearest distance computation, the count of red light and speed cameras in a geohash was done by grouping the location co-ordinates by geohash and merged with crime using a left join as shown in Figure 17. Left join was used because there is a possibility of a crime geohash presence with no cameras.

```

148 #Merge speed cams and red cams in that geohash area
149 scamslocs = pd.read_csv('C:/Data Analytics/Sem 3/Dataset/LocationData/Speed_Camera_Locations_withZip.csv')
150 rcamlocs = pd.read_csv('C:/Data Analytics/Sem 3/Dataset/LocationData/Red_Camera_Locations_withZip.csv')
151 scamslocs.head()
152 rcamlocs.head()
153 #Convert location co-ordinates to geohash
154 scamslocs['geohash'] = scamslocs.apply(lambda x: pgh.encode(x.LATITUDE, x.LONGITUDE, precision=5), axis=1)
155 rcamlocs['geohash'] = rcamlocs.apply(lambda x: pgh.encode(x.LATITUDE, x.LONGITUDE, precision=5), axis=1)
156 #Find the count red light cameras and speed cameras in a geohash
157 slocs = scamslocs.groupby(['geohash']).size().reset_index()
158 rlocs = rcamlocs.groupby(['geohash']).size().reset_index()
159 #Rename columns
160 slocs = slocs.rename(columns={0:'SpeedCamCount'})
161 rlocs = rlocs.rename(columns={0:'RedCamCount'})
162 #Merge with crime data
163 hs = schoolrate_df.merge(slocs,on=['geohash'],how='left')
164 hs = hs.merge(rlocs,on=['geohash'],how='left')
165 hs.fillna(0,inplace = True) #Fill Speedcams and RedCams with 0 in case of no cam locations in that area
166 hs.shape
167 hs.Year.value_counts()
168 hs.head()

```

Figure 17: Merge safety camera counts with crime data

- **Merge high school**

As shown in the Figure 18, pre-processed high school data was merged with crime

data based on the year and geohash attributes.

```
193 #Merge education data
194 edu = pd.read_csv('C:/Data Analytics/Sem 3/Dataset/LocationData/Distances/EducationMerged.csv')
195 edu.head()
196 list(edu.columns)
197 #Calculate average columns from Year1 and Year 2 data for each academic year
198 edu['Avg_Dropout_Rate'] = (edu['One_Year_Dropout_Rate_Year_1_Pct'] + edu['One_Year_Dropout_Rate_Year_2_Pct'])/2
199 edu['Avg_FreshmanTrack_Rate'] = (edu['Freshmen_On_Track_School_Pct_Year_2'] + edu['Freshmen_On_Track_School_Pct_Year_1'])/2
200 edu['Avg_CollegeEnrollment_Rate'] = (edu['College_Enrollment_School_Pct_Year_2'] + edu['College_Enrollment_School_Pct_Year_1'])/2
201 edu['Avg_College_Persistence_Rate'] = (edu['College_Persistence_School_Pct_Year_2'] + edu['College_Persistence_School_Pct_Year_1'])/2
202 #Select only required columns
203 education=edu[['geohash', 'Year', 'Avg_Misconduct_Rate',
204 'Avg_Suspension_Rate',
205 'Avg_Student_Attendance_Rate',
206 'Avg_Teacher_Attendance_Rate',
207 'Avg_Suspension_Days',
208 'SchoolCount',
209 'Avg_Dropout_Rate',
210 'Avg_FreshmanTrack_Rate',
211 'Avg_CollegeEnrollment_Rate',
212 'Avg_College_Persistence_Rate', 'Mobility_Rate_Pct']]
213 #Merge with crime data
214 data = hs.merge(education,on=['geohash', 'Year'],how='left')
215 data.shape
```

Figure 18: Merge high school factors with crime data

- **Merge Weather**

Weather data was merged with crime based on the Year, month and Day attribute as denoted in Figure 19.

```
221 #Merge Weather Data
222 weather = pd.read_csv('C:/Data Analytics/Sem 3/Dataset/LocationData/Distances/weatherapi.csv')
223 weather.head()
224 #Group by year and month
225 weatherm = weather.groupby(['Year', 'Month']).mean().reset_index()
226 #Merge with the crime data
227 alldata = data.merge(weatherm,on=['Year', 'Month'],how='left')
```

Figure 19: Merge weather attributes with crime data

- After merging, there were around 800 missing values which were dropped. The final dataset after merging activity consists of 33565 rows and 28 attributes as highlighted in Figure 20.

```
228 #Check for NA values
229 alldata.isna().sum()
230 alldata.Year.value_counts()
231 alldata=alldata.dropna()
232 alldata.Year.value_counts()
233 alldata.shape #33565,28 After merging all data
```

Figure 20: Merged Data - Check for NA values

- The dataset after merging, can be described as represented in the below Table 1:

Table 1: Crime Prediction Dataset Description

Attribute Code	Description	Domain
crimescount	Count of crime incidents reported	1 - 130
PrimaryType	Type of crime	Assault, Narcotics, Homicide and Violations
Year	Year	2015 - 2018 years
Month	Month	1 - 12 months
geohash	Representation of nearby locations grouped as one area	Alphanumeric value with precision 5
WEEKDAY	Flag indicating whether the day is a weekday or not	0 or 1
Holiday	Flag indicating whether the day is a holiday or not	0 or 1
Time	Time of the day	Morning, Afternoon, Evening or Night
NearestPoliceDist	Distance in kilometers	0.78 - 6.33
NearestRedCamDist	Distance in kilometers	0.51-8.04
RedCamCount	Distance in kilometers	0 - 12
Avg_Student_Attendance_Rate	Attendance rate of student	70% - 96%
Avg_Teacher_Attendance_Rate	Attendance rate of teacher	89% - 95%
Mobility_Rate_Pct	Mobility rate	2% - 37%
SchoolCount	count of schools	1 - 48
wind	average speed of wind in km/hr	34 - 37
prcp	precipitation in mm	32 - 45
avgTemp	temperature in celsius degrees	55 - 77

3.2 Feature Engineering

Effective feature engineering before implementing the models on the data help improve the performance of the models and reduce any possible (Bocca et al.; 2016). These engineering techniques were done in three parts namely, one hot encoding to treat the categorical variables, normalization to treat the numerical features and lastly feature selection to select the best features.

3.2.1 One Hot Encoding

Using the pre-processing library for one hot encoding as well as pandas `get_dummies()` in python as shown in Figure 21, the categorical variables were transformed to binary encoded attributes.

```

330 best_rf_dummies = pd.get_dummies(best_rf)
331 best_rf_dummies.head() #86 columns

```

Figure 21: One hot encoding using `get_dummies()`

Another way to one-hot encode which gives the output in same format as One Hot encoder (OHE) library as shown in Figure 22:

```

42 #One hot encoding for only categorical columns
43 ohe = OneHotEncoder(sparse=False)
44 cat = best.select_dtypes('object')
45 columns_to_encode = cat.columns
46 encoded_columns = ohe.fit_transform(cat[columns_to_encode])

```

Figure 22: One hot encoding using OHE library

Encoded features expressed as binary form (0 and 1) attributes are shown in the below Figure 23:

```

narcotics_dummies = pd.get_dummies(narcotics)
narcotics_dummies.head() #95 columns

```

	geohash_dp3sy	geohash_dp3sz	geohash_dp3t5	geohash_dp3t7	geohash_dp3td	geohash_dp3te	geohash_dp3tf	geohash_dp3tg	geohash_dp3th	geohash_dp3tj	geohash_dp3tk	geohash_dp3tm
7	1	0	0	0	0	0	0	0	0	0	0	0
7	1	0	0	0	0	0	0	0	0	0	0	0
7	1	0	0	0	0	0	0	0	0	0	0	0
7	1	0	0	0	0	0	0	0	0	0	0	0
7	1	0	0	0	0	0	0	0	0	0	0	0

Figure 23: Features after One hot encoding

3.2.2 Normalization

Normalization has been done using the MinMaxScaler library in python as highlighted in the Figure 24.

```

36 #Scaling only on numerical columns
37 numericcols = ['float64', 'int64']
38 numericbest = best.select_dtypes(include=numericcols)
39 from sklearn.preprocessing import StandardScaler, OneHotEncoder
40 sc = MinMaxScaler()
41 NUM = sc.fit_transform(numericbest)

```

Figure 24: Feature scaling using MinMaXScaler library

3.3 Feature Selection

Feature selection has been done using Recursive feature elimination (RFE) combined with random forest which has been followed in machine learning implementations in the past (Granitto et al.; 2006). This technique ranks the features by its importance and elimination is done by RFE as shown in the code using ranks function in Figure 25 and

rfe-rf models as shown in the Figure 26. A function was defined to compute the ranks using RFE.⁷

```
245 from sklearn.preprocessing import MinMaxScaler
246 ranks = {}
247 # Create our function which stores the feature rankings to the ranks dictionary
248 def ranking(ranks, names, order=1):
249     minmax = MinMaxScaler()
250     ranks = minmax.fit_transform(order*np.array([ranks]).T).T[0]
251     ranks = map(lambda x: round(x,2), ranks)
252     return dict(zip(names, ranks))
```

Figure 25: Ranks Function for RFE-RF features

3.3.1 Top 20 features with One Hot Encoding

Including the categorical features encoded as binary values, a round of feature selection was performed and the ranks are listed as shown in the Figure 27

```
254 # Construct our Random Forest Regression model
255 from sklearn.feature_selection import RFE
256 rr = RandomForestRegressor()
257 rr.fit(X,Y)
258 #stop the search when only the last feature is left
259 rfe = RFE(rr, n_features_to_select=20, verbose =3 )
260 rfe.fit(X,Y)
261 ranks["RFE_pub"] = ranking(list(map(float, rfe.ranking_)), colnames, order=-1)
262 print(ranks)
```

Figure 26: Code for fetching top 20 features using RFE-RF

```
ranks
{'RFE_pub': {'Avg_CollegeEnrollment_Rate': 0.93,
'Avg_College_Persistence_Rate': 1.0,
'Avg_Dropout_Rate': 0.91,
'Avg_FreshmanTrack_Rate': 0.92,
'Avg_Misconduct_Rate': 0.89,
'Avg_Student_Attendance_Rate': 1.0,
'Avg_Suspension_Days': 1.0,
'Avg_Suspension_Rate': 1.0,
'Avg_Teacher_Attendance_Rate': 0.99,
'Holiday_False': 0.09,
'Holiday_True': 1.0,
'Mobility_Rate_Pct': 1.0,
'Month_1': 0.76,
'Month_10': 0.68,
'Month_11': 0.73,
'Month_12': 1.0,
'Month_2': 0.88,
'Month_3': 0.8,
```

Figure 27: Ranks of top 20 features

3.3.2 Top 10 numerical features

With RFE-RF (Recursive feature elimination-random forest) method, top 10 numerical features were extracted and their ranks are listed in Figure 29 and for the code, refer

⁷<https://www.kaggle.com/arthurtok/feature-ranking-rfe-random-forest-linear-models>

Figure 28.

```
314 # Construct our Random Forest Regression model
315 from sklearn.feature_selection import RFE
316 rr = RandomForestRegressor()
317 rr.fit(X,Y)
318 #stop the search when only the last feature is left
319 rfe = RFE(rr, n_features_to_select=10, verbose =3 )
320 rfe.fit(X,Y)
321 ranks["RFE_pub"] = ranking(list(map(float, rfe.ranking_)), colnames, order=-1)
322 ranks
```

Figure 28: Code for fetching top 10 features using RFE-RF

```
ranks
{'RFE_pub': {'Avg_CollegeEnrollment_Rate': 0.67,
'Avg_College_Persistence_Rate': 0.44,
'Avg_Dropout_Rate': 0.11,
'Avg_FreshmanTrack_Rate': 0.56,
'Avg_Misconduct_Rate': 0.89,
'Avg_Student_Attendance_Rate': 1.0,
'Avg_Suspension_Days': 0.78,
'Avg_Suspension_Rate': 0.33,
'Avg_Teacher_Attendance_Rate': 1.0,
'Mobility_Rate_Pct': 1.0,
'NearestPoliceDist': 1.0,
'NearestRedCamDist': 1.0,
'NearestSpeedCamDist': 0.22,
'RedCamCount': 1.0,
'SchoolCount': 1.0,
'SpeedCamCount': 0.0,
'avgTemp': 1.0,
'prcp': 1.0,
'wind': 1.0}}
```

Figure 29: Ranks of top 10 features

3.4 Modelling

Modelling was done using the python scikit libraries for machine learning. XGBoost regressor, random forest regressor, keras, tensorflow and linear regression libraries were used for modelling.

3.4.1 Data Split

The models have been tested for both the versions of the split i.e. with train-test (80:20) as well as cross validation (k folds:3-30) as highlighted in the Figure 30. Cross validation techniques have enabled efficient sampling of data for the models and eventually generating better results (Ingilevich and Ivanov; 2018) and Kadar (2019).

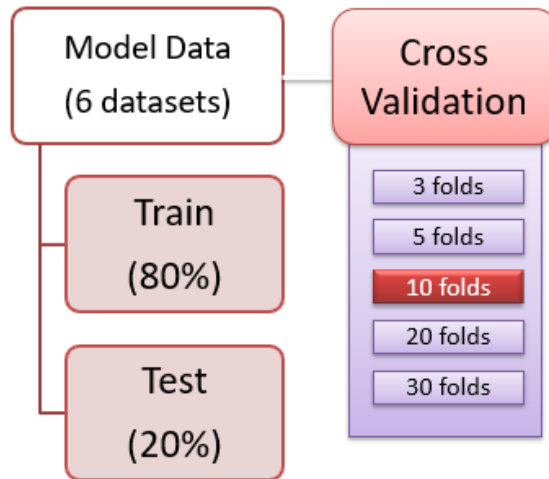


Figure 30: Split of Data

Train test split and k fold libraries were from scikit learn model selection for this task, refer Figure 31.

```

330 best_rf_dummies = pd.get_dummies(best_rf)
331 best_rf_dummies.head() #86 columns
332 #list(best_rf_dummies.columns) (Big output cell)
333 model_data = best_rf_dummies.values
334 #separate target and independent variables
335 X = model_data[:,1:]
336 Y = model_data[:,0]
337 #Set seed
338 np.random.seed(80)
339 #Split train and test data in 80:20 ratio
340 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.20, random_state=20)

```

Figure 31: Code to split data

3.4.2 Random Forest

Random forest was applied with default parameters, optimized parameters and finally optimized parameters with 10 fold cross validation.

- Refer the code shown in Figure 32 for random forest with default parameters

```

354 clf = RandomForestRegressor(n_estimators=10, random_state=20, n_jobs=-1)
355 # Train the classifier
356 clf.fit(X_train, Y_train)
357 #Training model accuracy
358 trainac = clf.predict(X_train)
359 print("Train accuracy details of Random Forest")
360 print("RMSE is",np.sqrt(mean_squared_error(Y_train,trainac)))
361 print("R2 ",r2_score(Y_train, trainac))
362 print("MAE ", mean_absolute_error(Y_train, trainac))
363 print("MSE ",mean_squared_error(Y_train,trainac))
364 #Testing model accuracy
365 y_pred = clf.predict(X_test)
366 # View The Accuracy Of best features (20 Features) Model
367 print(clf.score(X_test,Y_test))
368 print('Mean Absolute Error:', metrics.mean_absolute_error(Y_test, y_pred))
369 print('Mean Squared Error:', metrics.mean_squared_error(Y_test, y_pred))
370 print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(Y_test, y_pred)))

```

Figure 32: Random Forest with default parameters

- Refer the code shown in Figure 33 for random forest with optimized parameters

```

653 tuned_model = RandomForestRegressor(bootstrap= True,
654     max_depth= 50,
655     max_features= 'auto',
656     min_samples_leaf= 4,
657     min_samples_split= 5,
658     n_estimators= 377, random_state = 20,n_jobs=-1)
659 #Fit the tuned model
660 tuned_model.fit(X_train, Y_train)
661 trainac = tuned_model.predict(X_train)
662 #Train accuracy
663 print("Train accuracy details on the tuned of Random Forest")
664 print("RMSE is",np.sqrt(mean_squared_error(Y_train,trainac)))
665 print("R2 ",r2_score(Y_train, trainac))
666 print("MAE ", mean_absolute_error(Y_train, trainac))
667 print("MSE ",mean_squared_error(Y_train,trainac))
668 #Predict crime count
669 y_pred = tuned_model.predict(X_test)
670 # View The Accuracy Of best features (86 Features) Model with tuned parameters
671 tuned_model.score(X_test,Y_test)
672 print("RMSE is",np.sqrt(mean_squared_error(Y_test,y_pred)))
673 print("R2 ",r2_score(Y_test, y_pred))
674 print("MAE ", mean_absolute_error(Y_test, y_pred))
675 print("MSE ",mean_squared_error(Y_test,y_pred))

```

Figure 33: Random Forest - optimized model fit

- Random Forest with 10 fold cross validation

```

426 #Random Forest
427 kfold = model_selection.KFold(n_splits=10, random_state=200,shuffle=True)
428 model_kfoldrand = RandomForestRegressor(n_estimators=10, random_state=20, n_jobs=-1)
429 results_kfoldrand = model_selection.cross_val_score(model_kfoldrand, X, Y, cv=kfold)
430 print("Accuracy: %.2f%%" % (results_kfoldrand.mean()*100.0))
431 for train_index, test_index in kfold.split(X):
432     print("TRAIN:", train_index, "TEST:", test_index)
433     X_trainkf, X_testkf = X[train_index], X[test_index]
434     y_trainkf, y_testkf = Y[train_index], Y[test_index]
435     clf = model_kfoldrand.fit(X_trainkf, y_trainkf)
436     print("Residual sum of squares: %.2f"
437           % np.mean((model_kfoldrand.predict(X_testkf) - y_testkf) ** 2))
438     #Explained variance score: 1 is perfect prediction
439     print('Explained Variance score: %.2f' % model_kfoldrand.score(X_testkf, y_testkf))
440     kfoldrf = model_kfoldrand.predict(X_testkf)
441     print("RMSE is",np.sqrt(mean_squared_error(y_testkf,kfoldrf)))
442     print("R2 ",r2_score(y_testkf, kfoldrf))
443     print("MAE ", mean_absolute_error(y_testkf, kfoldrf))
444     print("MSE ",mean_squared_error(y_testkf,kfoldrf))

```

Figure 34: Random Forest - 10 fold cross validation

3.4.3 XGBoost

XGBoost was also applied with default parameters, optimized parameters and finally optimized parameters with 10 fold cross validation as shown in the codes in Figure 35, Figure 36 and Figure 37.

- XGBoost with default parameters

```

373 #XGBoost
374 xgb = xgbo.XGBRegressor(n_estimators=100, learning_rate=0.1)
375 xgb.fit(X_train,Y_train)
376 predictions = xgb.predict(X_test)
377 print("Variance Score is", explained_variance_score(predictions,Y_test))
378 print("R2 ",r2_score(Y_test, predictions))
379 print("MAE ", mean_absolute_error(Y_test, predictions))
380 print("MSE ",mean_squared_error(Y_test,predictions))
381 print("RMSE is",np.sqrt(mean_squared_error(Y_test,predictions)))
382 trainac = xgb.predict(X_train)
383 print("Train accuracy details on the tuned - XGBoost")
384 print("RMSE is",np.sqrt(mean_squared_error(Y_train,trainac)))
385 print("R2 ",r2_score(Y_train, trainac))
386 print("MAE ", mean_absolute_error(Y_train, trainac))
387 print("MSE ",mean_squared_error(Y_train,trainac))

```

Figure 35: XGBoost - default parameters

- XGBoost with Optimized parameters

```

753 #Tuned model for XGB
754 tuned_modelxgb = xgbo.XGBRegressor(bootstrap= True,
755     colsample_bytree= 1.0,
756     gamma= 1.5,
757     learning_rate= 0.05,
758     max_depth= 8,
759     min_child_weight= 10,
760     n_estimators = 200,
761     subsample= 0.75, random_state = 20,n_jobs=-1)
762 tuned_modelxgb.fit(X_train, Y_train)
763 trainac = tuned_modelxgb.predict(X_train)
764 print("Train accuracy details on the tuned - XGBoost")
765 print("RMSE is",np.sqrt(mean_squared_error(Y_train,trainac)))
766 print("R2 ",r2_score(Y_train, trainac))
767 print("MAE ", mean_absolute_error(Y_train, trainac))
768 print("MSE ",mean_squared_error(Y_train,trainac))
769 y_pred = tuned_modelxgb.predict(X_test)
770 # View The Accuracy Of best features (86 Features) Model with tuned parameters
771 tuned_modelxgb.score(X_test,Y_test)
772 print("RMSE is",np.sqrt(mean_squared_error(Y_test,y_pred)))
773 print("R2 ",r2_score(Y_test, y_pred))
774 print("MAE ", mean_absolute_error(Y_test, y_pred))
775 print("MSE ",mean_squared_error(Y_test,y_pred))

```

Figure 36: XGBoost with optimized parameters

- XGBoost with tuning and 10 fold cross validation

```

781 #kfold on xgboost tuned
782 from sklearn.model_selection import cross_val_score
783 kfold = model_selection.KFold(n_splits=10, random_state=200,shuffle=True)
784 model_kfoldxgb = xgbo.XGBRegressor(bootstrap= True,
785     colsample_bytree= 1.0,
786     gamma= 1.5,
787     learning_rate= 0.05,
788     max_depth= 8,
789     min_child_weight= 10,
790     n_estimators = 200,
791     subsample= 0.75, random_state = 20,n_jobs=-1)
792 results_kfoldxgb = cross_val_score(model_kfoldxgb, X, Y, cv=kfold)
793 print("Accuracy: %.2f%%" % (results_kfoldxgb.mean()*100.0))
794 for train_index, test_index in kfold.split(X):
795     print("TRAIN:", train_index, "TEST:", test_index)
796     X_trainkf, X_testkf = X[train_index], X[test_index]
797     y_trainkf, y_testkf = Y[train_index], Y[test_index]
798     clf_xgb = model_kfoldxgb.fit(X_trainkf, y_trainkf)
799     print("Residual sum of squares: %.2f"
800         # % np.mean((model_kfoldxgb.predict(X_testkf) - y_testkf) ** 2))
801         #Explained variance score: 1 is perfect prediction
802         print('Variance score: %.2f' % model_kfoldxgb.score(X_testkf, y_testkf))
803         kfoldxgb = model_kfoldxgb.predict(X_testkf)
804         print("RMSE is",np.sqrt(mean_squared_error(y_testkf,kfoldxgb)))
805         print("R2 ",r2_score(y_testkf, kfoldxgb))
806         print("MAE ", mean_absolute_error(y_testkf, kfoldxgb))
807         print("MSE ",mean_squared_error(y_testkf,kfoldxgb))
808         #Train accuracy
809         trainac = model_kfoldxgb.predict(X_trainkf)
810         print("RMSE is",np.sqrt(mean_squared_error(y_trainkf,trainac)))
811         print("R2 ",r2_score(y_trainkf, trainac))
812         print("MAE ", mean_absolute_error(y_trainkf, trainac))
813         print("MSE ",mean_squared_error(y_trainkf,trainac))

```

Figure 37: XGBoost with optimized parameters and 10 fold cross validation

3.4.4 Artificial Neural Network (ANN)

- Refer the code in Figure 38 for implementing ANN without hidden layer

```

92 #Define ann layers
93 model = Sequential()
94 model.add(Dense(128, input_dim=X_train.shape[1], kernel_initializer='normal', activation='relu'))
95 model.add(Dense(1, kernel_initializer='normal'))
96 model.compile(loss='mse', optimizer='rmsprop', metrics=['mse'])
97 model.summary()
98 #fit the ann model
99 history = model.fit(X_train, y_train, batch_size=64, epochs=50, verbose=2, validation_split=.2)
100 #plot loss vs. epoch curve
101 plt.figure(figsize=(10,5))
102 plt.plot(history.history['loss'],marker='o',color='orange')
103 plt.plot(history.history['val_loss'],marker='^',color='blue')
104 plt.title('Value Loss')
105 plt.ylabel('loss')
106 plt.xlabel('epoch')
107 plt.legend(['train', 'test'], loc='upper right')
108 plt.show()
109 #plot the train vs test mse
110 plt.figure(figsize=(10,5))
111 plt.plot(history.history['mean_squared_error'],marker='o',color='red')
112 plt.plot(history.history['val_mean_squared_error'],marker='^',color = 'green')
113 plt.title('Value Mean Squared Error')
114 plt.ylabel('MSE')
115 plt.xlabel('epoch')
116 plt.legend(['train', 'test'], loc='upper right')
117 plt.show()
118 #predict test data with ann
119 pred_crimes = model.predict(X_test)
120 mse_pred_score = metrics.mean_squared_error(pred_crimes, y_test)
121 print('mse_pred_score {}'.format(mse_pred_score))
122 rmse_pred_score = np.sqrt(mse_pred_score)
123 print('rmse_pred_score {}'.format(rmse_pred_score))
124 r2_pred_score = r2_score(y_test, pred_crimes, multioutput='uniform_average')
125 print('r2_pred_score - Coefficient of Determination {}'.format(r2_pred_score))
126 print("MAE ", mean_absolute_error(y_test, pred_crimes))

```

Figure 38: ANN model with one layer

- Refer the code in Figure 39 for implementing ANN with multiple layers


```

143 model = Sequential()
144 model.add(Dense(128, input_dim=X_train.shape[1], kernel_initializer='normal', activation='relu'))
145 model.add(Dense(64, kernel_initializer='he_uniform', activation='relu'))
146 model.add(Dense(32, kernel_initializer='he_uniform', activation='relu'))
147 model.add(Dense(1, kernel_initializer='normal'))
148 model.compile(loss='mse', optimizer='rmsprop', metrics=['mse'])
149 model.summary()
150 #Set seed and fit the model
151 np.random.seed(80)
152 history = model.fit(X_train, y_train, batch_size=128, epochs=50, validation_split=.2, verbose=2)
153 plt.figure(figsize=(10,5))#Plot loss vs. epoch
154 plt.plot(history.history['loss'],marker='o',color='orange')
155 plt.plot(history.history['val_loss'],marker='^',color='blue')
156 plt.title('Value Loss')
157 plt.ylabel('loss')
158 plt.xlabel('epoch')
159 plt.legend(['train', 'test'], loc='upper right')
160 plt.show()
161 #plot the train vs test mse
162 plt.figure(figsize=(10,5))
163 plt.plot(history.history['mean_squared_error'],marker='o',color='red')
164 plt.plot(history.history['val_mean_squared_error'],marker='^',color = 'green')
165 plt.title('Value Mean Squared Error')
166 plt.ylabel('MSE')
167 plt.xlabel('epoch')
168 plt.legend(['train', 'test'], loc='upper right')
169 plt.show()
170 pred_crimes = model.predict(X_test) #Predict crime count
171 mse_pred_score = metrics.mean_squared_error(pred_crimes, y_test)
172 print('mse_pred_score {}'.format(mse_pred_score))
173 rmse_pred_score = np.sqrt(mse_pred_score)
174 print('rmse_pred_score {}'.format(rmse_pred_score))
175 r2_pred_score = r2_score(y_test, pred_crimes, multioutput='uniform_average')
176 print('r2_pred_score - Coefficient of Determination {}'.format(r2_pred_score))
177 print("MAE ", mean_absolute_error(y_test, pred_crimes))

```

Figure 39: ANN model with multiple layers

3.4.5 Multiple Linear Regression

- For multiple linear regression, tuning is not applicable and hence 10 fold cross validation is applied and checked as shown in the code block Figure 40.

```

54 X = final.drop(columns='crimescount')
55 y = final.crimescount
56 # Splitting the dataset into the Training set and Test set
57 from sklearn.model_selection import train_test_split
58 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 20)
59 # Fitting Multiple Linear Regression to the Training set
60 from sklearn.linear_model import LinearRegression
61 from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
62 regressor = LinearRegression()
63 regressor.fit(X_train, y_train)
64 # Predicting the Test set results
65 y_pred = regressor.predict(X_test)
66 from sklearn.metrics import r2_score
67 #score=r2_score(y_test,y_pred)
68 print("R2 ",r2_score(y_test,y_pred))
69 print("MAE ", mean_absolute_error(y_test,y_pred))
70 print("MSE ",mean_squared_error(y_test,y_pred))
71 print("RMSE is",np.sqrt(mean_squared_error(y_test,y_pred)))
72 y_pred
73 residual = y_test-y_pred
74 residual
75 fig, ax = plt.subplots(figsize=(10,5))
76 _ = ax.scatter(residual, y_pred)
77 # Set common labels
78 ax.tick_params(axis="x", labelszize=15)
79 ax.tick_params(axis="y", labelszize=15)
80 ax.set_xlabel('Residuals', fontsize=15)
81 ax.set_ylabel('Predictions', fontsize=15)
82 #The residual vs predictions is clumped and the behaviour is not random, thus homoscedasticity assumption not satis
83 import scipy as sp
84 fig, ax = plt.subplots(figsize=(10,5))
85 _, (__, ___, r) = sp.stats.probplot(residual, plot=ax, fit=True)
86 ax.tick_params(axis="x", labelszize=15)
87 ax.tick_params(axis="y", labelszize=15)
88 ax.set_title('Normal Probability Plot for Errors', fontsize=16)

```

Figure 40: Multiple Linear regression code

3.4.6 Hyper Parameter Optimization

Tuning has been done using Randomized search cv for each model as shown in the below Figure 41 and Figure 42

- Random Forest Tuning

```

597 #Hyper parameters for Random Forest
598 from sklearn.model_selection import RandomizedSearchCV
599 # Number of trees in random forest
600 n_estimators = [int(x) for x in np.linspace(start = 10, stop = 500, num = 5)]
601 # Number of features to consider at every split
602 max_features = ['auto', 'sqrt']
603 # Maximum number of levels in tree
604 max_depth = [int(x) for x in np.linspace(10, 50, num = 5)]
605 max_depth.append(None)
606 # Minimum number of samples required to split a node
607 min_samples_split = [2, 5, 10]
608 # Minimum number of samples required at each leaf node
609 min_samples_leaf = [1, 2, 4]
610 # Method of selecting samples for training each tree
611 bootstrap = [True, False]
612 # Create the random grid
613 random_grid = {'n_estimators': n_estimators,
614               'max_features': max_features,
615               'max_depth': max_depth,
616               'min_samples_split': min_samples_split,
617               'min_samples_leaf': min_samples_leaf,
618               'bootstrap': bootstrap}
619 print(random_grid)
620 # Use the random grid to search for best hyperparameters
621 # First create the base model to tune
622 import datetime
623 print(datetime.datetime.now())
624 rf = RandomForestRegressor()
625 # Random search of parameters, using 10 fold cross validation,
626 # search across 1000 different combinations, and use all available cores
627 rf_random = RandomizedSearchCV(estimator = rf, param_distributions = random_grid, n_iter = 100, cv = 3, verbose=2, random_state=20, n_jobs = -1)
628 # Fit the random search model
629 rf_random.fit(X_train, Y_train)

```

Figure 41: Code for tuning Random Forest using Randomized search cv

- XGBoost Tuning

```

714 # A parameter grid for XGBoost
715 params = {
716     'learning_rate': [0.02, 0.05, 0.08, 0.10],
717     'n_estimators': [20, 50, 100, 150, 200],
718     'min_child_weight': [1, 5, 10],
719     'gamma': [0, 0.5, 1, 1.5, 2],
720     'subsample': [0.6, 0.75, 0.8, 1.0],
721     'colsample_bytree': [0.6, 0.8, 1.0],
722     'max_depth': [3, 4, 5, 7, 8]
723 }
724 xgb = xgbo.XGBRegressor()
725 xgbkf = model_selection.KFold(n_splits=10, random_state=200, shuffle=True)
726 random_search = RandomizedSearchCV(xgb, param_distributions=params, n_iter=100, n_jobs=-1,
727 cv=xgbkf.split(X, Y), verbose=3, random_state=20)
728 print(datetime.datetime.now())
729 random_search.fit(X, Y)
730 print(datetime.datetime.now())

```

Figure 42: Code for tuning XGBoost using Randomized search cv

References

Bocca, F. F., Henrique, L. and Rodrigues, A. (2016). The effect of tuning , feature engineering , and feature selection in data mining applied to rainfed sugarcane yield modelling, *Computers and Electronics in Agriculture* **128**: 67–76.

URL: <http://dx.doi.org/10.1016/j.compag.2016.08.015>

Granitto, P. M., Furlanello, C., Biasioli, F. and Gasperi, F. (2006). Recursive feature elimination with random forest for PTR-MS analysis of agroindustrial products, **83**: 83–90.

Ingilevich, V. and Ivanov, S. (2018). Crime rate prediction in the urban environment using social factors, *Procedia Computer Science* **136**: 472–478.

URL: <https://doi.org/10.1016/j.procs.2018.08.261>

Kadar, C. (2019). Public decision support for low population density areas : An imbalance-aware hyper-ensemble for spatio-temporal crime prediction, *Decision Support Systems* **119**(September 2018): 107–117.

URL: <https://doi.org/10.1016/j.dss.2019.03.001>