

Introduction:

This Assignment focus on Data Integration and Data reshaping. Where we are given 2 main tasks to work on. In Task 1 we are provided with several datasets in various formats like:

- JSON
- XML
- PDF
- EXCEL
- Shape Files
- Text Files Now we have to load all these files into python and convert them to pandas data frame:

In Task 2 we have to understand the effect of various normalisation/transformation methods i.e. standardization, MinMax normalisation, log, power box-cox transformation on various columns which are price, Distance_to_sc, travel_min_to_CBD, and Distance_to_hospital and explain the effect. To develop a linear model to predict the price using Distance_to_sc, travel_min_to_CBD, and Distance_to_hospital.

TASK 1: Data Integration:

In this task we have 7 files that contains the information regarding the housing in Victoria Each of these file contains a different data set about housing information in Victoria, Australia. Where in

- JSON and XML file we have the information about the house like property_id, lat, lng, addr_street, suburb, price, property_type, year, bedrooms, bathrooms and parking_space .
- PDF Contains the information about the hospitals in victoria id , lat , lng , name
- EXCEL : Information regards the supermarket in victoria id , lat , lng , type
- Shape File : Contains the information of suburb and the polygon of lat and lng in a suburb area.
- Text File : Contains the train route, time , station, distance and various information of PTV in Victoria.

All the distance Calculated are converted into metres

Importing Libraries:

```
In [ ]: 1 import json
2 import pandas as pd
3 #!/pip install tabula
4 #!/pip install tabula-py
5 import tabula
6 from tabula import wrapper
7 from tabula import read_pdf
8 import xml.etree.cElementTree as et
9 #!/pip install BeautifulSoup4
10 import os
11 import sys
12 import pandas as pd
13 from bs4 import BeautifulSoup
14 from math import sin, cos, sqrt, atan2, radians
15 #!/pip install pyshp
16 import numpy as np
17 import shapefile
18 from shapely.geometry import Point
19 from shapely.geometry.polygon import Polygon
```

Reading JSON FILE:

We have read the JSON file and stored it into a variable after that we have converted it to dataframe

```
In [ ]: 1 # Opening JSON file
2 f = open('real_state.json',)
3 # returns JSON object as
4 # a dictionary
5 data = json.load(f)
6 # Closing file`
7 f.close()
8
9 real_stateJS = pd.DataFrame(data)
```

Reading XML

- Read the XML file
- Open XML file and remove the unwanted character and Re-write the XML file and re read it again
- parsing over the XML and creating a dictionary for each row and then converted it into a dataframe

```
In [ ]: 1 f=open('real_state.xml', 'r')
2 data1 = f.read()
3 f.close()
4 data1=data1[2:]
5 data1=data1[0:-1]
6 finalxml = open('real_states.xml', "w",encoding='UTF-8')
7 finalxml.write(data1)
8 finalxml.close()
9
10
```

```
In [ ]: 1 parsed_xml = et.parse("real_states.xml")
2 xml_dict = {}
3 # parsing in xml and getting dict for each row
4 for i in parsed_xml.getroot():
5     xml_list = []
6     for j in i:
7         xml_list.append(j.text)
8     xml_dict[i.tag] = xml_list
9
10 # converting dictionary to dataframe
11 real_stateXML = pd.DataFrame(xml_dict)
```

Merging JSON and XML:

In this part I have merged the data for JSON and XML as both the files contain data related to the house address and house's. along with that we have re-indexed the index and removed the duplicate rows

```
In [ ]: 1 # concating both the dataframe to one final dataframe
2 real_state=pd.concat([real_stateJS,real_stateXML], axis=0,ignore_index=True)
3
4 # Lat Long data type change
5 real_state[['lat', 'lng']] = real_state[['lat', 'lng']].astype(float)
6
7 # property id datatype change
8 real_state['property_id']=real_state['property_id'].astype(int)
9
10 # removing duplicate
11 real_state=real_state.drop_duplicates('property_id')
12 # reseting index
13 real_state=real_state.reset_index(drop=True)
```

Creating required column & setting Default Value

- In this part I have created the required column in the main dataframe
- Set their default value and converted it into the dataframe
- Re-arranged the columns based on the assignment specification and created a final structure for the dataframe

```
In [ ]: 1 real_state['suburb']='not available'
2 real_state['Shopping_center_id']='not available'
3 real_state['Distance_to_sc']=0
4 real_state['Train_station_id']=0
5 real_state['Distance_to_train_station']=0
6 real_state['travel_min_to_CBD']=0
7 real_state['Transfer_flag']=-1
8 real_state['Hospital_id']='not available'
9 real_state['Distance_to_hospital']=0
10 real_state['Supermarket_id']='not available'
11 real_state['Distance_to_supermaket']=0
12 columnsTitles = ['property_id','lat','lng','addr_street','suburb','price','p
13                   'bathrooms','parking_space','Shopping_center_id','Distance_
14                   'Distance_to_train_station','travel_min_to_CBD','Transfer_f
15                   'Distance_to_hospital','Supermarket_id','Distance_to_superm
16
17 real_state= real_state.reindex(columns=columnsTitles)
```

Reading PDF

Reading all the data in the PDF page by page and mearging all the dataframe to one dataframe.

```
In [ ]: 1 #Convert your file
2 df1 = tabula.read_pdf("hospitals.pdf",pages=1)
3 df2 = tabula.read_pdf("hospitals.pdf",pages=2)
4 df3 = tabula.read_pdf("hospitals.pdf",pages=3)
5 df4 = tabula.read_pdf("hospitals.pdf",pages=4)
6 df5 = tabula.read_pdf("hospitals.pdf",pages=5)
7 hospital=pd.concat([df1,df2,df3,df4,df5],axis=0)
8 hospital.dropna(subset = ["id"], inplace=True)
9 hospital=hospital.reset_index()
10
```

Reading HTML

Reading the HTML file into python and converting it into a dataframe

```

In [ ]: 1 path = 'shoppingcenters.html'
        2 # empty list
        3 data = []
        4 # for getting the header from
        5 # the HTML file
        6 list_header = []
        7 soup = BeautifulSoup(open(path), 'html.parser')
        8 header = soup.find_all("table")[0].find("tr")
        9
        10 for items in header:
        11     try:
        12         list_header.append(items.get_text())
        13     except:
        14         continue
        15
        16 # for getting the dataa
        17 HTML_data = soup.find_all("table")[0].find_all("tr")[1:]
        18
        19 for element in HTML_data:
        20     sub_data = []
        21     for sub_element in element:
        22         try:
        23             sub_data.append(sub_element.get_text())
        24         except:
        25             continue
        26     data.append(sub_data)
        27
        28
        29
        30 # Storing the data into pandas DataFrame
        31 shoppingCentre = pd.DataFrame(data = data, columns = list_header)

```

Reading EXCEL:

reading excel file for the supermarket data and loading it into a dataframe

```

In [ ]: 1 superMarket=pd.read_excel("supermarkets.xlsx")

```

Reading Train Information Text File:

Reading the stop file to find out the nearest staion and station id.

```

In [ ]: 1 #stops= pd.read_csv("stops.txt", sep = ',')
        2 df_trips=pd.read_csv("trips.txt")
        3 df_stops=pd.read_csv("stops.txt")
        4 df_routes=pd.read_csv("routes.txt")
        5 df_stop_times=pd.read_csv("stop_times.txt")
        6 df_calender=pd.read_csv("calendar.txt")
        7

```

Reading shape File

Reading the shape files for the suburb boundry polygon and fetching the records for suburb and from the list of polygon's

```
In [ ]: 1 sf = shapefile.Reader("./VIC_LOCALITY_POLYGON_shp") # note, no suffix, all 3
        2 recs = sf.records()
        3 shapes = sf.shapes()
```

Finding suburb's: and Suburb id:

In this part we have looped over the dataframe and find out the each records lat long and stored it then looped over the shape file polygon and found the record for required lat long and append the suburb name for that

```
In [ ]: 1 # Finding the suburb using the shape file records and shapes we have created
        2 # Then checking the points i.e. lat and long for property id in that polygon
        3 #index and returning the suburb for that same index in records file
        4 for i in range(len(real_state)):
        5     pts=Point(real_state.lng[i],real_state.lat[i],)
        6     for j in range(len(shapes)):
        7         pointList=shapes[j].points
        8         polygon = Polygon(pointList)
        9         if polygon.contains(pts):
        10             real_state.at[i, 'suburb']=recs[j][6]
        11
```

```
In [ ]: 1 # checking if all the values are replaced with the suburb
        2 real_state[real_state['suburb']=='not available']
```

Finding nearest Shopping Centre and Shopping Centre ID

- Using haversine distance formula we have found out the distance for the nearest shopping centre and shopping id .
- then looping over the property id fetching lat long for each record and then calculating distance with each shopping centre
- Finally finding the shopping centre id and shopping centre distance with least distance.

```
In [ ]: 1 def distance_to_sc(x1,y1,sc_id):
2       Radius = 6378.0
3       x2 = radians(shoppingCentre['lat'][shoppingCentre['sc_id'] == sc_id])
4       y2 = radians(shoppingCentre['lng'][shoppingCentre['sc_id'] == sc_id])
5       dist_long = y2 - y1
6       dist_lat = x2 - x1
7       mat = sin(dist_lat / 2)**2 + cos(x1) * cos(x2) * sin(dist_long / 2)**2
8       mat1 = 2 * atan2(sqrt(mat), sqrt(1 - mat))
9       distance = Radius * mat1
10      return distance, sc_id
11
```

```
In [ ]: 1 # Iterating over dataframe and finding the nearest Shopping Centre and
2 # and appending the distance from the property along with Shopping centre id
3 for i in range(len(real_state['property_id'])):
4     final_list = []
5     dis_list = []
6     name_list = []
7     for j in shoppingCentre['sc_id']:
8         dis,name = distance_to_sc(radians(real_state.lat[i])\
9                                   ,radians(real_state.lng[i]),j)
10        dis_list.append(dis)
11        name_list.append(name)
12    final_list = list(zip(dis_list,name_list))
13    final_dis,final_name = min(final_list)[0],min(final_list)[1]
14    real_state.at[i,'Distance_to_sc'] = (round(final_dis,4)*1000)
15    real_state.at[i,'Shopping_center_id'] = final_name
```

```
In [ ]: 1 # checking if all the values are replaced with the Shopping Centre id
2 #shopping centre distance is calculated
3 real_state[real_state['Shopping_center_id']=='not available']
```

```
In [ ]: 1 real_state[real_state['Distance_to_sc']==0]
```

Finding nearest Train station

Using haversine distance we have created a function to calculate the distance to the nearest train station with each record in the dataframe

```
In [ ]: 1 # Function to calculate the distance
2 def distance_to_station(x1,y1,stop_id):
3     Radius = 6378.0
4     x2 = radians(df_stops['stop_lat'][df_stops['stop_id'] == stop_id])
5     y2 = radians(df_stops['stop_lon'][df_stops['stop_id'] == stop_id])
6     dist_long = y2 - y1
7     dist_lat = x2 - x1
8     mat = sin(dist_lat / 2)**2 + cos(x1) * cos(x2) * sin(dist_long / 2)**2
9     mat1 = 2 * atan2(sqrt(mat), sqrt(1 - mat))
10    distance = Radius * mat1
11    return distance, stop_id
12
```

```

In [ ]: 1 # Iterating over dataframe and finding the nearest station
        2 # and appending the distance from the property along with station id to each
        3 for i in range(len(real_state['property_id'])):
        4     final_list = []
        5     dis_list = []
        6     name_list = []
        7     for j in df_stops['stop_id']:
        8         dis,name = distance_to_station(radians(real_state.lat[i])\
        9                                     ,radians(real_state.lng[i]),j)
        10        dis_list.append(dis)
        11        name_list.append(name)
        12        final_list = list(zip(dis_list,name_list))
        13        final_dis,final_name = min(final_list)[0],min(final_list)[1]
        14        real_state.at[i,'Distance_to_train_station'] = (round(final_dis,4)*1000)
        15        real_state.at[i,'Train_station_id'] = final_name

```

```

In [ ]: 1 # checking if all the values are replaced with the Station ID
        2 # distance to station is calculated
        3 real_state[real_state['Train_station_id']==0]
        4

```

```

In [ ]: 1 real_state[real_state['Distance_to_train_station']==0]

```

Minimum Average Travel Time to CBD and Transfer Flag

- In this part we have to calculate the average travel time to CBD i.e Flinder Street Station from each property record using the nearest train Station for weekdays i.e. Monday to Friday where the trains departures between 7AM-9AM. if there are 3 trip
- If there are any direct transfers between the closest station and Flinders street station, only the average of direct transfers should be calculated.
- We have created a column that holds a Boolean value which indicate's whether there is a direct trip to the Flinders street station from the closest station between 7-9am on the weekdays
- It is set to 0 if there is a direct trip i.e. no transfer between trains is required to get from the closest train station to the Flinders station
- And it is set to 1 when 1 there is transfer between trains is required to get from the closest station to flinders streen between 7-9AM on week day
- Default value remains -1


```
In [ ]: 1 import pandas as pd
2 from datetime import datetime
3 df_combined=df_trips.join(df_routes.set_index('route_id'), on='route_id')\
4     .join(df_stop_times.set_index('trip_id'), on='trip_id')\
5     .join(df_calender.set_index('service_id'), on='service_id')
6
7
8 start = datetime.strptime('07:00:00','%H:%M:%S')
9 end= datetime.strptime('11:00:00','%H:%M:%S')
10 df_combined=df_combined[(df_combined["monday"]==1) & (df_combined["tuesday"]
11     & (df_combined["thursday"]==1) & (df_combined["frida
12     &(df_combined['trip_headsign']=='City (Flinders Stre
```

```
In [ ]: 1 acceptedStart = datetime.strptime('07:00:00','%H:%M:%S')
2 acceptedEnd= datetime.strptime('09:00:00','%H:%M:%S')
3 lst_trip=df_combined.trip_id.unique()
4
5
6 for index, row in df_combined.iterrows():
7     try:
8         at=datetime.strptime(row['arrival_time'],'%H:%M:%S')
9         if((at>=acceptedStart )& (at<=acceptedEnd)):
10             pass
11         else:
12             df_combined.drop(index, inplace=True)
13
14     except:
15         pass
```

```
In [ ]: 1 len(df_combined)
```

In []:

```

1  def Average(lst):
2      if(len(lst)==0):
3          return 0
4      else:
5          return sum(lst) / len(lst)
6
7
8
9  def findAverage(startStopid,endStopId):
10     averageTime=[]
11     flag=1
12     if(endStopId==startStopid==19854):
13         return (0,0)
14     for eachTrip in lst_trip:
15         trip=df_combined[df_combined["trip_id"]==eachTrip]
16         stops=trip["stop_id"].tolist()
17
18         if((startStopid in stops ) & (endStopId in stops)):
19             index_start=stops.index(startStopid)
20             index_end=stops.index(endStopId)
21             if(index_start<index_end):
22                 try:
23                     time_dpt_start=datetime.strptime(trip[trip["stop_id"]==s
24                     time_arrv_end= datetime.strptime(trip[trip["stop_id"]==e
25
26                     if((time_dpt_start>=acceptedStart) & (time_dpt_start <=a
27                         difference= time_arrv_end-time_dpt_start
28
29                     averageTime.append(difference.total_seconds()/60)
30                     flag=0
31                 except:
32                     pass
33
34     return (int(Average(averageTime)),flag)
35
36
37

```

In []:

```
1 real_state.head
```

In []:

```

1  #print(real_state.columns)
2  for i in range(len(real_state['property_id'])):
3      res=findAverage( real_state.Train_station_id[i],flinders)
4      real_state.travel_min_to_CBD[i]=res[0]
5      real_state.Transfer_flag[i]=res[1]
6      print(res[0])
7      print(res[1])
8
9
10

```

In []:

```
1 real_state
```

Finding Distance to Hospital and the hospital id

Using haversine distance we have created a function to calculate the distance to the nearest hospital with each record in the dataframe

```
In [ ]: 1 def distance_to_hospital(x1,y1,id1):
2         Radius = 6378.0
3         x2 = radians(hospital['lat'][hospital['id'] == id1])
4         y2 = radians(hospital['lng'][hospital['id'] == id1])
5         dist_long = y2 - y1
6         dist_lat = x2 - x1
7         mat = sin(dist_lat / 2)**2 + cos(x1) * cos(x2) * sin(dist_long / 2)**2
8         mat1 = 2 * atan2(sqrt(mat), sqrt(1 - mat))
9         distance = Radius * mat1
10        return distance, id1

In [ ]: 1 # Iterating over dataframe and finding the nearest Hospital and
2         # and appending the distance from the property along with Hospital id to eac
3
4
5         for i in range(len(real_state['property_id'])):
6             final_list = []
7             dis_list = []
8             name_list = []
9             for j in hospital['id']:
10
11                 dis,name = distance_to_hospital(radians(real_state.lat[i])\
12                                                  ,radians(real_state.lng[i]),j)
13
14                 dis_list.append(dis)
15                 name_list.append(name)
16             final_list = list(zip(dis_list,name_list))
17             final_dis,final_name = min(final_list)[0],min(final_list)[1]
18             real_state.at[i,'Distance_to_hospital'] = (round(final_dis,4)*1000)
19             real_state.at[i,'Hospital_id'] = final_name

In [ ]: 1 # checking if all the values are replaced with the Hospital id
2         #and distance to hospital is calculated
3         real_state[real_state['Hospital_id']=='not available']
4

In [ ]: 1 real_state[real_state['Distance_to_hospital']==0]
```

Finding Distance to Super Market and SuperMarket ID

Using haversine distance we have created a function to calculate the distance to the nearest super market with each record in the dataframe

```
In [ ]: 1 def distance_to_supermarket(x1,y1,id):
2       Radius = 6378.0
3       x2 = radians(superMarket['lat'][superMarket['id'] == id])
4       y2 = radians(superMarket['lng'][superMarket['id'] == id])
5       dist_long = y2 - y1
6       dist_lat = x2 - x1
7       mat = sin(dist_lat / 2)**2 + cos(x1) * cos(x2) * sin(dist_long / 2)**2
8       mat1 = 2 * atan2(sqrt(mat), sqrt(1 - mat))
9       distance = Radius * mat1
10      return distance, id
11
```

```
In [ ]: 1 # Iterating over dataframe and finding the nearest SuperMarket and
2 # and appending the distance from the property along with Supermarket ID to
3
4 for i in range(len(real_state['property_id'])):
5     final_list = []
6     dis_list = []
7     name_list = []
8     for j in superMarket['id']:
9         dis,name = distance_to_supermarket(radians(real_state.lat[i])\
10                                           ,radians(real_state.lng[i]),j)
11         dis_list.append(dis)
12         name_list.append(name)
13     final_list = list(zip(dis_list,name_list))
14     final_dis,final_name = min(final_list)[0],min(final_list)[1]
15     real_state.at[i,'Distance_to_supermaket'] = (round(final_dis,4)*1000)
16     real_state.at[i,'Supermarket_id'] = final_name
```

```
In [ ]: 1 # checking if all the values are replaced with the Super Market ID
2 #and distance to superMarket is calculated
3 real_state[real_state['Supermarket_id']=='not available']
```

```
In [ ]: 1 real_state[real_state['Distance_to_supermaket']==0]
```

```
In [ ]: 1 real_state
```

Writing the Final Dataframe to CSV FILE:

```
In [ ]: 1 real_state.to_csv('30823293_solution.csv',index=False)
```

TASK 2: Data Reshaping

In this task we have to understand the effect of various normalisation/tranformation methods i.e. standardization, MinMax normalisation, log, power box-cox transformation on various columns which are price, Distance_to_sc, travel_min_to_CBD, and Distance_to_hospital and explain the effect. To develop a linear model to predict the price using Distance_to_sc, travel_min_to_CBD, and Distance_to_hospital. For this we have done.

```
In [ ]: 1 new_real_state=real_state[["price","Distance_to_sc","travel_min_to_CBD","Dis
2 new_real_state['price']=new_real_state['price'].astype(float)
```

Plot Description of all column:

As we can see that the plot for Price and Distance to Hospital is a Skewed and both are Right Skewed this is because the main given data is skewed and some of the values are concentrated towards Right

```
In [ ]: 1 from scipy import stats
2 import numpy as np
3 from matplotlib import pyplot as plt
4
5 fig = plt.figure(figsize=(12, 8))
6 ax1 = fig.add_subplot(221)
7 ax2 = fig.add_subplot(222)
8 ax3 = fig.add_subplot(223)
9 ax4 = fig.add_subplot(224)
10
11 ax1.hist((new_real_state['price'],0.5),bins=100)
12 ax2.hist((new_real_state['Distance_to_sc'],0.5),bins=100)
13 ax3.hist((new_real_state['travel_min_to_CBD'],0.5),bins=100)
14 ax4.hist((new_real_state['Distance_to_hospital'],0.5),bins=100)
15 plt.show()
```

Transformation:

I have used log transformation in price column and root transformation in Distance_to_sc and Distance_to_hospital to normalize the data in these columns

```
In [ ]: 1 from scipy import stats
2 import numpy as np
3 from matplotlib import pyplot as plt
4
5 fig = plt.figure(figsize=(12, 8))
6 ax1 = fig.add_subplot(221)
7 ax2 = fig.add_subplot(222)
8 ax3 = fig.add_subplot(223)
9 ax4 = fig.add_subplot(224)
10
11
12
13
14
15 ax1.hist(np.log(new_real_state['price']+1),bins=100)
16 ax2.hist(np.power(new_real_state['Distance_to_sc'],0.6),bins=100)
17 ax3.hist((new_real_state['travel_min_to_CBD'],0.5),bins=100)
18 ax4.hist(np.power(new_real_state['Distance_to_hospital'],0.5),bins=100)
19 plt.show()
```

```
In [ ]: 1 Transformed_real_state=new_real_state.copy()
        2
        3 Transformed_real_state['price']=(np.log(new_real_state['price']+1))
        4 Transformed_real_state['Distance_to_hospital']=(np.power(new_real_state['Dis
        5 Transformed_real_state['Distance_to_sc']=(np.power(new_real_state['Distance_
        6 Transformed_real_state
```

Plotting Scatter Plot:

In this part I have plotted a scatter plot to display the linearty

```
In [ ]: 1 plt.scatter(Transformed_real_state.price,Transformed_real_state.Distance_to_
        2 plt.show()
```

```
In [ ]: 1 plt.scatter(Transformed_real_state.price,Transformed_real_state.Distance_to_
        2 plt.show()
```

```
In [ ]: 1 plt.scatter(Transformed_real_state.price,Transformed_real_state.travel_min_t
        2 plt.show()
```

References:

[1] <https://www.geeksforgeeks.org/convert-html-table-into-csv-file-in-python/>
(<https://www.geeksforgeeks.org/convert-html-table-into-csv-file-in-python/>)

[2] <https://stackoverflow.com/questions/36399381/whats-the-fastest-way-of-checking-if-a-point-is-inside-a-polygon-in-python> (<https://stackoverflow.com/questions/36399381/whats-the-fastest-way-of-checking-if-a-point-is-inside-a-polygon-in-python>)

[3] <https://stackoverflow.com/questions/19412462/getting-distance-between-two-points-based-on-latitude-longitude> (<https://stackoverflow.com/questions/19412462/getting-distance-between-two-points-based-on-latitude-longitude>)

[4] <https://stackoverflow.com/questions/28259301/how-to-convert-an-xml-file-to-nice-pandas-dataframe> (<https://stackoverflow.com/questions/28259301/how-to-convert-an-xml-file-to-nice-pandas-dataframe>)