

# **Data Technician**

Name: Sara Haider

Course Date: 22/05/2025

#### **Table of contents**

Day 2: Task 1	3
Day 3: Task 1	4
Exercise 1: Loading and Exploring the Data	4
Exercise 2: Indexing and Slicing	6
Exercise 3: Data Manipulation	8
Exercise 4: Aggregation and Grouping	9
Exercise 5: Advanced Operations	11
Exercise 6: Exporting Data	12
Exercise 7: If finished early try visualising the results	12
Day 4: Task 1	13
4: Task 2	15
Course Notes	19
Additional Information	19

#### Day 2: Task 1

It is a common software development interview question to create the below with a certain programming language. Create the below using Python syntax, test it and past the completed syntax and output below.

#### FizzBuzz:

Go through the integers from 1 to 100. If a number is divisible by 3, print "fizz." If a number is divisible by 5, print "buzz." If a number is both divisible by 3 and by 5, print "fizzbuzz." Otherwise, print just the number.

```
[2] # Go through the integers from 1 to 100.
                           for x in range(1, 101):
                               # If a number is both divisible by 3 and by 5, print "fizzbuzz"
                               if x \% 3 == 0 and x \% 5 == 0:
                                   print("fizzbuzz")
                               # If a number is divisible by 3, print "fizz"
                               elif x % 3 == 0:
                                   print("fizz")
                               # If a number is divisible by 5, print "buzz"
                               elif x % 5 == 0:
                                   print("buzz")
                               # Otherwise, print just the number
   Paste your
                               else:
completed work
                                   print(x)
  to the right
                       → 1
                           2
                           fizz
                           4
                           buzz
                           fizz
                           8
                           fizz
                           buzz
                           11
                           fizz
                           13
                           14
                           fizzbuzz
```

#### Day 3: Task 1

Download the 'student.csv', complete the below exercises as a group and paste your input and output. Although this is a group activity, everyone should have the below answered so it supports your portfolio:

#### **Exercise 1: Loading and Exploring the Data**

- 1. Question: "Write the code to read a CSV file into a Pandas DataFrame."
- 2. Question: "Write the code to display the first 5 rows of the DataFrame."
- 3. Question: "Write the code to get the information about the DataFrame."
- 4. Question: "Write the code to get summary statistics for the DataFrame."
- 1. Import by clicking files>upload>click on student file

```
import pandas as pd

df = pd.read_csv('student.csv')
print(df)
```

```
id
              name class
                          mark
                                gender
                                female
    1
          John Deo Four
0
                            75
                                  male
1
    2
          Max Ruin Three
                            85
2
    3
            Arnold Three
                            55
                                  male
3
      Krish Star Four
                            60 female
    4
4
    5
         John Mike Four
                            60 female
                                  male
5
         Alex John Four
                            55
6
    7 My John Rob Fifth
                                  male
                            78
```

2.

```
print (df.head(5))
```

```
<del>.</del>₹
       id
                               mark
                                     gender
                        class
                  name
                                     female
        1
              John Deo
    0
                         Four
                                 75
        2
             Max Ruin Three
                                       male
    1
                                 85
    2
        3
                Arnold Three
                                 55
                                       male
                                 60 female
    3
        4 Krish Star
                         Four
    4
        5
            John Mike
                         Four
                                 60 female
```

3.

```
print(df.info())
→ <class 'pandas.core.frame.DataFrame'>
   RangeIndex: 35 entries, 0 to 34
   Data columns (total 5 columns):
    # Column Non-Null Count Dtype
       id
              35 non-null
                            int64
    1 name 34 non-null object
    2 class 34 non-null object
    3 mark
               35 non-null
                             int64
       gender 33 non-null
                           object
   dtypes: int64(2), object(3)
   memory usage: 1.5+ KB
   None
```

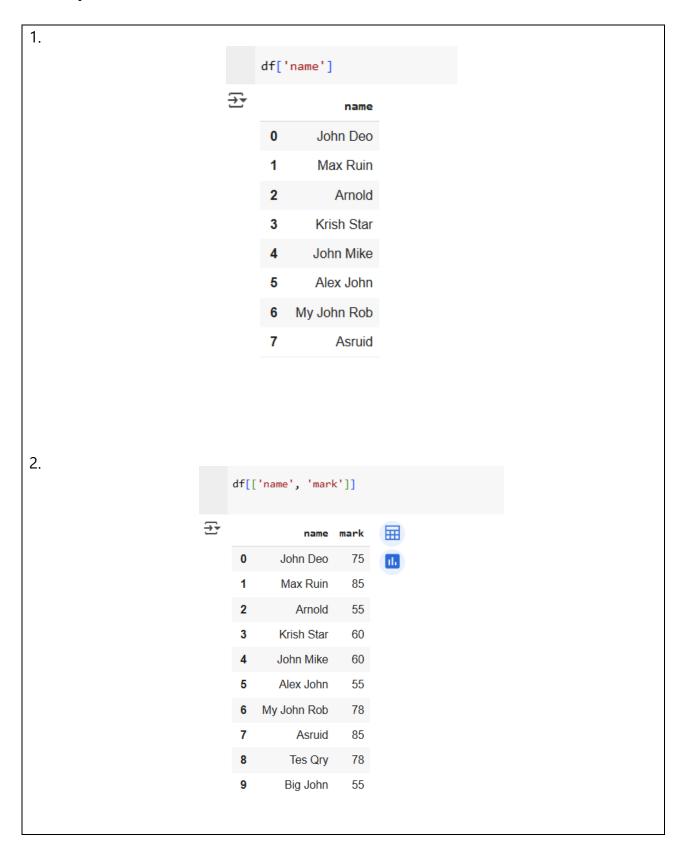
4.

```
print(df.describe())
```

```
id
                       mark
       35.000000
                  35.000000
count
mean
       18.000000
                 74.657143
std
       10.246951
                 16.401117
min
        1.000000 18.000000
25%
        9.500000 62.500000
50%
       18.000000 79.000000
75%
       26.500000 88.000000
       35.000000 96.000000
max
```

### **Exercise 2: Indexing and Slicing**

- 1. Question: "Write the code to select the 'name' column."
- 2. Question: "Write the code to select the 'name' and 'mark' columns."
- 3. Question: "Write the code to select the first 3 rows."
- 4. Question: "Write the code to select all rows where the 'class' is 'Four'."



3.

df.iloc[:3]

	id	name	class	mark	gender
0	1	John Deo	Four	75	female
1	2	Max Ruin	Three	85	male
2	3	Arnold	Three	55	male

4.

df[df['class'] == 'Four']

₹ id name class mark gender 0 1 John Deo Four 75 female 3 60 4 Krish Star Four female 5 John Mike Four 60 female 6 55 5 Alex John Four male 9 10 Big John Four 55 female 88 15 16 Gimmy Four male 21 Babby John 69 20 Four female **30** 31 Marry Toeey Four 88 male

### **Exercise 3: Data Manipulation**

- 1. Question: "Write the code to add a new column 'passed' that indicates whether the student passed (mark >= 60)."
- 2. Question: "Write the code to rename the 'mark' column to 'score'."
- 3. Question: "Write the code to drop the 'passed' column."

```
1.
                df['passed'] = df['mark'] >= 60
                print(df)
           ₹
                    id
                                                    gender
                                name
                                      class
                                              mark
                                                            passed
                                                    female
                     1
                            John Deo
                                       Four
                                                75
                                                               True
                1
                     2
                            Max Ruin
                                      Three
                                                85
                                                      male
                                                               True
                2
                     3
                              Arnold
                                                55
                                                      male
                                                              False
                                      Three
                3
                         Krish Star
                                                   female
                     4
                                       Four
                                                60
                                                              True
                     5
                           John Mike
                                                   female
                4
                                       Four
                                                60
                                                               True
                5
                     6
                          Alex John
                                                      male
                                                              False
                                      Four
                                                55
                6
                     7
                        My John Rob Fifth
                                                78
                                                              True
                                                      male
                7
                                      Five
                     8
                              Asruid
                                                85
                                                      male
                                                               True
                8
                     9
                             Tes Qry
                                        Six
                                                78
                                                       NaN
                                                               True
                9
                                                   female
                    10
                            Big John
                                       Four
                                                55
                                                              False
                                                   female
                10
                   11
                              Ronald
                                        Six
                                                89
                                                              True
                    12
                               Recky
                                        Six
                                                94
                                                   female
                                                               True
                11
                12
                    13
                                 Kty
                                      Seven
                                                88
                                                   female
                                                               True
                13
                    14
                                Bigy
                                                88
                                                   female
                                                               True
                                      Seven
2.
            df.rename(columns={'mark': 'score'}, inplace=True)
            print(df)
                id
                                class score
                                               gender
                                                       passed
                           name
           0
                1
                       John Deo
                                  Four
                                           75
                                               female
                                                         True
            1
                 2
                       Max Ruin
                                 Three
                                           85
                                                 male
                                                         True
            2
                 3
                         Arnold
                                Three
                                           55
                                                 male
                                                        False
            3
                    Krish Star
                                           60 female
                4
                                  Four
                                                         True
           4
                5
                      John Mike
                                  Four
                                           60 female
                                                         True
            5
                6
                      Alex John
                                  Four
                                           55
                                                 male
                                                        False
            6
                7
                   My John Rob Fifth
                                           78
                                                 male
                                                         True
            7
                8
                         Asruid
                                  Five
                                           85
                                                 male
                                                         True
            8
                9
                        Tes Qry
                                   Six
                                           78
                                                  NaN
                                                         True
            9
               10
                       Big John
                                  Four
                                           55
                                              female
                                                        False
                                   Six
                                           89 female
            10
              11
                         Ronald
                                                         True
                                           94 female
               12
                          Recky
                                   Six
                                                         True
            11
            12
               13
                            Kty Seven
                                           88 female
                                                         True
                           Rimy Sovon
                                           ΩΩ
                                              fomalo
            12
                                                         Tnuo
```



```
3.
                df.drop('passed', axis=1, inplace=True)
                print(df)
                                                     gender
           ₹
                    id
                                name
                                      class
                                             score
                0
                     1
                           John Deo
                                       Four
                                                 75
                                                     female
                1
                     2
                                                 85
                                                       male
                           Max Ruin
                                     Three
                2
                     3
                             Arnold Three
                                                 55
                                                       male
                3
                         Krish Star
                                                    female
                     4
                                       Four
                                                 60
                                                    female
                4
                     5
                          John Mike
                                       Four
                                                 60
                5
                          Alex John
                                                       male
                     6
                                       Four
                                                 55
                                                       male
                6
                     7
                        My John Rob
                                      Fifth
                                                 78
                7
                     8
                             Asruid
                                      Five
                                                 85
                                                       male
                     9
                            Tes Ory
                                        Six
                                                 78
                                                        NaN
                8
                9
                           Big John
                                                 55
                                                     female
                    10
                                       Four
                10
                    11
                             Ronald
                                        Six
                                                 89
                                                     female
```

#### **Exercise 4: Aggregation and Grouping**

- 1. Question: "Write the code to group the DataFrame by the 'class' column and calculate the mean 'mark' for each group."
- 2. Question: "Write the code to count the number of students in each class."
- 3. Question: "Write the code to calculate the average mark for each gender."

```
df.groupby('class')['mark'].mean()
1.
                               ₹
                                                 mark
                                     class
                                     Eight 79.000000
                                            78.000000
                                     Fifth
                                      Five
                                            80.000000
                                            68.750000
                                     Four
                                     Nine
                                            41.500000
                                     Seven 77.600000
                                            82.571429
                                      Six
                                     Three 73.666667
                                    dtype: float64
```



```
2.
                        df['class'].value_counts()
                   ₹
                                count
                         class
                                   10
                         Seven
                                   8
                         Four
                          Six
                         Three
                                    3
                                    2
                         Nine
                         Five
                                   2
                         Fifth
                         Eight
                                    1
3.
                df.groupby('gender')['mark'].mean()
                                mark
                  gender
                          77.312500
                  female
                           71.588235
                  male
```

#### **Exercise 5: Advanced Operations**

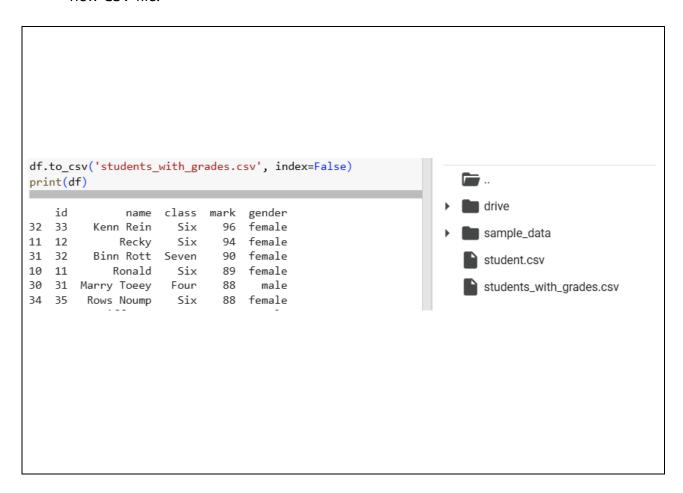
- 1. Question: "Write the code to create a pivot table with 'class' as rows, 'gender' as columns, and 'mark' as values."
- 2. Question: "Write the code to create a new column 'grade' where marks >= 85 are 'A', 70-84 are 'B', 60-69 are 'C', and below 60 are 'D'."
- 3. Question: "Write the code to sort the DataFrame by 'mark' in descending order."

```
1.
              pivot = df.pivot_table(values='mark', index='class', columns='gender', aggfunc='mean')
              print(pivot)
              gender female male
              class
                        NaN 79.0
              Eight
              Fifth
                       NaN 78.0
              Five
                        NaN 80.0
                      63.8 77.0
              Four
                     65.0 18.0
81.4 73.8
              Nine
              Seven
                      89.2 54.0
              Six
              Three
                       NaN 70.0
2.
                                     df['grade'] = pd.cut(
                                         df['mark'],
                                         bins=[0, 59, 69, 84, 100],
                                        labels=['D', 'C', 'B', 'A']
                                     print(df[['name', 'mark', 'grade']])
                                               name mark grade
                                           John Deo
                                          Max Ruin
                                                     85
                                                    55
60
                                            Arnold
                                        Krish Star
                                                            C
                                         John Mike
                                                     60
55
                                          Alex John
                                     6 My John Rob
                                                     78
85
                                                            В
                                            Asruid
                                                     78
55
                                            Tes Qry
                                                            В
                                          Big John
                                                            D
                                           Ronald
Recky
                                     10
                                                     89
                                                            Α
                                     11
                                                            Α
                                             Kty
Bigy
                                                     88
                                                            Α
                                     12
                                                     88
                                                            Α
                                     13
                                     14
                                          Tade Row
                                                     88
                                                            Α
                                             Gimmy
3.
                            df.sort_values(by='mark', ascending=False, inplace=True)
                            print(df[['name', 'mark']])
                                      name mark
                            32 Kenn Rein
                                             96
                            11
                                    Recky
                                             94
                            31 Binn Rott
                                             90
                            10
                                   Ronald
                                             89
                            30 Marry Toeey
                                             88
                            34 Rows Noump
                                             88
                                 Giff Tow
                            24
                            14
                                  Tade Row
                                             88
                            15
                                   Gimmy
                            12
                                      Kty
                                             88
                            13
                                     Bigy
                            27
                                Rojj Base
                                             86
                                    Asruid
                                  Max Ruin
```

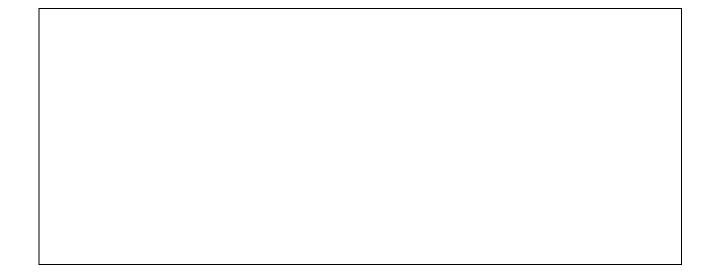


#### **Exercise 6: Exporting Data**

 Question: "Write the code to save the DataFrame with the new 'grade' column to a new CSV file."



### **Exercise 7: If finished early try visualising the results**



#### Day 4: Task 1

Using the 'GDP (nominal) per Capita.csv' which can be downloaded from the shared Folder, complete the below exercises and paste your input and output. Work individually, but we will work and support each other in the room.

- Read and save the 'GDP (nominal) per Capita' data to a data frame called "df" in Jupyter notebook
- Print the first 10 rows
- Print the last 5 rows
- Print 'Country/Territory' and 'UN\_Region' columns

```
import pandas as pd
df = pd.read_csv('GDP (nominal) per Capita.csv')
print(df.head(10))
  Unnamed: 0 Country/Territory UN_Region \, IMF_Estimate \, IMF_Year \, \
                      Monaco Europe 0 0
tenstein Europe 0 0
0
    1
          2 Liechtenstein Europe
1
                                            0 0
132372 2023
114581 2023
         3 Luxembourg Europe
3
         4
                   Ireland Europe
         4
5
6
8
9
  WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year
                     2021 234317 2021
0
            234316
1
             157755
                             2020
                                         169260 2021
                          2021
2021
2
             133590
                                        133745 2021
3
            100172
                                        101109 2021
                             2021
2021
4
            114090
                                        112653 2021
5
             89154
                                        89242 2021

      2021
      93525
      2021

      2021
      66822
      2021

      2019
      0
      0

      2021
      85250
      2021

6
              91992
7
              72794
8
              87158
9
               86569
```



#### print(df.tail(5)) Unnamed: 0 Country/Territory UN\_Region IMF\_Estimate IMF\_Year 218 219 Malawi Africa 496 2023 219 220 South Sudan Africa 467 2023 220 221 Sierra Leone Africa 415 2023 222 221 Afghanistan Asia 611 2020 222 223 Africa 249 2023 Burundi

	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year
218	635	2021	613	2021
219	1072	2015	400	2021
220	480	2021	505	2021
221	369	2021	373	2021
222	222	2021	311	2021

#### print(df[['Country/Territory', 'UN\_Region']])

	Country/Territory	UN_Region
0	Monaco	Europe
1	Liechtenstein	Europe
2	Luxembourg	Europe
3	Ireland	Europe
4	Bermuda	Americas
218	Malawi	Africa
219	South Sudan	Africa
220	Sierra Leone	Africa
221	Afghanistan	Asia
222	Burundi	Africa

[223 rows x 2 columns]

#### 4: Task 2

Back with 'GDP (nominal) per Capita'. As a group, import and work your way through the Day\_4\_Python\_Activity.ipynb notebook which can be found on the shared Folder. There are questions to answer, but also opportunities to have fun with the data – paste your input and output below.

Once complete, and again as a group, work with some more data and have some fun – there is no set agenda for this section, other than to embed the skills developed this week. Paste your input and output below and upon return we'll discuss progress made.

Additional data found here.

```
# Number of countries per region
 region = df['UN_Region'].value_counts()
 print(region)
 UN Region
            55
 Africa
 Asia
            51
 Americas 48
           48
 Europe
 Oceania
            20
 Name: count, dtype: int64
#What is European Union [n 1]?
 eu = df[df["Country/Territory"] == "European Union[n 1]"]
print(eu)
    Unnamed: 0 Country/Territory UN_Region IMF_Estimate IMF_Year \
     36 European Union[n 1] Europe
    WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year
                                 2021 31875
# Countries in Europe below average
europe_ba = df[(df["UN_Region"] == "Europe") & (df["WorldBank_Estimate"] < df["WorldBank_Estimate"].mean())]</pre>
print(europe_ba[["Country/Territory", "WorldBank_Estimate"]])
        Country/Territory WorldBank Estimate
69
                 Croatia
                                    17685
71
                  Poland
                                    18000
74
                 Hungary
                                    18728
77
                 Romania
                                    14858
86
                Bulgaria
                                    12222
                                    12195
89
                  Russia
102
             Montenegro
                                     9466
105
                 Serbia
                                     9230
111 Bosnia and Herzegovina
                                     7143
                 Belarus
117
          North Macedonia
          Albania
                                     6493
119
126
                 Moldova
                                     5231
132
                                     5270
                  Kosovo
142
                 Ukraine
                                     4836
```



```
#Which countries in Europe has higher GDP than UK?
gdp = df.loc[df["Country/Territory"] == "United Kingdom", "WorldBank_Estimate"].values[0]
# Find European countries with higher GDP
higher_gdp= df[(df["UN_Region"] == "Europe") & (df["WorldBank_Estimate"] > gdp)][["Country/Territory", "WorldBank_Estimate"]]
print(higher_gdp)
  Country/Territory WorldBank_Estimate
            Monaco
                             234316
     Liechtenstein
                             157755
      Luxembourg
                             133590
         Ireland
3
                            100172
                             89154
5
            Norway
     Switzerland
Isle of Man
6
                              91992
8
                              87158
                              68728
12
           Iceland
13 Channel Islands
                              75153
    Faroe Islands
14
                             69010
                              68008
15
           Denmark
     Denmark
Netherlands
                              57768
17
       Austria
19
                              53638
21
            Sweden
                              61029
       Finland
Belgium
Germany
22
                              53655
23
                              51247
27
                              51204
#What is the average GDP for each region?
 # Group by region and calculate average GDP
 avg_gdp_reg = df.groupby("UN_Region")["WorldBank_Estimate"].mean()
 print(avg_gdp_reg)
 UN_Region
              2470.836364
 Africa
           18565.125000
 Americas
Asia 13921.313725
Europe 45193.687500
Oceania 15113.650000
World 12235.000000
Name: WorldBank_Estimate, dtype: float64
#Which countries below average by IMF world estimate?
 # Calculate the average world IMF
 avg_imf = df["IMF_Estimate"].mean()
 # Filter and display countries below the average
 print(df[df["IMF_Estimate"] < avg_imf][["Country/Territory", "IMF_Estimate"]])</pre>
    Country/Territory IMF_Estimate
              Monaco 0
      Liechtenstein
 1
             Bermuda
        Bermuda
Isle of Man
 8
     Cayman Islands
             Malawi
        South Sudan
 219
                             467
       Sierra Leone
        Afghanistan
             Burundi
 [159 rows x 2 columns]
```



```
#IMF estimate 0 values
# Countries where IMF estimate is 0
print(df[df["IMF_Estimate"] == 0][["Country/Territory"]])
           Country/Territory
                   Monaco
             Liechtenstein
1
4
                   Bermuda
8
               Isle of Man
9
             Cayman Islands
           Channel Islands
13
14
              Faroe Islands
18
                 Greenland
30 British Virgin Islands
      US Virgin Islands
36
             New Caledonia
41
                      Guam
57 Sint Maarten (Dutch part)
60
    Northern Mariana Islands
64 Saint Martin (French part)
    Turks and Caicos Islands
70
      French Polynesia
75
              Cook Islands
76
                  Anguilla
81
                   Curação
                Montserrat
             American Samoa
103
                     Cuba
195
                  Zanzibar
203
                     Syria
                North Korea
211
#Which country has highest UN Estimate?
print(df.loc[df["UN_Estimate"].idxmax(), ["Country/Territory", "UN_Estimate"]])
Country/Territory
                       Monaco
UN_Estimate
                       234317
Name: 0, dtype: object
#Which country has highest Worldbank Estimate?
print(df.loc[df["WorldBank_Estimate"].idxmax(), ["Country/Territory", "WorldBank_Estimate"]])
Country/Territory
                   Monaco
WorldBank_Estimate
                   234316
Name: 0, dtype: object
#Which country has highest IMF Estimate?
 print(df.loc[df["IMF_Estimate"].idxmax(), ["Country/Territory", "IMF_Estimate"]])
 Country/Territory Luxembourg
                          132372
 IMF_Estimate
 Name: 2, dtype: object
```



```
#Filling 0 Values by average
 import numpy as np
 # replace 0 with null values
 df.replace(0, np.nan, inplace=True)
 print(df.head())
   Unnamed: 0 Country/Territory UN_Region IMF_Estimate IMF_Year \
           a.
1
                   Monaco Europe NaN NaN
Liechtenstein Europe NaN NaN
Luxembourg Europe 132372.0 2023.0
Ireland Europe 114581.0 2023.0
Bermuda Americas NaN NaN
 0
 1
            3
 2
 3
             4
 4
            5
    WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year
     234316.0 2021.0 234317.0 2021
157755.0 2020.0 169260.0 2021
 1
            133590.0 2021.0 133745.0 2021
100172.0 2021.0 101109.0 2021
114090.0 2021.0 112653.0 2021
 2
 3
 4
 # Calculate the average of 'Worldbank Estimate' and 'UN Estimate' columns
 df["Average_Estimate"] = df[["WorldBank_Estimate", "UN_Estimate"]].mean(axis=1)
 print(df[["Country/Territory", "Average_Estimate"]])
     Country/Territory Average_Estimate
     Monaco 234316.5
Liechtenstein 163507.5
 0
 1
         Luxembourg 133667.5
Ireland 100640.5
Bermuda 113371.5
 2
 3
 4
        Malawi
South Sudan
Sierra Leone
Afghanistan
                               624.0
 218
                                 736.0
492.5
371.0
266.5
 219
 220
 221
               Burundi
 [223 rows x 2 columns]
  # Fill null values with the average
  df["IMF_Estimate"].fillna(df["IMF_Estimate"].mean(), inplace=True)
 print(df[["Country/Territory", "IMF_Estimate"]].head())
   Country/Territory IMF_Estimate
Monaco 17377.736041
        Liechtenstein 17377.736041
  1
         Luxembourg 132372.000000
  2
  3
               Ireland 114581.000000
                Bermuda 17377.736041
#Missing Values
 missing = df.isnull().sum()
 print(missing)
 Unnamed: 0
 Country/Territory
                          0
 UN_Region
 IMF_Estimate
 IMF_Year
                         26
 WorldBank Estimate
 WorldBank_Year
 UN Estimate
                          9
 UN_Year
                          0
 Average_Estimate
 dtype: int64
```



## **Course Notes**

It is recommended to take notes from the course, use the space below to do so, or use the revision guide shared with the class:

We have included a range of additional links to further resources and information that you may find useful, these can be found within your revision guide.

#### **END OF WORKBOOK**

Please check through your work thoroughly before submitting and update the table of contents if required.

Please send your completed work booklet to your trainer.

