



# Data Technician

**Name:**

**Course Date:**

## **Table of contents**

Day 2: Task 1 .....	3
Day 3: Task 1 .....	3
Exercise 1: Loading and Exploring the Data .....	4
Exercise 2: Indexing and Slicing .....	4
Exercise 3: Data Manipulation .....	6
Exercise 4: Aggregation and Grouping .....	8
Exercise 5: Advanced Operations .....	8
Exercise 6: Exporting Data .....	9
Exercise 7: If finished early try visualising the results .....	10
Day 4: Task 1 .....	10
Day 4: Task 2 .....	12
Course Notes .....	26
Additional Information .....	27



## 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.

Paste  
your  
complet  
ed work  
to the  
right

```
1 for x in range(1,101):
2     if x%3 == 0 and x%5 ==0:
3         print("The number is",x,"it is fizzbuzz.")
4     elif x%5 == 0:
5         print("The number is",x,"it is buzz.")
6     elif x%3 == 0:
7         print("The number is",x,"it is fizz.")
8     else:
9         print("The number is",x,"it is indivisible by 3 or 5.")
10    continue
```

The number is 1 , it is indivisible by 3 or 5.  
The number is 2 , it is indivisible by 3 or 5.  
The number is 3 , it is fizz.  
The number is 4 , it is indivisible by 3 or 5.  
The number is 5 , it is buzz.  
The number is 6 , it is fizz.  
The number is 7 , it is indivisible by 3 or 5.  
The number is 8 , it is indivisible by 3 or 5.  
The number is 9 , it is fizz.  
The number is 10 , it is buzz.  
The number is 11 , it is indivisible by 3 or 5.  
The number is 12 , it is fizz.  
The number is 13 , it is indivisible by 3 or 5.  
The number is 14 , it is indivisible by 3 or 5.  
The number is 15 , it is fizzbuzz.  
The number is 16 , it is indivisible by 3 or 5.  
The number is 17 , it is indivisible by 3 or 5.  
The number is 18 , it is fizz.  
The number is 19 , it is indivisible by 3 or 5.  
The number is 20 , it is buzz.

## 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 df_student = pd.read_csv(r"D:\桌面\bootcamp\week 6\student.csv")
2 print("The first 5 rows of the DataFrame is:\n",df_student.head(5))
3 print("\n-----\n")
4 print("The information of the DataFrame is:")
5 df_student.info()
6 print("\n-----\n")
7 #print("The summary statistics for the DataFrame is:\n",df_student.describe())
8 print("The summary statistics for the DataFrame is:\n",df_student['mark'].describe()) #Focused on the summary statistics for
```

The first 5 rows of the DataFrame is:

	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
3	4	Krish Star	Four	60	female
4	5	John Mike	Four	60	female

-----

The information of the DataFrame is:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 35 entries, 0 to 34
Data columns (total 5 columns):
 #   Column  Non-Null Count  Dtype  
--- 
 0   id      35 non-null    int64  
 1   name    34 non-null    object  
 2   class   34 non-null    object  
 3   mark    35 non-null    int64  
 4   gender  33 non-null    object  
dtypes: int64(2), object(3)
memory usage: 1.5+ KB
```

-----

The summary statistics for the DataFrame is:

```
count    35.000000
mean     74.657143
std      16.401117
min      18.000000
25%     62.500000
50%     79.000000
75%     88.000000
max     96.000000
Name: mark, dtype: float64
```

## 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'."



```

1 df_student = pd.read_csv(r"D:\桌面\bootcamp\week 6\student.csv")
2 print("The 'name' column in the DataFrame is:\n",df_student['name'])
3 print("\n-----\n")
4 print("The 'name' and 'mark' columns in the DataFrame are:\n",df_student[['name','mark']])
5 print("\n-----\n")
6 print("The first 3 rows in the DataFrame are:\n",df_student.loc[:2])
7 print("The first 3 rows in the DataFrame are:\n",df_student.iloc[:3]) #Alternatively, using df.iLoc to achieve the task
8 print("\n-----\n")
9 print("Rows where the 'class' is 'Four' in the DataFrame are:\n",df_student[df_student['class']=="Four"])

```

The 'name' column in the DataFrame is:

```

0      John Deo
1      Max Ruin
2      Arnold
3      Krish Star
4      John Mike
5      Alex John
6      My John Rob
7      Asruid
8      Tes Qry
9      Big John
10     Ronald
11     Recky
12     Kty
13     Bigy
14     Tade Row
15     Gimmy
16     Tumyu
17     Honny
18     Tinny
19     Jackly
20     Babby John
21     Reggid
22     Herod
23     Tiddy Now
24     Giff Tow
25     Crelea
26     NaN
27     Rojj Base
28     Tess Played
29     Reppy Red
30     Marry Toeey
31     Binn Rott
32     Kenn Rein
33     Gain Toe
34     Rows Noump
Name: name, dtype: object
-----
```

The 'name' and 'mark' columns in the DataFrame are:

	name	mark
0	John Deo	75
1	Max Ruin	85
2	Arnold	55
3	Krish Star	60
4	John Mike	60
5	Alex John	55
6	My John Rob	78
7	Asruid	85
8	Tes Qry	78
9	Big John	55
10	Ronald	89
11	Recky	94
12	Kty	88
13	Bigy	88
14	Tade Row	88
15	Gimmy	88
16	Tumyu	54
17	Honny	75
18	Tinny	18
19	Jackly	65
20	Babby John	69
21	Reggid	55
22	Herod	79
23	Tiddy Now	78
24	Giff Tow	88
25	Crelea	79
26	NaN	81
27	Rojj Base	86
28	Tess Played	55
29	Reppy Red	79
30	Marry Toeey	88
31	Binn Rott	90
32	Kenn Rein	96
33	Gain Toe	69
34	Rows Noump	88



```
-----
The first 3 rows in the DataFrame are:
   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
The first 3 rows in the DataFrame are:
   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
-----
Rows where the 'class' is 'Four' in the DataFrame are:
   id      name  class  mark gender
0  1  John Deo    Four   75 female
3  4  Krish Star  Four   60 female
4  5  John Mike   Four   60 female
5  6  Alex John   Four   55  male
9 10  Big John   Four   55 female
15 16  Gimmy     Four   88  male
20 21 Babby John  Four   69 female
30 31 Marry Toey  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_student = pd.read_csv(r"D:\桌面\bootcamp\week 6\student.csv")
2 df_student['passed'] = df_student['mark'] >= 60
3 print("The students who passed the exam are marked as followed:\n",df_student)
4 print("\n-----\n")
5 df_renamed = df_student.copy().rename(columns={'mark':'score'})
6 print("Renaming the 'mark' column to 'score' in the DataFrame:\n",df_renamed)
7 print("\n-----\n")
8 df_drop = df_renamed.drop(columns=['passed'])
9 print("Dropping the column named 'passed' in the DataFrame:\n",df_drop)
```

The students who passed the exam are marked as followed:

	id	name	class	mark	gender	passed
0	1	John Deo	Four	75	female	True
1	2	Max Ruin	Three	85	male	True
2	3	Arnold	Three	55	male	False
3	4	Krish Star	Four	60	female	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
10	11	Ronald	Six	89	female	True
11	12	Recky	Six	94	female	True
12	13	Kty	Seven	88	female	True
13	14	Bigy	Seven	88	female	True
14	15	Tade Row	NaN	88	male	True
15	16	Gimmy	Four	88	male	True
16	17	Tumyu	Six	54	male	False
17	18	Honny	Five	75	male	True
18	19	Tinnny	Nine	18	male	False
19	20	Jacklv	Nine	65	female	True



```

20 21 Babby John Four 69 female True
21 22 Reggid Seven 55 female False
22 23 Herod Eight 79 male True
23 24 Tiddy Now Seven 78 male True
24 25 Giff Tow Seven 88 male True
25 26 Crelea Seven 79 male True
26 27 NaN Three 81 NaN True
27 28 Rojj Base Seven 86 female True
28 29 Tess Played Seven 55 male False
29 30 Reppy Red Six 79 female True
30 31 Marry Toeey Four 88 male True
31 32 Binn Rott Seven 90 female True
32 33 Kenn Rein Six 96 female True
33 34 Gain Toe Seven 69 male True
34 35 Rows Noump Six 88 female True
-----
```

Renaming the 'mark' column to 'score' in the DataFrame:

```

id      name class score gender passed
0 1 John Deo Four 75 female True
1 2 Max Ruin Three 85 male True
2 3 Arnold Three 55 male False
3 4 Krish Star Four 60 female 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
10 11 Ronald Six 89 female True
11 12 Becky Six 94 female True
12 13 Kty Seven 88 female True
13 14 Bigy Seven 88 female True
14 15 Tade Row NaN 88 male True
15 16 Gimmy Four 88 male True
16 17 Tumyu Six 54 male False
17 18 Honny Five 75 male True
18 19 Tinny Nine 18 male False
19 20 Jackly Nine 65 female True
20 21 Babby John Four 69 female True
21 22 Reggid Seven 55 female False
22 23 Herod Eight 79 male True
23 24 Tiddy Now Seven 78 male True
24 25 Giff Tow Seven 88 male True
25 26 Crelea Seven 79 male True
26 27 NaN Three 81 NaN True
27 28 Rojj Base Seven 86 female True
28 29 Tess Played Seven 55 male False
29 30 Reppy Red Six 79 female True
30 31 Marry Toeey Four 88 male True
31 32 Binn Rott Seven 90 female True
32 33 Kenn Rein Six 96 female True
33 34 Gain Toe Seven 69 male True
34 35 Rows Noump Six 88 female True
-----
```

Dropping the column named 'passed' in the DataFrame:

```

id      name class score gender
0 1 John Deo Four 75 female
1 2 Max Ruin Three 85 male
2 3 Arnold Three 55 male
3 4 Krish Star Four 60 female
4 5 John Mike Four 60 female
5 6 Alex John Four 55 male
6 7 My John Rob Fifth 78 male
7 8 Asruid Five 85 male
8 9 Tes Qry Six 78 NaN
9 10 Big John Four 55 female
10 11 Ronald Six 89 female
11 12 Becky Six 94 female
12 13 Kty Seven 88 female
13 14 Bigy Seven 88 female
14 15 Tade Row NaN 88 male
15 16 Gimmy Four 88 male
16 17 Tumyu Six 54 male
17 18 Honny Five 75 male
18 19 Tinny Nine 18 male
19 20 Jackly Nine 65 female
20 21 Babby John Four 69 female
21 22 Reggid Seven 55 female
22 23 Herod Eight 79 male
23 24 Tiddy Now Seven 78 male
24 25 Giff Tow Seven 88 male
25 26 Crelea Seven 79 male
26 27 NaN Three 81 NaN
27 28 Rojj Base Seven 86 female
```



```

28 29 Tess Played Seven 55 male
29 30     Reppy Red Six 79 female
30 31 Marry Toeey Four 88 male
31 32     Binn Rott Seven 90 female
32 33 Kenn Rein Six 96 female
33 34 Gain Toe Seven 69 male
34 35 Rows Noump Six 88 female

```

## Exercise 4: Aggregation and Grouping

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

```

1 df_student = pd.read_csv(r"D:\桌面\bootcamp\week 6\student.csv")
2 print(df_student.groupby('class', as_index=False)[['mark']].mean().sort_values('mark', ascending = False))
3 print("\n-----\n")
4 print(df_student.groupby('class',as_index=False).size().sort_values('size',ascending = False))
5 print("\n-----\n")
6 #print(df_student.groupby('Class',as_index=False)['id'].count().sort_values('id',ascending = False))
7 #print("\n-----\n")
8 print(df_student.groupby('gender',as_index=False)[['mark']].mean().sort_values('mark',ascending = False))

class      mark
6   Six  82.571429
2   Five  80.000000
0  Eight  79.000000
1  Fifth  78.000000
5  Seven  77.600000
7  Three  73.666667
3   Four  68.750000
4   Nine  41.500000
-----
class  size
5  Seven    10
3   Four     8
6   Six      7
7  Three     3
2   Five     2
4   Nine     2
0  Eight     1
1  Fifth     1
-----
gender      mark
0  female  77.312500
1   male   71.588235

```

## Exercise 5: Advanced Operations

- Question: "Write the code to create a pivot table with 'class' as rows, 'gender' as columns, and 'mark' as values."
- 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'."
- Question: "Write the code to sort the DataFrame by 'mark' in descending order."



```

1 df_student = pd.read_csv(r"D:\桌面\bootcamp\week 6\student.csv")
2 pd.pivot_table(df_student, index='class', columns='gender', values='mark',aggfunc='mean')

```

gender	female	male
class		
Eight	NaN	79.0
Fifth	NaN	78.0
Five	NaN	80.0
Four	63.8	77.0
Nine	65.0	18.0
Seven	81.4	73.8
Six	89.2	54.0
Three	NaN	70.0

```

1 df_student = pd.read_csv(r"D:\桌面\bootcamp\week 6\student.csv")
2 def assign_grade(mark):
3     if mark >= 85:
4         return 'A'
5     elif mark >= 70:
6         return 'B'
7     elif mark >= 60:
8         return 'C'
9     else:
10        return 'D'
11
12 df_student['grade'] = df_student['mark'].apply(assign_grade)
13 df_student.sort_values('mark', ascending = False)

```

id	name	class	mark	gender	grade	
32	33	Kenn Rein	Six	96	female	A
11	12	Recky	Six	94	female	A
31	32	Binn Rott	Seven	90	female	A
10	11	Ronald	Six	89	female	A
24	25	Giff Tow	Seven	88	male	A
15	16	Gimmy	Four	88	male	A
14	15	Tade Row	NaN	88	male	A
13	14	Bigy	Seven	88	female	A
12	13	Kty	Seven	88	female	A
34	35	Rows Noump	Six	88	female	A

## Exercise 6: Exporting Data

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

```
1 df_student.to_csv('students_with_grades.csv', index=False)
```

id	name	class	mark	gender	grade
1	John Deo	Four	75	female	B
2	Max Ruin	Three	85	male	A
3	Arnold	Three	55	male	D
4	Krish Sta	Four	60	female	C
5	John Mike	Four	60	female	C
6	Alex John	Four	55	male	D
7	My John	RFifth	78	male	B
8	Asrud	Five	85	male	A
9	Tes Qry	Six	78		B
10	Big John	Four	55	female	D
11	Ronald	Six	89	female	A
12	Recky	Six	94	female	A
13	Kty	Seven	88	female	A
14	Bigy	Seven	88	female	A
15	Tade Row		88	male	A
16	Gimmy	Four	88	male	A
17	Tumyu	Six	54	male	D
18	Hony	Five	75	male	B
19	Tinny	Nine	18	male	D
20	Jackly	Nine	65	female	C
21	Babby Joh	Four	69	female	C
22	Reggid	Seven	55	female	D
23	Herod	Eight	79	male	B
24	Tiddy Now	Seven	78	male	B
25	Giff Tow	Seven	88	male	A
26	Crelea	Seven	79	male	B
97		Three	81		R



## 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



```

1 df = pd.read_csv(r"D:\桌面\bootcamp\week 6\Day 4 Resources\GDP (nominal) per Capita.csv")
2 print("The fist 10 rows are:/n",df.head(10))
3 print("\n-----\n")
4 print("The last 5 rows are:/n",df.tail(5))
5 print("\n-----\n")
6 print("‘Country/Territory’ and ‘UN_Region’ columns shown as follow:/n",df[['Country/Territory','UN_Region']])

```

The fist 10 rows are:/n      Unnamed: 0 Country/Territory UN\_Region IMF\_Estimate IMF\_Year \
0      1      Monaco      Europe      0      0
1      2      Liechtenstein      Europe      0      0
2      3      Luxembourg      Europe      132372      2023
3      4      Ireland      Europe      114581      2023
4      5      Bermuda      Americas      0      0
5      6      Norway      Europe      101103      2023
6      7      Switzerland      Europe      98767      2023
7      8      Singapore      Asia      91100      2023
8      9      Isle of Man      Europe      0      0
9      10      Cayman Islands      Americas      0      0

	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year
0	234316	2021	234317	2021
1	157755	2020	169260	2021
2	133590	2021	133745	2021
3	100172	2021	101109	2021
4	114090	2021	112653	2021
5	89154	2021	89242	2021
6	91992	2021	93525	2021
7	72794	2021	66822	2021
8	87158	2019	0	0
9	86569	2021	85250	2021

-----

The last 5 rows are:/n      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
221      222      Afghanistan      Asia      611      2020
222      223      Burundi      Africa      249      2023

	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

-----

‘Country/Territory’ and ‘UN\_Region’ columns shown as follow:

	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]



## Day 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.](#)

```
1 # number of countries per region
2
3 #df.groupby("UN_Region")["Country/Territory"].count()
4 #df.groupby("UN_Region")["Country/Territory"].size()
5 pd.pivot_table(df, index='UN_Region', values='Country/Territory',aggfunc='count')
```

Country/Territory	
UN_Region	
Africa	55
Americas	48
Asia	51
Europe	48
Oceania	20
World	1

```
1 #What is European Union[n 1]?
2
3 df[df["Country/Territory"] == "European Union[n 1]"]
```

```
Country/Territory UN_Region IMF_Estimate IMF_Year WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year UN_Year_Clean
36 European Union[n 1] Europe 39940 2023 38411 2021 31875 2021 2021
```

```
1 # Countries in Europe below average
2
3 #Based on IMF Statistics
4 df_Europe_below_imfavg = df_Europe[df_Europe["IMF_Estimate"] < df_Europe["IMF_Estimate"].mean()]
5 df_Europe_below_imfavg["IMF_Estimate_Mean"] = df_Europe["IMF_Estimate"].mean()
6 print(df_Europe_below_imfavg[['Country/Territory', 'UN_Region', 'IMF_Estimate', 'IMF_Estimate_Mean']], "\n")
7
8 #Based on WorldBank Statistics
9 df_Europe_below_wbavg = df_Europe[df_Europe["WorldBank_Estimate"] < df_Europe["WorldBank_Estimate"].mean()]
10 df_Europe_below_wbavg["WorldBank_Estimate_Mean"] = df_Europe["WorldBank_Estimate"].mean()
11 print(df_Europe_below_wbavg[['Country/Territory', 'UN_Region', 'WorldBank_Estimate', 'WorldBank_Estimate_Mean']], "\n")
12
13 #Based on UN Statistics
14 df_Europe_below_unavg = df_Europe[df_Europe["UN_Estimate"] < df_Europe["UN_Estimate"].mean()]
15 df_Europe_below_unavg["UN_Estimate_Mean"] = df_Europe["UN_Estimate"].mean()
16 print(df_Europe_below_unavg[['Country/Territory', 'UN_Region', 'UN_Estimate', 'UN_Estimate_Mean']])
```

```
Country/Territory UN_Region IMF_Estimate IMF_Estimate_Mean
1 Monaco Europe 0 34446.75
2 Liechtenstein Europe 0 34446.75
9 Isle of Man Europe 0 34446.75
14 Channel Islands Europe 0 34446.75
15 Faroe Islands Europe 0 34446.75
51 Slovenia Europe 32214 34446.75
52 Czech Republic Europe 31368 34446.75
53 Spain Europe 31223 34446.75
54 Estonia Europe 31209 34446.75
57 Lithuania Europe 28094 34446.75
59 Portugal Europe 26012 34446.75
60 Latvia Europe 25136 34446.75
62 Slovakia Europe 23457 34446.75
```



```

1 ## Which countries in Europe has higher GDP than UK?
2 df_uk = df_Europe[df_Europe["Country/Territory"] == 'United Kingdom']
3 #Based on IMF Statistics
4 uk_imf = df_uk['IMF_Estimate'].values[0]
5 df_country_aboveuk_imf = df_Europe[df_Europe["IMF_Estimate"]>uk_imf]
6 df_country_aboveuk_imf["UK_IMF_Estimate"] = uk_imf
7 print("Based on IMF Statistics,countries in Europe has higher GDP than UK are shown as follow:")
8 print(df_country_aboveuk_imf[['Country/Territory','UN_Region','IMF_Estimate','UK_IMF_Estimate']],"\\n")
9 print("-----")
10
11 #Based on WorldBank Statistics
12 uk_wb = df_uk['WorldBank_Estimate'].values[0]
13 df_country_aboveuk_wb = df_Europe[df_Europe["WorldBank_Estimate"]>uk_wb]
14 df_country_aboveuk_wb["UK_WorldBank_Estimate"] = uk_wb
15 print("Based on WorldBank Statistics,countries in Europe has higher GDP than UK are shown as follow:")
16 print(df_country_aboveuk_wb[['Country/Territory','UN_Region','WorldBank_Estimate','UK_WorldBank_Estimate']],"\\n")
17 print("-----")
18
19 #Based on UN Statistics
20 uk_un = df_uk['UN_Estimate'].values[0]
21 df_country_aboveuk_un = df_Europe[df_Europe["UN_Estimate"]>uk_un]
22 df_country_aboveuk_un["UK_UN_Estimate"] = uk_un
23 print("Based on UN Statistics,countries in Europe has higher GDP than UK are shown as follow:")
24 print(df_country_aboveuk_un[['Country/Territory','UN_Region','UN_Estimate','UK_UN_Estimate']],"\\n")
25 print("-----")

```

Based on IMF Statistics,countries in Europe has higher GDP than UK are shown as follow:

	Country/Territory	UN_Region	IMF_Estimate	UK_IMF_Estimate
3	Luxembourg	Europe	132372	46371
4	Ireland	Europe	114581	46371
6	Norway	Europe	101103	46371
7	Switzerland	Europe	99767	46371

## Which countries below average by IMF world estimate?

```

1 df_countries_below_imf = df[df['IMF_Estimate']<df['IMF_Estimate'].mean()]
2 df_countries_below_imf["IMF_world_avg_Estimate"] = df['IMF_Estimate'].mean()
3 df_countries_below_imf[['Country/Territory','IMF_Estimate','IMF_world_avg_Estimate']]

```

C:\Users\cxy\AppData\Local\Temp\ipykernel\_35076\2372232952.py:2: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)  
df\_countries\_below\_imf["IMF\_world\_avg\_Estimate"] = df['IMF\_Estimate'].mean()

	Country/Territory	IMF_Estimate	IMF_world_avg_Estimate
1	Monaco	0	15351.632287
2	Liechtenstein	0	15351.632287
5	Bermuda	0	15351.632287
9	Isle of Man	0	15351.632287
10	Cayman Islands	0	15351.632287
...	...	...	...
219	Malawi	496	15351.632287
220	South Sudan	467	15351.632287
221	Sierra Leone	415	15351.632287

## IMF estimate 0 values

```

1 df_imf_0 = df[(df['IMF_Estimate'] == 0) | (df['IMF_Year'] == 0)]
2 df_imf_0[['Country/Territory','UN_Region','IMF_Estimate','IMF_Year']]

```

	Country/Territory	UN_Region	IMF_Estimate	IMF_Year
5	Bermuda	Americas	0	0
9	Isle of Man	Europe	0	0
10	Cayman Islands	Americas	0	0
14	Channel Islands	Europe	0	0
15	Faroe Islands	Europe	0	0
19	Greenland	Americas	0	0
31	British Virgin Islands	Americas	0	0
37	US Virgin Islands	Americas	0	0
39	New Caledonia	Oceania	0	0
42	Guam	Oceania	0	0
58	Sint Maarten (Dutch part)	Americas	0	0
61	Northern Mariana Islands	Oceania	0	0
65	Saint Martin (French part)	Americas	0	0



## Which country has highest UN Estimate?

```
: 1 highest_un = df.loc[df['UN_Estimate'].idxmax()]
 2
 3 # Display the country and the estimate
 4 print("Country with highest UN Estimate:", highest_un['Country/Territory'])
 5 print("UN Estimate value:", highest_un['UN_Estimate'])
```

Country with highest UN Estimate: Monaco  
UN Estimate value: 234317

## Which country has highest Worlbank Estimate?

```
1 highest_un = df.loc[df['WorldBank_Estimate'].idxmax()]
2
3 # Display the country and the estimate
4 print("Country with highest WorldBank Estimate:", highest_un['Country/Territory'])
5 print("WorldBank Estimate value:", highest_un['WorldBank_Estimate'])
```

Country with highest WorldBank Estimate: Monaco  
WorldBank Estimate value: 234316

## Which country has highest IMF Estimate?

```
: 1 highest_un = df.loc[df['IMF_Estimate'].idxmax()]
 2
 3 # Display the country and the estimate
 4 print("Country with highest IMF Estimate:", highest_un['Country/Territory'])
 5 print("IMF Estimate value:", highest_un['IMF_Estimate'])
```

Country with highest IMF Estimate: Luxembourg  
IMF Estimate value: 132372

```
1 # replace 0 with null values
2
3 import numpy as np
4
5 df.replace(0, np.nan, inplace=True)
6 df
```

	Country/Territory	UN_Region	IMF_Estimate	IMF_Year	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year	UN_Year_Clean
1	Monaco	Europe	NaN	NaN	234316.0	2021.0	234317.0	2021	2021
2	Liechtenstein	Europe	NaN	NaN	157755.0	2020.0	169260.0	2021	2021
3	Luxembourg	Europe	132372.0	2023.0	133590.0	2021.0	133745.0	2021	2021
4	Ireland	Europe	114581.0	2023.0	100172.0	2021.0	101109.0	2021	2021
5	Bermuda	Americas	NaN	NaN	114090.0	2021.0	112653.0	2021	2021
...	...	...	...	...	...	...	...	...	...
219	Malawi	Africa	496.0	2023.0	635.0	2021.0	613.0	2021	2021
220	South Sudan	Africa	467.0	2023.0	1072.0	2015.0	400.0	2021	2021
221	Sierra Leone	Africa	415.0	2023.0	480.0	2021.0	505.0	2021	2021
222	Afghanistan	Asia	611.0	2020.0	369.0	2021.0	373.0	2021	2021
223	Burundi	Africa	249.0	2023.0	222.0	2021.0	311.0	2021	2021

223 rows × 9 columns

```
1 # Calculate the average of 'Worldbank_Estimate' and 'UN_Estimate' columns
2 avg_Worldbank_Estimate = df['WorldBank_Estimate'].mean()
3 print("The average of 'WorldBank_Estimate' column is:", avg_Worldbank_Estimate)
4 avg_UN_Estimate = df['UN_Estimate'].mean()
5 print("The average of 'UN_Estimate' column is:", avg_UN_Estimate)
6
7
8 df['avg_worldbank_un'] = df[['WorldBank_Estimate', 'UN_Estimate']].mean(axis=1)
9 df
```

The average of 'WorldBank\_Estimate' column is: 19540.80555555555

The average of 'UN\_Estimate' column is: 18514.528037383177

	Country/Territory	UN_Region	IMF_Estimate	IMF_Year	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year	UN_Year_Clean	avg_worldbank_un
1	Monaco	Europe	17377.736041	2022	234316.0	2021.0	234317.0	2021	2021	234316.5
2	Liechtenstein	Europe	17377.736041	2022	157755.0	2020.0	169260.0	2021	2021	163507.5
3	Luxembourg	Europe	132372.000000	2023	133590.0	2021.0	133745.0	2021	2021	133667.5
4	Ireland	Europe	114581.000000	2023	100172.0	2021.0	101109.0	2021	2021	100640.5
5	Bermuda	Americas	17377.736041	2022	114090.0	2021.0	112653.0	2021	2021	113371.5
...	...	...	...	...	...	...	...	...	...	...
219	Malawi	Africa	496.000000	2023	635.0	2021.0	613.0	2021	2021	624.0
220	South Sudan	Africa	467.000000	2023	1072.0	2015.0	400.0	2021	2021	736.0
221	Sierra Leone	Africa	415.000000	2023	480.0	2021.0	505.0	2021	2021	492.5



```

: 1 # Fill the null values in 'imf' column with the calculated average
 2 df['IMF_Estimate'] = df['IMF_Estimate'].fillna(df['IMF_Estimate'].mean())
 3 imf_mean_year = round(df['IMF_Year'].mean())
 4 df['IMF_Year'] = df['IMF_Year'].fillna(imf_mean_year).astype(int).astype(str)
 5 df

```

	Country/Territory	UN_Region	IMF_Estimate	IMF_Year	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year	UN_Year_Clean
1	Monaco	Europe	17377.736041	2022	234316.0	2021.0	234317.0	2021	2021
2	Liechtenstein	Europe	17377.736041	2022	157755.0	2020.0	169260.0	2021	2021
3	Luxembourg	Europe	132372.000000	2023	133590.0	2021.0	133745.0	2021	2021
4	Ireland	Europe	114581.000000	2023	100172.0	2021.0	101109.0	2021	2021
5	Bermuda	Americas	17377.736041	2022	114090.0	2021.0	112653.0	2021	2021
...	...	...	...	...	...	...	...	...	...
219	Malawi	Africa	496.000000	2023	635.0	2021.0	613.0	2021	2021
220	South Sudan	Africa	467.000000	2023	1072.0	2015.0	400.0	2021	2021
221	Sierra Leone	Africa	415.000000	2023	480.0	2021.0	505.0	2021	2021
222	Afghanistan	Asia	611.000000	2020	369.0	2021.0	373.0	2021	2021
223	Burundi	Africa	249.000000	2023	222.0	2021.0	311.0	2021	2021

223 rows × 9 columns

```

1 # Drop the temporary 'avg_worldbank_un' column if not needed
2 df.drop(columns=['avg_worldbank_un'], inplace=True)
3 df

```

	Country/Territory	UN_Region	IMF_Estimate	IMF_Year	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year	UN_Year_Clean
1	Monaco	Europe	17377.736041	2022	234316.0	2021.0	234317.0	2021	2021
2	Liechtenstein	Europe	17377.736041	2022	157755.0	2020.0	169260.0	2021	2021
3	Luxembourg	Europe	132372.000000	2023	133590.0	2021.0	133745.0	2021	2021
4	Ireland	Europe	114581.000000	2023	100172.0	2021.0	101109.0	2021	2021
5	Bermuda	Americas	17377.736041	2022	114090.0	2021.0	112653.0	2021	2021
...	...	...	...	...	...	...	...	...	...
219	Malawi	Africa	496.000000	2023	635.0	2021.0	613.0	2021	2021
220	South Sudan	Africa	467.000000	2023	1072.0	2015.0	400.0	2021	2021
221	Sierra Leone	Africa	415.000000	2023	480.0	2021.0	505.0	2021	2021
222	Afghanistan	Asia	611.000000	2020	369.0	2021.0	373.0	2021	2021
223	Burundi	Africa	249.000000	2023	222.0	2021.0	311.0	2021	2021

223 rows × 9 columns

## Checking Missing Values

```

]: 1 print(df.isnull())
 2 print(df.isnull().sum())

```

```

Country/Territory UN_Region IMF_Estimate IMF_Year WorldBank_Estimate \
1 False False False False False
2 False False False False False
3 False False False False False
4 False False False False False
5 False False False False False
...
219 ... ... ... ...
220 False False False False False
221 False False False False False
222 False False False False False
223 False False False False False

WorldBank_Year UN_Estimate UN_Year UN_Year_Clean
1 False False False False
2 False False False False
3 False False False False
4 False False False False
5 False False False False
...
219 ... ...
220 False False False False
221 False False False False
222 False False False False
223 False False False False

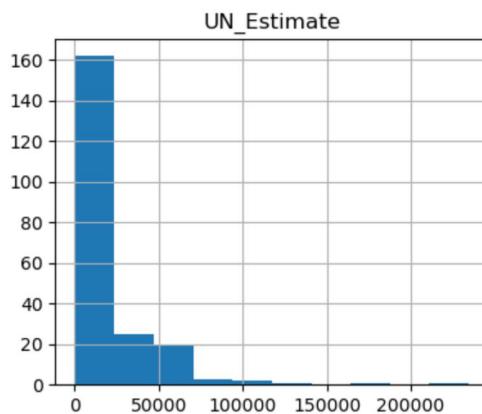
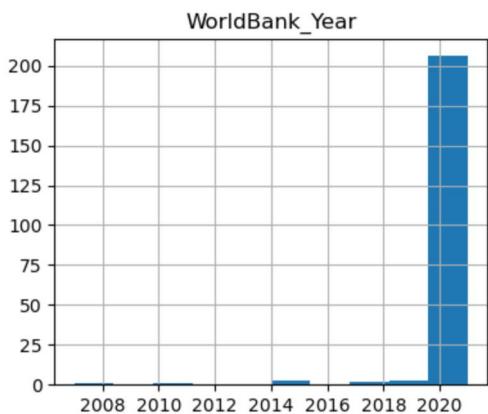
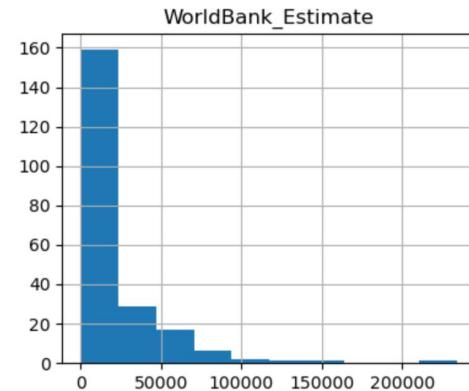
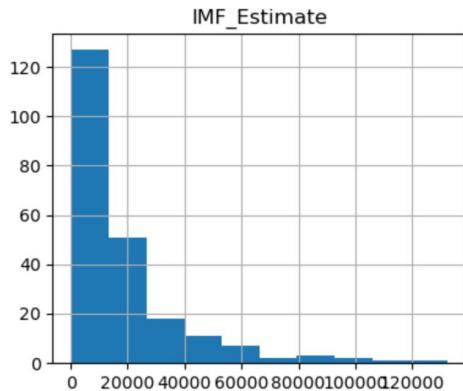
```



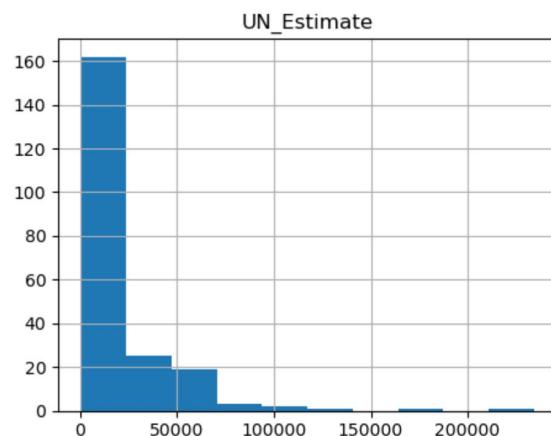
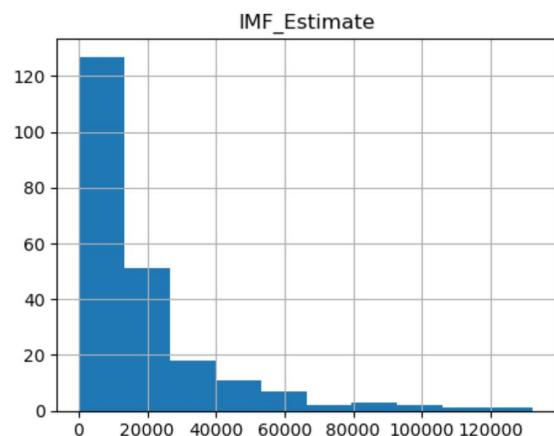
```
[223 rows x 9 columns]
Country/Territory      0
UN_Region              0
IMF_Estimate           0
IMF_Year                0
WorldBank_Estimate      7
WorldBank_Year           7
UN_Estimate             9
UN_Year                  0
UN_Year_Clean            0
dtype: int64
```

### Histogram

```
: 1 df.hist(figsize=(10,8))
2 plt.show()
```



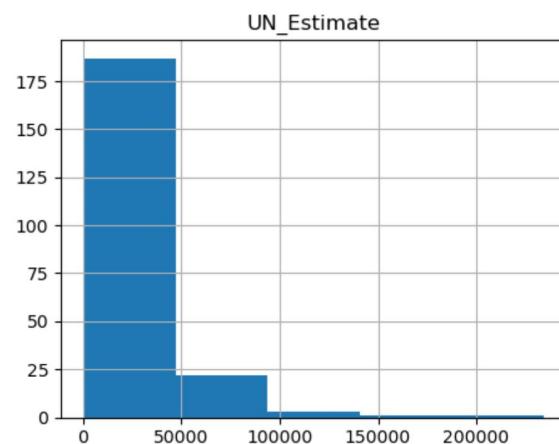
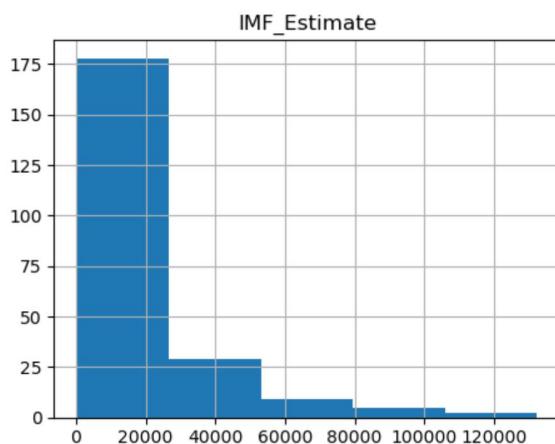
```
: 1 df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].hist(figsize=(12,9))
2
3 plt.show()
```



```

1 df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].hist(bins=5, figsize=(12,9))
2
3 plt.show()

```



```

1 df[["WorldBank_Estimate"]].agg(["min", "max"])
2
3 plt.show()

```

```

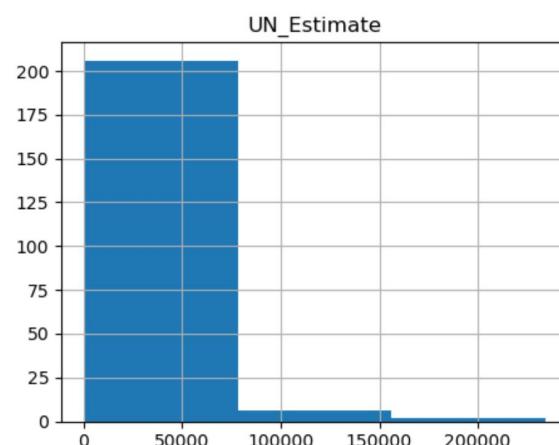
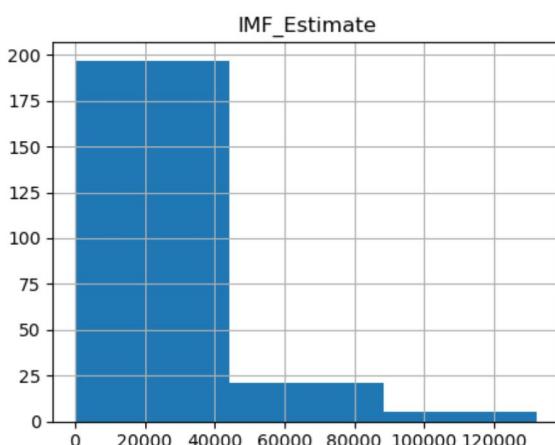
min    222.0
max  234316.0
Name: WorldBank_Estimate, dtype: float64

```

```

1 df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].hist(bins=3, figsize=(12,9))
2
3 plt.show()

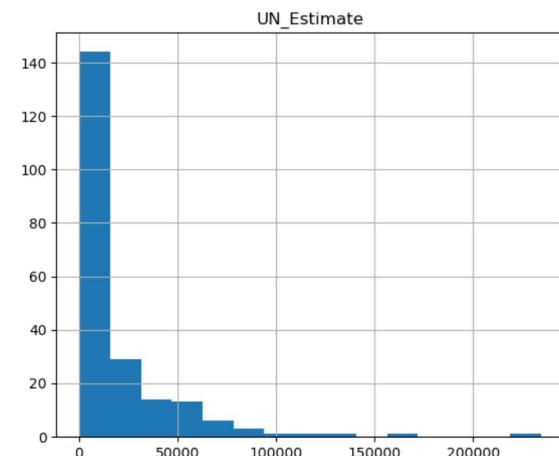
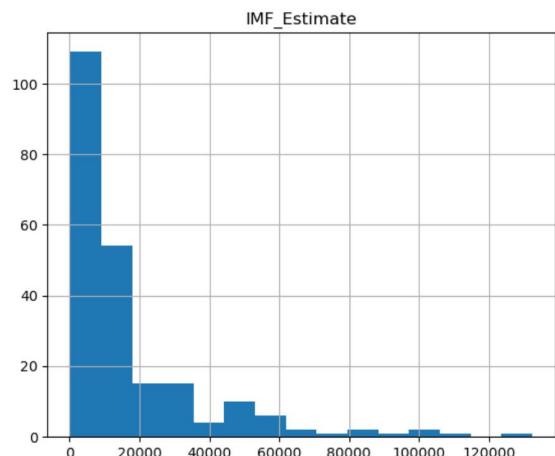
```



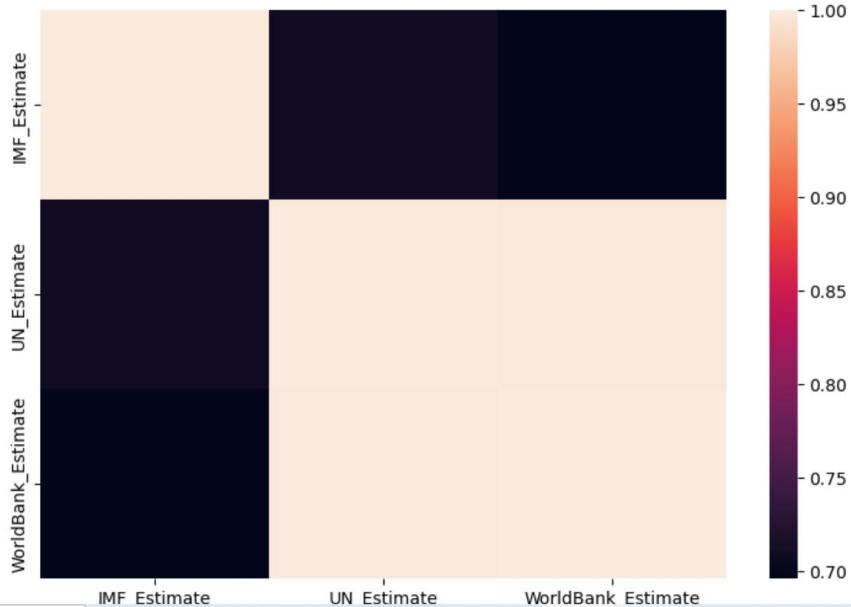
```

1 df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].hist(bins=15, figsize=(15,12))
2
3 #23400/15 = 15300
4 plt.show()

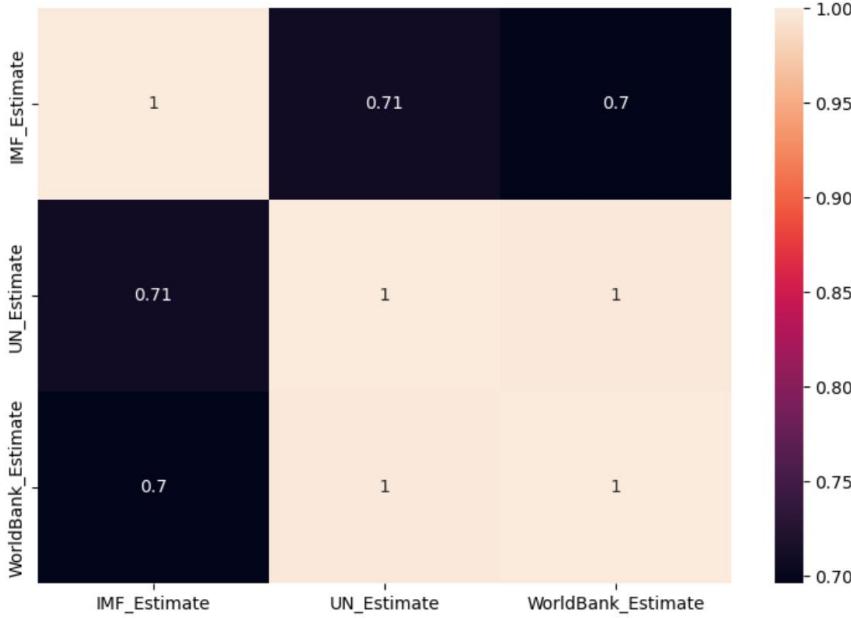
```



```
: 1 corr = df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].corr()
 2
 3 plt.figure(figsize=(9,6))
 4 sns.heatmap(corr)
 5
 6 plt.show()
```



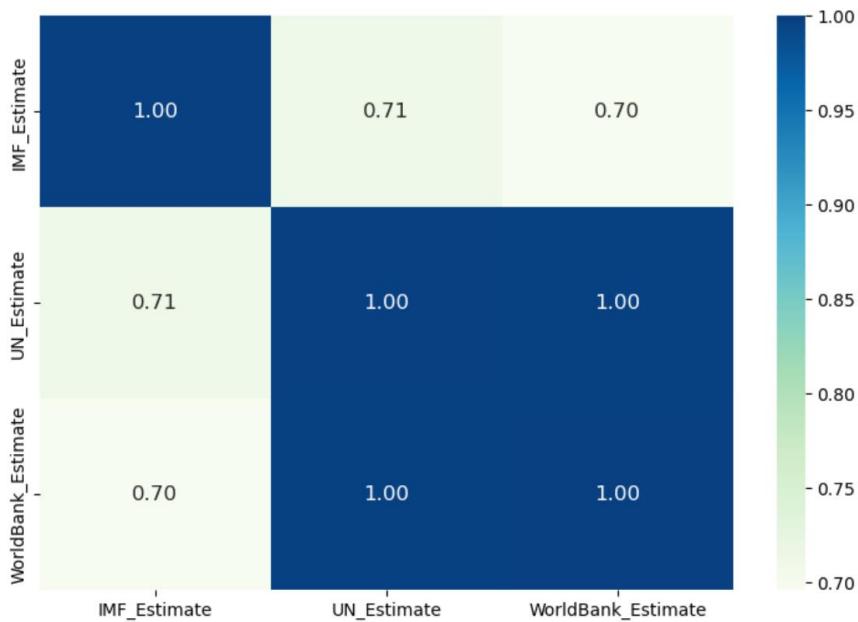
```
1 corr = df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].corr()
2
3 plt.figure(figsize=(9,6))
4
5 sns.heatmap(corr, annot=True)
6
7 plt.show()
```



```

1 corr = df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].corr()
2
3 plt.figure(figsize=(9,6))
4
5 sns.heatmap(corr, annot=True, fmt=".2f", cmap = 'GnBu', annot_kws={"size": 12})
6
7 plt.show()

```



```

1 corr = df[["IMF_Estimate", "UN_Estimate", "WorldBank_Estimate"]].corr()
2
3 plt.figure(figsize=(9,6))
4
5 sns.heatmap(corr, annot=True, cmap = 'Purples')
6
7 plt.title("Correlation Map")
8
9
10 plt.show()

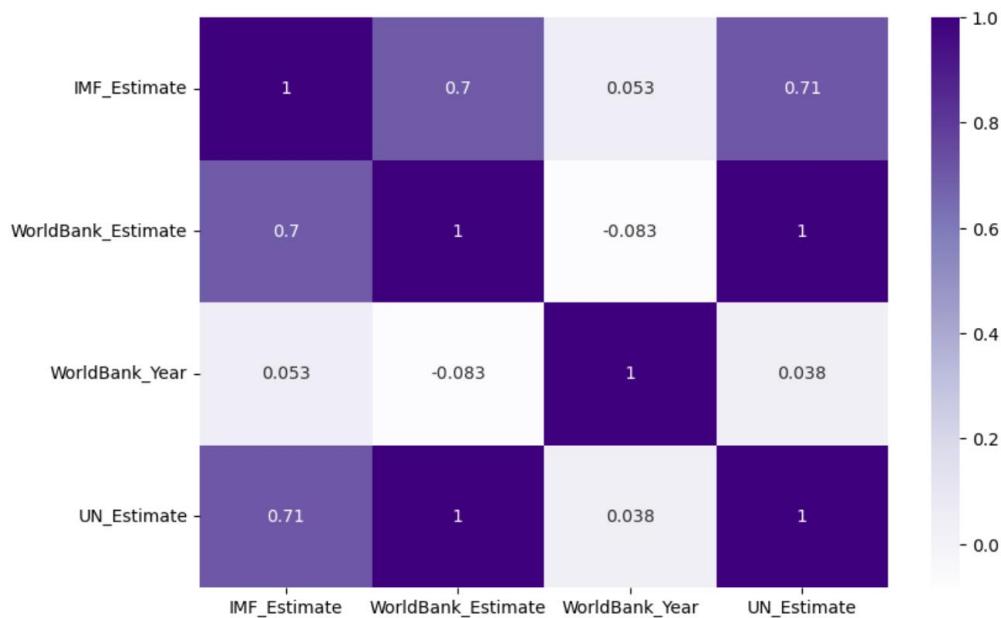
```



```

1 corr = df.select_dtypes(include=[int, float]).corr()
2
3 plt.figure(figsize=(9,6))
4
5 sns.heatmap(corr, annot=True, cmap = 'Purples')
6
7 plt.show()

```



```

1 df.head()

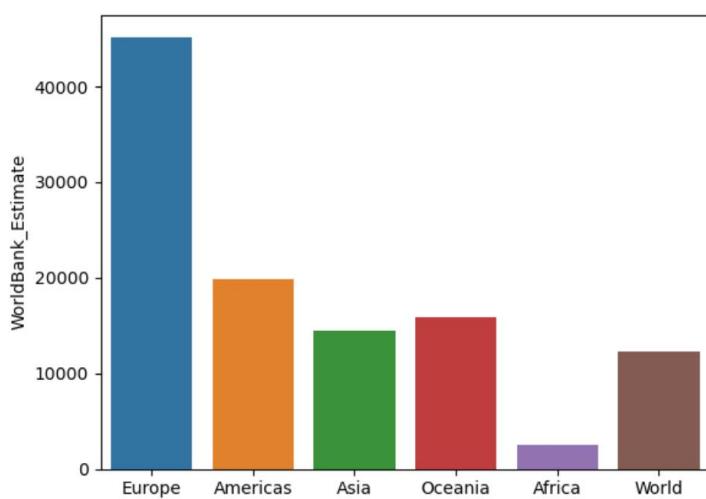
```

	Country/Territory	UN_Region	IMF_Estimate	IMF_Year	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year	UN_Year_Clean
1	Monaco	Europe	17377.736041	2022	234316.0	2021.0	234317.0	2021	2021
2	Liechtenstein	Europe	17377.736041	2022	157755.0	2020.0	169260.0	2021	2021
3	Luxembourg	Europe	132372.000000	2023	133590.0	2021.0	133745.0	2021	2021
4	Ireland	Europe	114581.000000	2023	100172.0	2021.0	101109.0	2021	2021
5	Bermuda	Americas	17377.736041	2022	114090.0	2021.0	112653.0	2021	2021

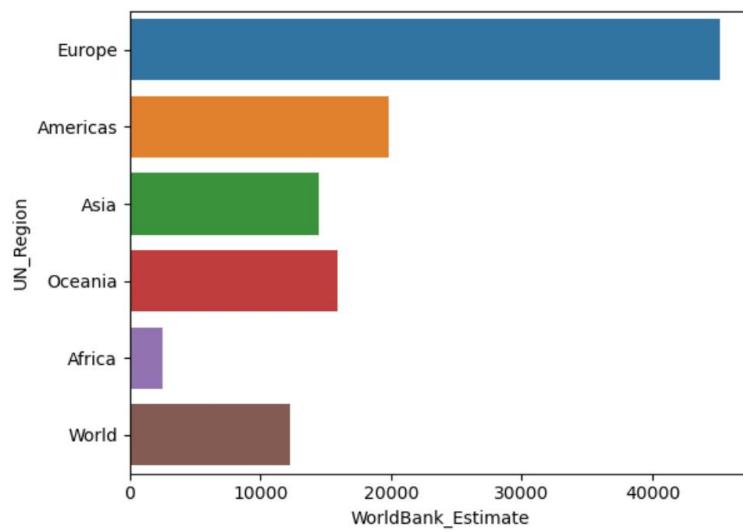
```

1 sns.barplot(x="UN_Region", y="WorldBank_Estimate", data=df, errorbar=None)
2
3 plt.show()

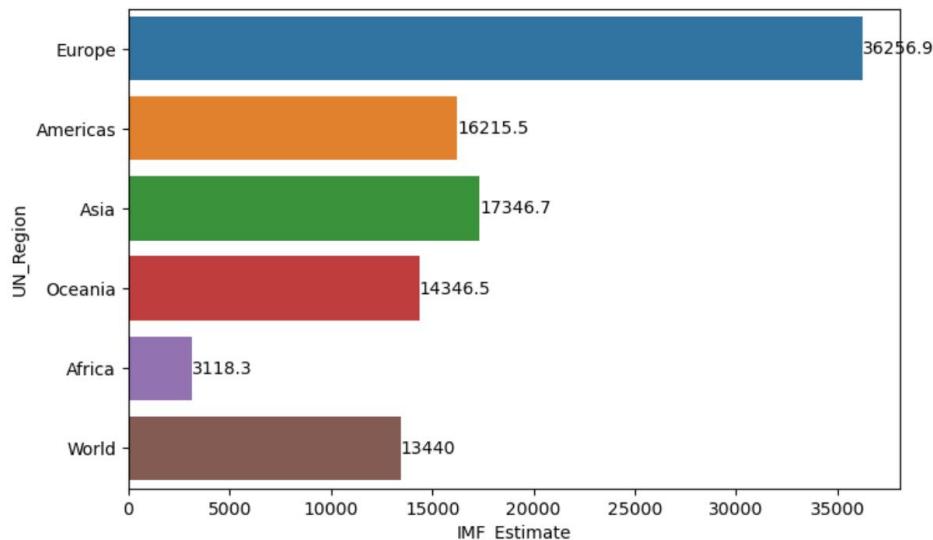
```



```
1 sns.barplot(x="WorldBank_Estimate", y="UN_Region", data=df, errorbar=None)
2
3 plt.show()
```



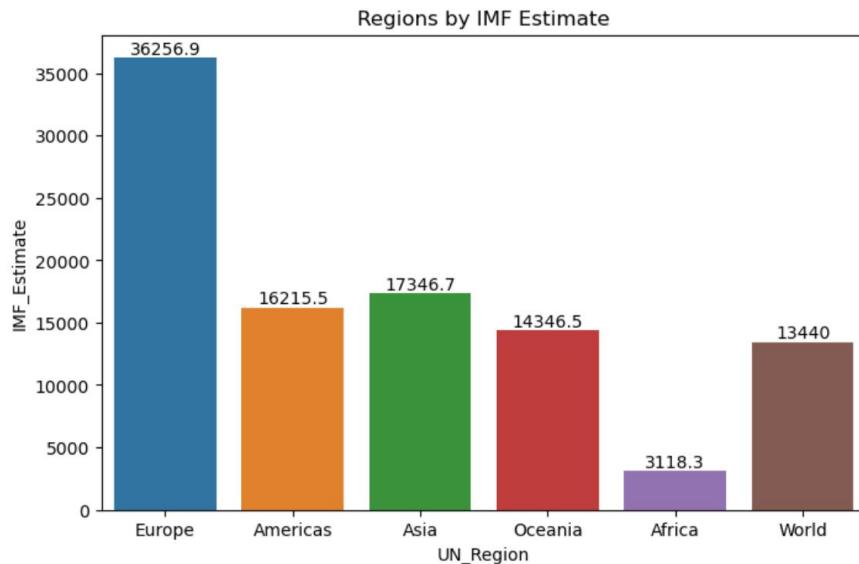
```
1 fig = plt.figure(figsize = (8,5))
2
3 ax = sns.barplot(x = "IMF_Estimate", y = "UN_Region",
4 data = df, errorbar = None)
5
6 ax.bar_label(ax.containers[0])
7
8 plt.show()
```



```

1 fig = plt.figure(figsize = (8,5))
2 ax = sns.barplot(x = "UN_Region", y = "IMF_Estimate",
3                   data = df, errorbar = None)
4
5 ax.bar_label(ax.containers[0])
6
7
8 ax.set_title("Regions by IMF Estimate")
9 plt.show()

```

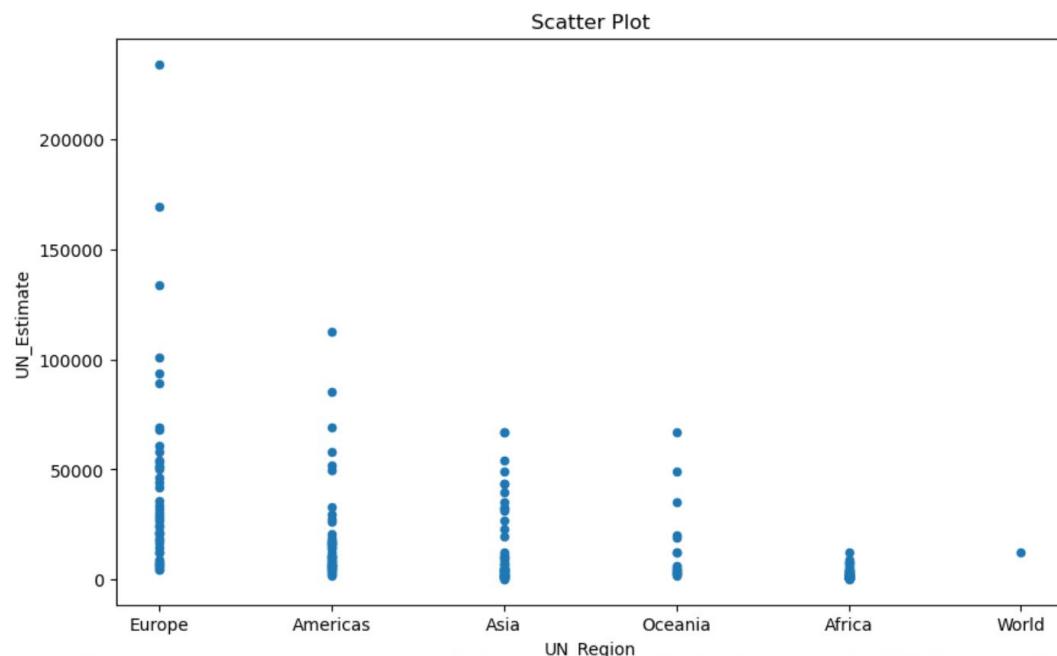


### Scatter Plot

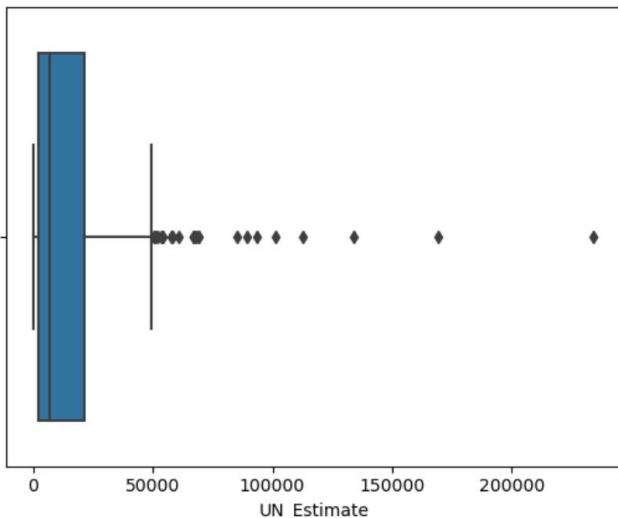
```

1 df.plot(x='UN_Region', y='UN_Estimate', kind='scatter',
2         figsize=(10,6),
3         title="Scatter Plot")
4
5 plt.show()

```



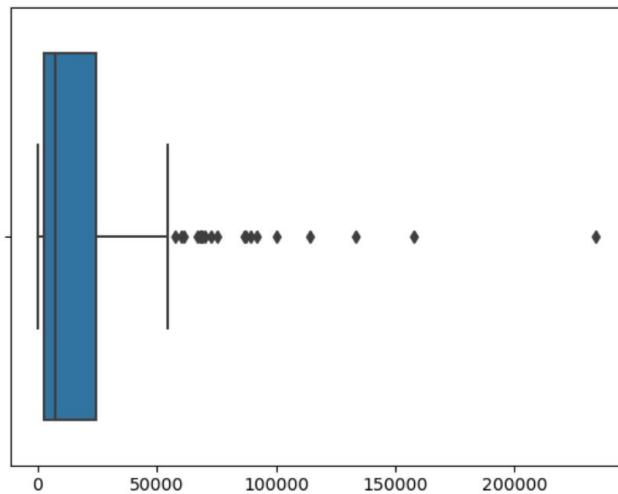
```
1 sns.boxplot(x=df["UN_Estimate"])
2
3 plt.show()
```



```
: 1 df[df["UN_Estimate"]>50000].head()
```

```
:  
Country/Territory UN_Region IMF_Estimate IMF_Year WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year UN_Year_Clean  
1 Monaco Europe 17377.736041 2022 234316.0 2021.0 234317.0 2021 2021  
2 Liechtenstein Europe 17377.736041 2022 157755.0 2020.0 169260.0 2021 2021  
3 Luxembourg Europe 132372.000000 2023 133590.0 2021.0 133745.0 2021 2021  
4 Ireland Europe 114581.000000 2023 100172.0 2021.0 101109.0 2021 2021  
5 Bermuda Americas 17377.736041 2022 114090.0 2021.0 112653.0 2021 2021
```

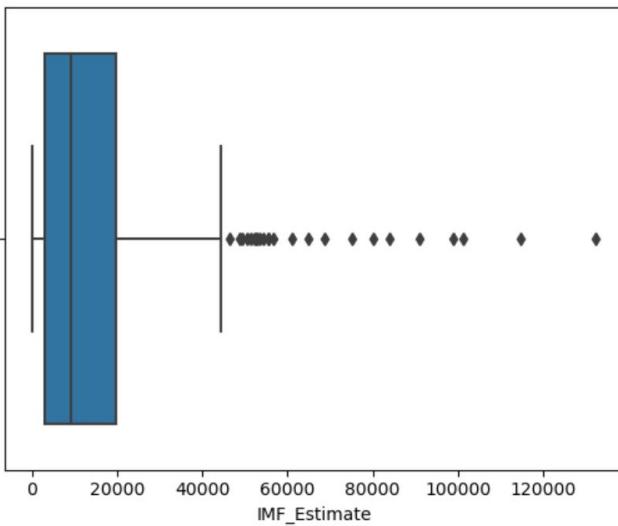
```
: 1 sns.boxplot(x=df["WorldBank_Estimate"])
2
3 plt.show()
```



```

1 sns.boxplot(x=df["IMF_Estimate"])
2
3 plt.show()

```



```
: 1 df[df["UN_Estimate"]>100000]
```

```
:
   Country/Territory UN_Region IMF_Estimate IMF_Year WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year UN_Year_Clean
1      Monaco        Europe    17377.736041    2022     234316.0      2021.0   234317.0    2021       2021
2 Liechtenstein      Europe    17377.736041    2022     157755.0      2020.0   169260.0    2021       2021
3 Luxembourg        Europe  132372.000000    2023     133590.0      2021.0   133745.0    2021       2021
4      Ireland        Europe  114581.000000    2023     100172.0      2021.0   101109.0    2021       2021
5     Bermuda      Americas  17377.736041    2022     114090.0      2021.0   112653.0    2021       2021
```

```
: 1 df.UN_Estimate.mean()
```

```
: 18514.528037383177
```

```
: 1 df.shape
```

```
: (223, 9)
```

### Create another dataframe called data excluding 5 countries with highest UN estimate

```
1 data = df[-(df["UN_Estimate"]>100000)]
```

```
1 data.head()
```

```
:
   Country/Territory UN_Region IMF_Estimate IMF_Year WorldBank_Estimate WorldBank_Year UN_Estimate UN_Year UN_Year_Clean
6      Norway        Europe  101103.000000    2023     89154.0      2021.0   89242.0    2021       2021
7 Switzerland      Europe  98767.000000    2023     91992.0      2021.0   93525.0    2021       2021
8     Singapore       Asia  91100.000000    2023     72794.0      2021.0   66822.0    2021       2021
9 Isle of Man        Europe  17377.736041    2022     87158.0      2019.0      NaN        0         0
10 Cayman Islands    Americas 17377.736041    2022     86569.0      2021.0   85250.0    2021       2021
```

```
1 data.shape
```

```
(218, 9)
```

```
1 data.UN_Estimate.mean()
```

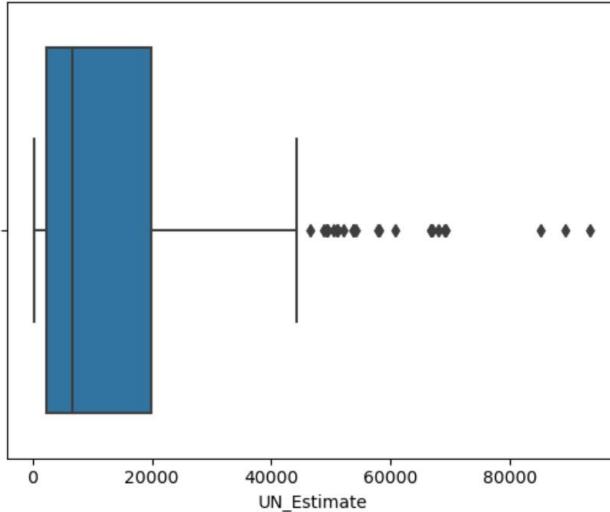
```
15363.755980861244
```

```
1 df.UN_Estimate.mean()
```

```
18514.528037383177
```



```
: 1 sns.boxplot(x=data["UN_Estimate"])
 2 plt.show()
```



```
1 lower_q = df["UN_Estimate"].quantile(0.25)
2 lower_q
```

2331.75

```
1 higher_q = df["UN_Estimate"].quantile(0.75)
2 higher_q
```

21359.25

```
1 iqr = higher_q - lower_q
2 iqr
```

19027.5

```
1 upper_boundary = higher_q + 1.5 * iqr
2 upper_boundary
```

49900.5

```
1 lower_boundary = lower_q - 1.5 * iqr
2 lower_boundary
```

-26209.5

```
1 df_filtered = df[(df["UN_Estimate"] < upper_boundary) & (df["UN_Estimate"] > lower_boundary)]
```

```
1 df_filtered.head()
```

Country/Territory	UN_Region	IMF_Estimate	IMF_Year	WorldBank_Estimate	WorldBank_Year	UN_Estimate	UN_Year	UN_Year_Clean
27	Hong Kong	52429.000000	2023	49801.0	2021.0	49259.0	2021	2021
29	Macau	50571.000000	2023	43874.0	2021.0	43555.0	2021	2021
30	United Arab Emirates	49451.000000	2023	44316.0	2021.0	43295.0	2021	2021
31	British Virgin Islands	17377.736041	2022	Nan	Nan	49444.0	2021	2021
32	New Zealand	48826.000000	2023	48781.0	2021.0	48824.0	2021	2021

```
1 df_filtered.shape
2 # there were 223 rows - 196 = 27 outliers dropped
```

(190, 9)

```
1 df_filtered.UN_Estimate.mean()
```

10488.947368421053

```
1 df.UN_Estimate.mean()
```

18514.528037383177

```
1 #how can we create a table with following
2 df_filtered.WorldBank_Estimate.mean()
```

10355.304347826086



```
1 df.WorldBank_Estimate.mean()
19540.805555555555
1 df_filtered.IMF_Estimate.mean()
11636.771413304836
1 df.IMF_Estimate.mean()
17377.736040609136
```

## 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.**

