

# Exercise 1 - Practical Part - Consolidated Report

Eliav Dayanof, 208674556  
Yair Ben Yakar, 319013090

May 21, 2025

## 1 Part 1 – Data Aquisition

### 1.1 Preview

- Preview of the datasets:
  - The first 10 rows of the demographics dataset:

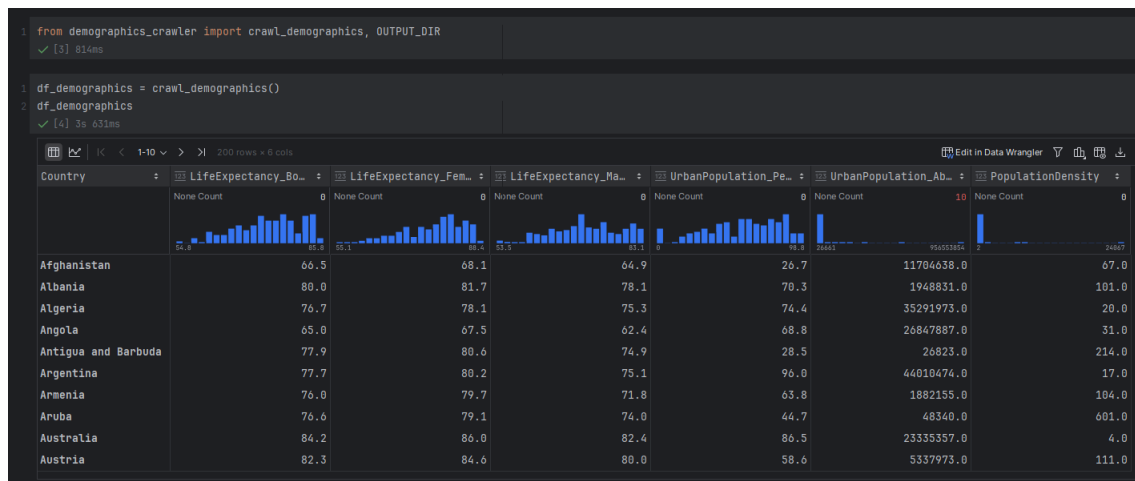


Figure 1: The first 10 rows of demographics data (before and after sorting, stays the same)

- The first 5 rows of the GDP dataset:

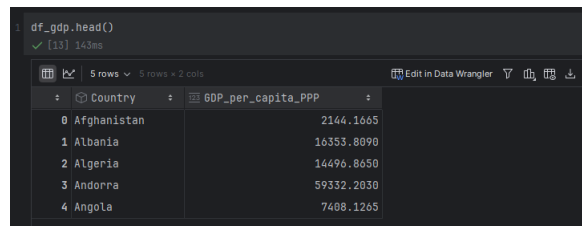


Figure 2: First 5 rows of GDP data (before and after sorting, stays the same)

- The first 5 rows of the Population dataset:

```
1 df_pop.head()
```

✓ [14] 153ms

5 rows × 2 cols

	Country	Population
0	Afghanistan	40000360
1	Africa	1413750475
2	Africa (UN)	1413753005
3	Albania	2849591
4	Algeria	44761051

Figure 3: First 5 rows of Population data (before and after sorting, stays the same)

- Summary statistics:

```
1 df_gdp.describe()
```

✓ [96]

8 rows × 1 cols

	GDP_per_capita_PPP
count	213.000000
mean	25822.084541
std	25794.941595
min	836.665600
25%	6214.017000
50%	16353.809000
75%	38862.090000
max	137947.340000

Figure 4: Summary statistics table for GDP

```
1 df_pop.describe()
```

✓ [98]

8 rows × 1 cols

	Population
count	2.000000e+02
mean	1.687752e+08
std	7.254974e+08
min	5.150000e+02
25%	5.224038e+05
50%	6.827910e+06
75%	3.425334e+07
max	7.954448e+09

Figure 5: Summary statistics table for Population

- Printed confirmation of DataFrame shapes and column names:

```
1 print(df_demographics.shape, list(df_demographics.columns), sep='\n')
```

✓ [19] 125ms

(200, 6)

['LifeExpectancy\_Both', 'LifeExpectancy\_Female', 'LifeExpectancy\_Male', 'UrbanPopulation\_Percentage', 'UrbanPopulation\_Absolute', 'PopulationDensity']

Figure 6: Demographics table shape and columns

```

1 gdp_dataset_expected_cols = ["Country", "GDP_per_capita_PPP"]
2 confirm_cols(df_gdp, gdp_dataset_expected_cols)
3 print(df_gdp.shape, df_gdp.columns)
✓ [16] 136ms

['Country', 'GDP_per_capita_PPP'] cols exist in dataframe
(213, 2) Index(['Country', 'GDP_per_capita_PPP'], dtype='object')

```

Figure 7: GDP table shape and columns

```

1 pop_dataset_expected_cols = ["Country", "Population"]
2 confirm_cols(df_pop, pop_dataset_expected_cols)
3 print(df_pop.shape, df_pop.columns)
✓ [15] 144ms

['Country', 'Population'] cols exist in dataframe
(268, 2) Index(['Country', 'Population'], dtype='object')

```

Figure 8: Population table shape and columns

## 1.2 Demographics Data Analysis

- Neumeric-fields Summary:

```

num_cols = df_demographics.select_dtypes("number").columns
summary = (
    df_demographics[num_cols]
    .agg(["mean", "std", "min", "max", "median"])
    .T
    .assign(missing=df_demographics[num_cols].isna().sum())
)
display(summary) # nice HTML in notebook
[100]

```

	mean	std	min	max	median	missing
LifeExpectancy_Both	7.397750e+01	7.004589e+00	54.8	85.8	74.85	0
LifeExpectancy_Female	7.659650e+01	7.126054e+00	55.1	88.4	78.05	0
LifeExpectancy_Male	7.137950e+01	6.979282e+00	53.5	83.1	71.20	0
UrbanPopulation_Percentage	5.727600e+01	2.514603e+01	0.0	98.8	60.55	0
UrbanPopulation_Absolute	2.483410e+07	8.517914e+07	26661.0	956553854.0	501189.50	10
PopulationDensity	3.768200e+02	1.866303e+03	2.0	24067.0	93.00	0

Figure 9: Demographics Data Neumeric-fields Summary

- Pearson correlation:

```

1 corr = df_demographics["LifeExpectancy_Both"].corr(
2     df_demographics["PopulationDensity"]
3 )
4 print(f"Pearson correlation (LifeExpectancy Both vs PopulationDensity): {corr:.4f}")
[102]

Pearson correlation (LifeExpectancy Both vs PopulationDensity): 0.1796

```

Figure 10: Pearson correlation of LifeExpectancy Both vs PopulationDensity

## 2 Part 2 – Cleaning Summary

### 2.1 Demographics Dataset

Issues encountered → Actions taken

- **Non-numeric values in columns** → All columns except "Country" were explicitly converted to numeric using `pd.to_numeric(errors='coerce')`.

- **Invalid Life Expectancy values** ( $< 40$  or  $> 100$ ) → Rows outside the valid range [40, 100] for LifeExpectancy\_Both were removed.
- **Inconsistent country names** (e.g., "the Gambia") → Standardized with: strip(), title(), and removed "the" prefix.
- **Potential name mismatches after standardization** → Documented in output/name\_mismatches.csv.

## Row counts

- Rows before cleaning: 200
- Rows after cleaning: 200

### 2.1.1 GDP per Capita Dataset

Issues encountered → Actions taken

- **Missing or malformed GDP values** → GDP\_per\_capita\_PPP was coerced to numeric, invalid/missing rows were dropped and logged to output/dropped\_gdp.csv (in reality there were none).
- **Outliers in GDP values** → Identified (but not removed) using Tukey's method, outlier count of 6 printed to console:

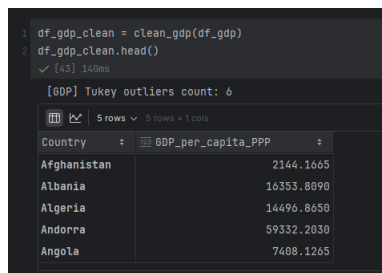


Figure 11: GDP After Cleaning – Outlier Count

- **Duplicate country entries** → Duplicates dropped, keeping the first occurrence per "Country" (in reality there were no duplicates, if there were, perhaps we would've chosen a smarter strategy).
- **Country name mismatches** → Mapped manually to match the Demographics dataset (e.g., "Cape Verde" → "Cabo Verde").
  - Note: We experimented with trying to map mismatches using fuzzy search techniques to find words with similar spellings and match the Demographics dataset (assuming it as ground truth) spelling (post standardization), but it didn't work out well so we resorted to doing it manually (with sorting and some help of code of course).
  - Also, another approach we considered was using edit distance.

## Row counts

- Rows before cleaning: 213
- Rows after cleaning: 213

## Population Dataset   Issues encountered → Actions taken

- **Non-numeric or missing values in Population** → Coerced to numeric, missing values dropped and recorded in output/dropped\_population.csv.
- **Extreme outliers due to scale differences** → Detected using log10 + Tukey's method, outlier count of 1 printed to console, but values retained:

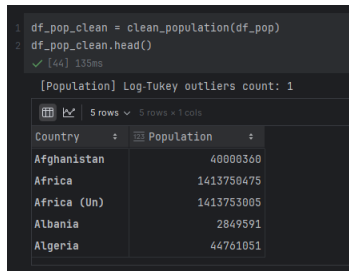


Figure 12: Population After Cleaning – Outlier Count

- **Duplicate entries by country** → Duplicates dropped (first occurrence kept).
- **Country name mismatches** → Mapped manually to match ground truth Demographics dataset (e.g., "Saint Vincent And The Grenadines" → "St. Vincent & Grenadines").

## Row counts

- Rows before cleaning: 260
- Rows after cleaning: 260

## 3 Part 3 – Feature Engineering

### Description of Transformations

- Log transformations were applied to GDP\_per\_capita\_PPP (LogGDPperCapita) and Population (LogPopulation) to compress the range and reduce skew.
- We did not log-transform LifeExpectancy\_Both, as it is already approximately normally distributed.
- Afterward, we applied z-score normalization to the following columns:
  - LifeExpectancy\_Both → LifeExpectancy\_z
  - LogGDPperCapita → LogGDPpc\_z
  - LogPopulation → LogPop\_z This produced three normalized features with zero mean and unit variance, suitable for further analysis.

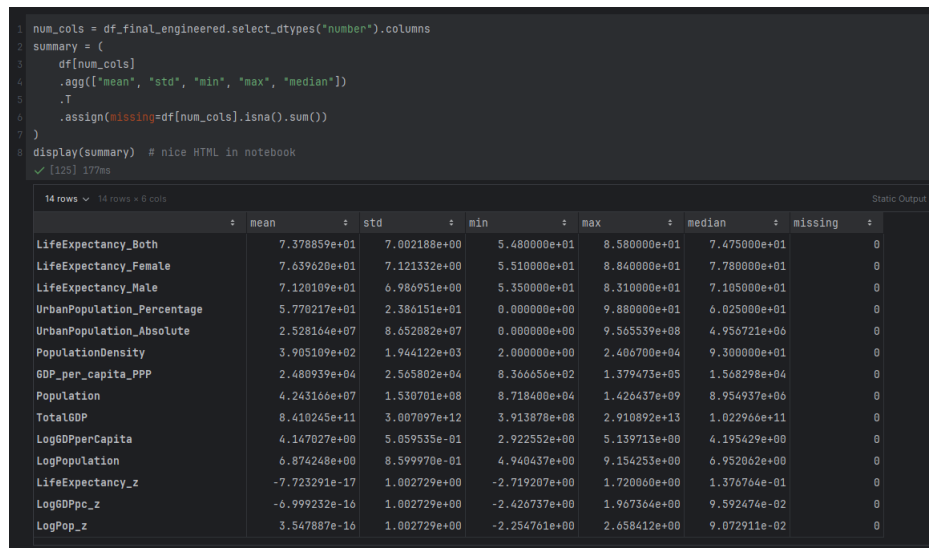


Figure 13: Updated Descriptive Statistics Table After Scaling

## Updated Descriptive Statistics Table After Scaling

## Evidence of Successful Crawling and Merging

- Final row count: 183 (as seen in the top left corner of the table below).

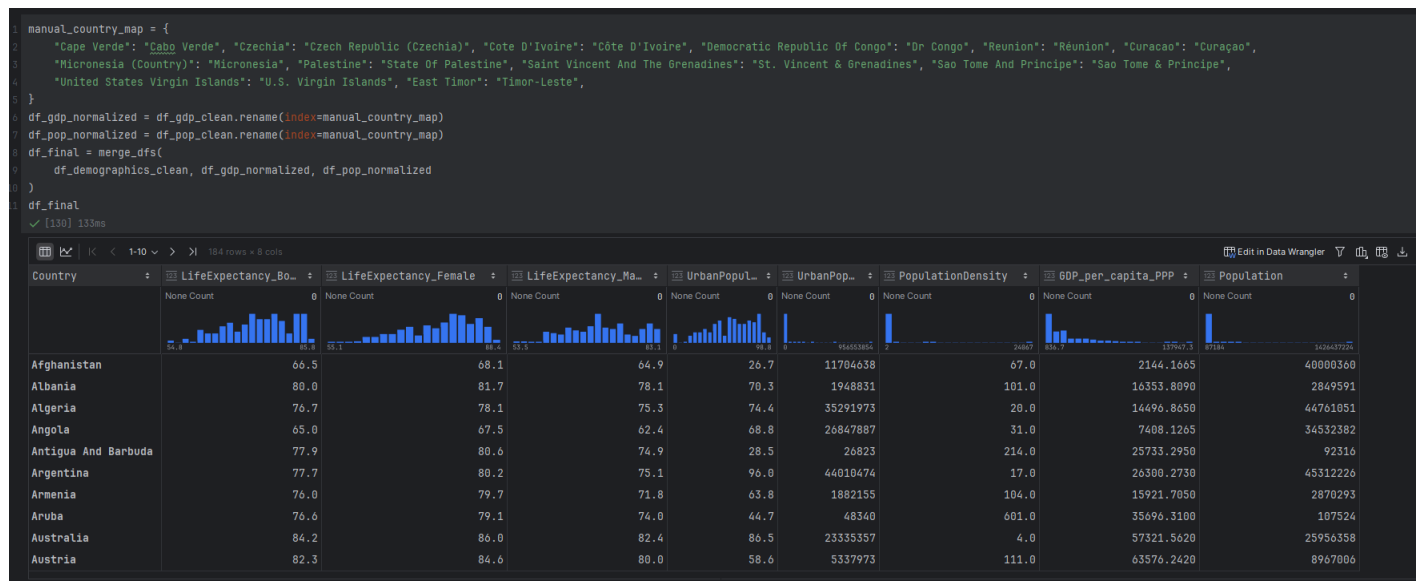


Figure 14: Final Merged Table with All Crawling Data