Life Expectancy Analysis Report

Introduction

This report presents a comprehensive analysis of global life expectancy statistics using data from the UN, World Bank Group, World Health Organization, and the CIA World Factbook. The analysis highlights life expectancy trends by age, sex, and over time, revealing disparities and improvements among countries. Additionally, the concept of Healthy Life Expectancy (HALE) is discussed.

1. Data Sources Overview

- **Table 1:** UN estimates of life expectancy at different ages (2023)
- **Table 2:** Changes in life expectancy from 2019 to 2023
- **Table 3:** World Bank Group historical trends (2014–2022)
- **Table 4:** World Health Organization (2019) life expectancy and HALE
- **Table 5:** CIA World Factbook (2022) by sex and sex gap
- **Table 6:** Life expectancy estimates in the OCED for 2022

2. Project Overview

This project focuses on analyzing life expectancy data for various countries, with emphasis on differences between males and females across time and age groups. The goal is to understand:

- Gender-based longevity gaps
- Time-based changes in life expectancy
- Trends across different geographic regions

Data will be scraped and compiled from **Wikipedia**, cleaned using Python and regular expressions, analyzed statistically, visualized using data plots, and stored in a **MongoDB** database. The project concludes with an optional Streamlit dashboard to display interactive insights.

3. Project Objectives

We aim to answer questions like:

- Which countries have the **highest** and **lowest** life expectancy?
- Where is the **male-female life expectancy gap** largest or smallest?
- How has life expectancy **changed from 2019 to 2023** in various countries?
- Are countries improving or declining in health over time?

4. Methodology

☐ Data Extraction

- Scrape tables using **BeautifulSoup** and **Pandas**.
- Export and store raw tables in Excel format.

☐ Data Cleaning and Regular Expressions

- Handle nulls, invalid or missing values.
- Normalize country names and formats.
- Use **Regex** to:
 - o Parse columns with multiple values into structured forms.

☐ Data Analysis

- Compute:
 - o Mean, median, and standard deviation of life expectancy.
 - o Gender gap by country and year.
 - o Year-over-year growth or decline.
- Compare:
 - o Continents (Asia, Europe, etc.)
 - o WHO vs World Bank vs CIA reports

☐ Data Visualization

- **Bar charts** for male vs. female life expectancy.
- **Line graphs** for year-over-year trends.
- **Heatmaps** to show global life expectancy distribution.
- **Histograms/Boxplots** for gap analysis.

☐ Data Storage

• Final processed data stored in **MongoDB** for flexible querying and scalability.

Bonus: Streamlit Web App

- A simple dashboard for exploring:
 - Top/bottom countries
 - o Male-female gap over time
 - Country-specific graphs

1. Cleaning & Data Overview phase

Sheet1:UN: Change of life expectancy from 2019 to 2023

This dataset presents the change in life expectancy across different countries between 2019 and 2023, based on United Nations estimates, highlighting trends of increase or decrease during this period.

Sheet2:World Health Organization (2019)

This dataset, published by the World Health Organization in December 2020, provides life expectancy at birth for all populations as of 2019. If countries have equal life expectancy values, they are further sorted by HALE (Healthy Life Expectancy) for the total population.

Sheet3:CIA World Factbook (2022)

This dataset, published by the CIA World Factbook in 2022, provides life expectancy data categorized by total population, male, female, and the gender gap between males and females.

Sheeet4:World Bank Group (2022)

The dataset provides World Bank Group estimates for life expectancy across different countries in 2022, focusing only on nations with populations exceeding 50,000 people. It includes comparative data for the years 2014, 2019, and 2022 to analyze trends, particularly noting that 2014 marked a local peak in life expectancy for some leading countries. Data rounding may cause minor inconsistencies (up to 0.01 year) when compared to the raw calculations. The main columns include the country name, life expectancy for the three reference years, and the population size in 2022

➤ Here are some of our cleaning process on "**sheet 4**":

```
import pandas as pd

sh4=pd.read_excel(r"C:\Users\sh138\OneDrive\Documents\sh4.xlsx",header=1)

sh4.columns= ['Countries and territories', 'All', 'Male', 'Female', 'Sex gap','2014', '2014:2019', '2019','2019:2020', '2020:2021',
```

- ✓ First, we read the data then use 'header' to handle the index and the columns.
- ✓ Secondly, we gave each column its proper name.



✓ Dropping a null column.

```
cols_to_coovert = ['All', 'Fale', 'Foosle', 'See gap', '2014', '2014:2010', '2019', '2019:2020', '2020', '2020:2021', '2021:2022', '2022', 'recovery from COVIO-18: 2019:2022']
for col in cols_to_convert:
    sh4[col] = (sh4[col],astype(str)
                .str.replace('-'.
                .str.replace('.',
                 .str.reploce(r'["\d.-]", '", regex-True))
   sh4[col] - pd.to_numeric(sh4[col], errors-'coerce')
sh4['All'] = sh4['All'].estype(float)
sh4['Hale'] = sh4['Hale'].astype(float)
sh4["female"] = sh4['female'].estype(float)
sh4['Sex gap'] = sh4['Sex gap'].astype(float)
sh4['2034') = sh4['2024'].astype(float)
sh4['2014:2019'] - sh4['2014:2019'].sstype(float)
sh4["2019"] = sh4["2019"], astype(float)
sh4['2019:2020'] = sh4['2019:2020'].astype(float)
sh4['2020'] - sh4['2020'].astype(float)
sh4['3020:3021'] + sh4['2020:2021'].astype(float)
sh4['2021:2022'] + sh4['2021:2022'].astype(float)
sh4['2022'] = sh4['2022'].astype(float)
sh4["recovery from COVID-19: 3819:3822"] - sh4["recovery from COVID-19: 2819:3833"].astype(float)
```

✓ There was some unknown symbols, we handled them then convert the data types to be appropriate

```
columns_to_fill = [
    "All", "Male", "Female", "Sex gap",
    "2014", "2014:2019", "2019", "2019:2020",
    "2020", "2020:2021", "2021:2022", "2022",
    "recovery from COVID-19: 2019:2022"
]

for col in columns_to_fill:
    sh4[col] = sh4[col].fillna(sh4[col].mean())
```

✓ Handling nulls values.

Sheet5:UN: Estimate of life expectancy for various ages in 2023

This dataset presents the United Nations' 2023 estimates of life expectancy at various ages for different countries, with location links directing to detailed demographic profiles for each country.

Sheet6:OECD (2022)

This dataset provides life expectancy estimates for OECD countries in 2022, with the default sorting based on life expectancy values from 2019.

2. Analysis phase

i. First table (Estimate of life expectancy for various ages in (2023):

```
Countries and territories 2023

Hong Kong, China 85.51

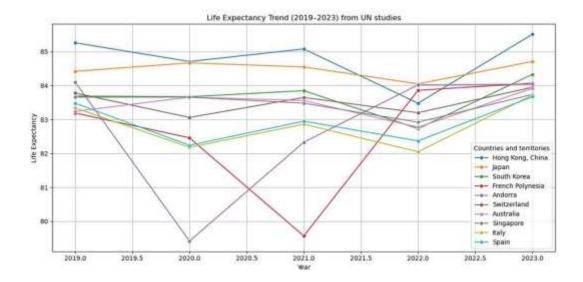
Japan 84.71

South Korea 84.33

French Polynesia 84.07

Andorra 84.04
```

o The highest 5 countries estimated life expectancy in 2023



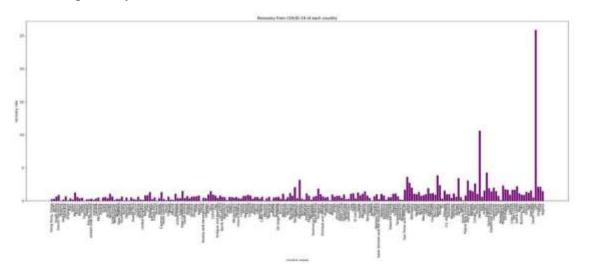
- > From this time line graph we found that:
- o In 2019, Hong Kong, China had the highest life expectancy, and French Polynesia had the least.
- o Andorra got the highest drop in 2020, followed by French Polynesia in 2021.
- o Hong Kong, China held the lead from 2019 to 2023, with the exception of 2022.

| | Countries and territories | Recovery from COVID-19: 2019:2023 |
|-----|---------------------------|-----------------------------------|
| 10 | Réunion | 1.23 |
| 25 | Bermuda | 1.06 |
| 42 | Maldives | 1.33 |
| 47 | Cayman Islands | 1.31 |
| 53 | Panama | 1.08 |
| ••• | *** | *** |
| 205 | Somalia | 1.57 |
| 207 | CAR | 25.88 |
| 208 | Lesotho | 2.13 |
| 209 | Chad | 2.08 |
| 210 | Nigeria | 1.45 |

o Now we can see the countries with high recovery after COVID-19 (Above average of life expectancy)

| | Countries and territories | Recovery from COVID-19: 2019:2023 | |
|-----|---------------------------|-----------------------------------|--|
| 0 | Hong Kong, China | 0.25 | |
| 1 | Japan | 0.29 | |
| 2 | South Korea | 0.64 | |
| 3 | French Polynesia | 0.88 | |
| 4 | Andorra | 0.06 | |
| *** | - | | |
| 192 | Madagascar | 0.13 | |
| 196 | Liberia | 0.93 | |
| 201 | Burkina Faso | 0.91 | |
| 202 | Benin | 0.89 | |
| 206 | South Sudan | 0.51 | |

And these are the countries with low recovery after COVID-19 (less of life expectancy)



 \circ $\;$ This bar plot shows the recovery from COVID-19 of each country.

ii. 2- The second table (World Bank Group (2022)):

| | Countries and territories | 2014:2022 |
|-----|--------------------------------|-----------|
| 139 | St. Vincent and the Grenadines | 5.50 |
| 86 | Lebanon | 4.55 |
| 166 | Yemen | 3.66 |
| 92 | Oman | 3.51 |
| 141 | Guatemalu | 3.29 |
| 143 | Ukraine | 2.60 |
| 131 | Paraguay | 2.41 |
| 95 | Colombia | 2.38 |
| 128 | Jamaica | 2.35 |
| 124 | Belize | 2.35 |

The countries with highest wiggle in life expectancy from 2014 to 2019.

| | Countries and territories | recovery from COVID-19: 2019:2022 | |
|----|---------------------------|-----------------------------------|--|
| 2 | Hong Kong SAR, China | 1.49 | |
| 19 | Iceland | 0.99 | |
| 37 | Greece | 1.00 | |
| 57 | USA | 1.35 | |
| 58 | Costa Rica | 2.11 | |

o Countries with the highest recovery from COVID-19.

| | Countries and territories | recovery from COVID-19: 2019:2022 |
|-----|---------------------------|-----------------------------------|
| 0 | Macao SAR, China | 0.40 |
| 1 | Japan | 0.36 |
| 3 | French Polynesia | 0.55 |
| 4 | Switzerland | 0.45 |
| 5 | Faroe Islands | 0.65 |
| | - | |
| 193 | Cote d'Ivoire | 0.40 |
| 197 | South Sudan | 0.34 |
| 198 | Central African Republic | 0.55 |
| 199 | Nigeria | 0.72 |
| 201 | Chad | 0.26 |

o Countries with the least recovery from COVID-19.

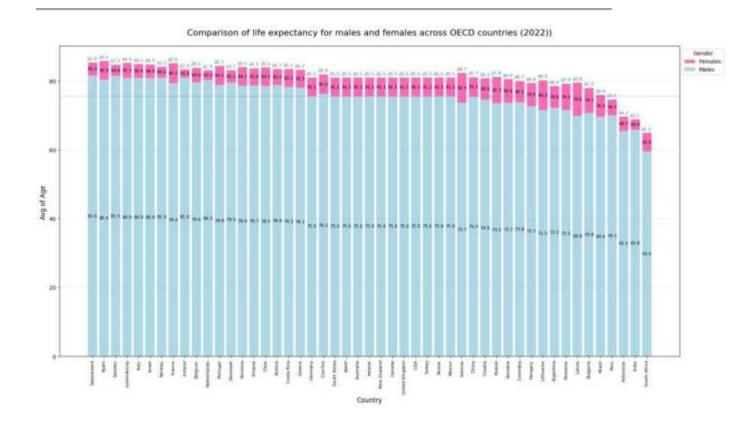
| | Countries and territories | Sex gap |
|-----|---------------------------|---------|
| 102 | Armenia | 10.58 |
| 143 | Ukraine | 10.38 |
| 103 | Belarus | 10.29 |
| 107 | Russia | 10.20 |
| 80 | Latvia | 9.80 |
| 115 | Georgia | 9.70 |
| 81 | Vietnam | 9.36 |
| 105 | Mongolia | 9.30 |
| 142 | Moldova | 9.10 |
| 117 | El Salvador | 8.98 |
| | | |

o The countries with the highest gap between male and female life expectancy.

| Japan - | 94.4 | H.F | 1845 | (94) | |
|--|---------|-------|-------|---------|-------------|
| Switzerland - | 84.0 | 60.3 | 83.9 | (40.5) | |
| Spain - | 1944 | 82.4 | 103 | (1002) | |
| Rely - | 100 | 62.3 | 197 | (200) | |
| South Korea - | 1949 | 83.5 | 20.0 | 79.8 | 12.5 |
| iceland - | 1882 | 811 | 83.3 | 1981 | |
| Sweden - | 612 | 10.0 | 483 | 011 | |
| Norway - | 10.0 | 1111 | 23.2 | (12.6) | |
| France - | 100 | 117 | 44 | (49) | |
| Australia - | (42.9) | 132 | 10.0 | 76.5 | |
| Strant - | 80.9 | 937 | 82.6 | 12.0 | 90.0 |
| Instand - | 62.6 | 82.6 | 82.4 | 984 | |
| Luxembourg - | 107 | 62.2 | 12.1 | #2.0 () | |
| Canada - | THIS A. | 47 | 96.6 | ONA! | |
| Netherlands - | 802 | 024 | 1664 | 149 | |
| New Zesland | 18000 | 63 | 10.3 | 784 | -77.5 |
| Fire and | 1(04) | H2.0 | 161.0 | (9.8) | 100 |
| Belgium - | (803) | 00.0 | FE.8 | (410) | |
| Autos | 19.8 | #12. | 383 | 2.55.) | |
| Portugal | MA. | 41.1 | 83.3 | 1.337 | |
| Greece - | iii) | 955 | 10,2 | 40.7 | |
| Stoverna - | 18.5 | 90.8 | 80,7 | 11.3 | 73.0 |
| Demnark - | 81.5 | #1.6 | *13 | 14.3 | |
| United Kingdom - | 810 | 10.4 | 76.7 | 78.4 | 3 |
| Attended to the control of the contr | *11.7 | 01.1 | 80.0 | 10.7 | Age Average |
| S this | 60.6 | 60.8 | 81.6 | 112 | Age. |
| Costa Rica | 803 | 80.6 | 10.1 | mito | -72.5 |
| Czechia | 1881 | Thi . | | 7361 | |
| Intonia - | 19.0 | ELE: | 772 | 362 | |
| USA - | (164) | 77.0 | 260 | 580 | |
| Croatia - | 900 | 11.41 | :007 | - 700 X | |
| Yarbey - | (10.0) | 76.2 | 187 | 78.4 | |
| Poland - | (160) | | | (204.) | 70,0 |
| Sitteraktie - | | 17.0 | | 37.9 | |
| China - | (3),7/2 | 78.9 | 18.3 | 79.3 | |
| Argentina - | 77,0 | 77.3 | | 75.4 | |
| Colombia - | 700 | 100 | 188 | 7887 | |
| Hungary - | 18.5 | 87 | | 763 | -67.5 |
| Lithuaria | (45) | 702 | 202 | (86) | |
| Peru - | 160 | 764 | | 104 | |
| Letvie - | 7807 | 16.5 | | 300 | |
| Renana - | 28 K | | | 76.5 | |
| Brazil - | | 75.8 | 740 | (20) | 10000 |
| Mexico - | This | | 180 | 368.) | -65.0 |
| Bulgaria - | | | | (201) | |
| Russia - | (84) | 79.2 | 36.7 | 776.6 | |
| tide - | (98) | 245 | 30.1 | 612 | |
| Indonesia - | | | | 07.0 | |
| South Africa - | 85.7 | 98.2 | 65.3 | 62.3 | -17.5 |
| | 2019 | 2020 | 2021 | | |

➤ In this heat map we compare the Average of ages in Countries based on OECD report and we can know that :

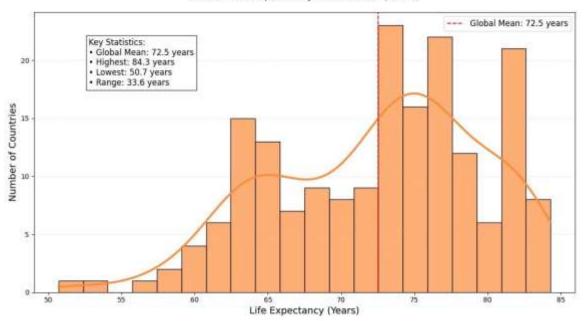
- 1) Japan and Switzerland has the highest value.
- 2) South Africa and Indonesia has lowest value
- 3) Some countries have firmness in values over years
- 4) Color changes mean that there is increase or decrease in values



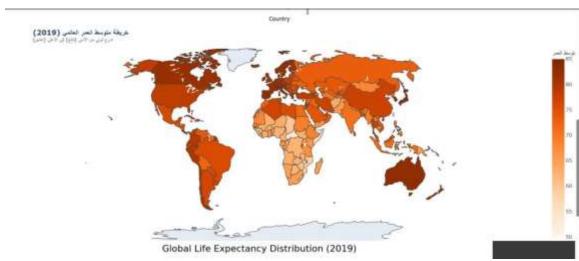
- ➤ In this plot we compare the ages between males and females in each country and we concluded that :
- 1) In Spain females has the highest value and South Africa has the lowest value
- 2) In Switzerland men has the highest value South Africa has the lowest value
- 3) Mostly the difference between genders isn't big and we want to say that this difference may occur because of different levels of health care between them or different cultures between countries
- 4) The blue line show the averages and white lines show where each country is belong above or below the average line.

iii. Third table (World Health Organization):

Global Life Expectancy Distribution (2019)



- In this plot we understand how life expectancy is distributed among countries and support analysis with statistical data that facilitate the results and we conclude that:
- 1) Most countries are located in range (65:80).
- 2) The distribution of ages is asymmetrical.
- 3) Some countries have low ages (50) that refers to problems of health.
- 4) Most of countries are close to the average of ages (72.5)
- 5) There is a clear disparities between countries because the range is a little big.

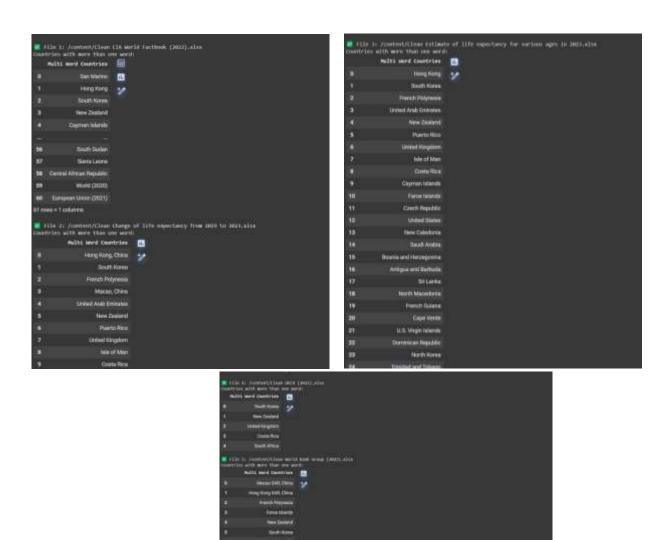


➤ We Creates an interactive world map (Choropleth Map) showing Life Expectancy per country in 2019 and we concluded that:

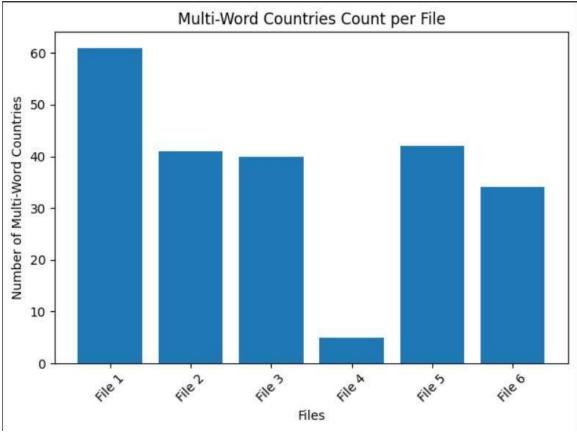
- 1) Western Europe, Australia show darker colors (higher life expectancy).
- 2) Some African and Asian countries show lighter colors (lower life expectancy).
- 3) The map reflects health care globally.
- 4) Some countries have shown significant improvements compared to 2000, indicating progress.

❖ The Regular expression

> These tables have the countries that have more than one word in each file after implement regular expression on columns of countries names in these tables.

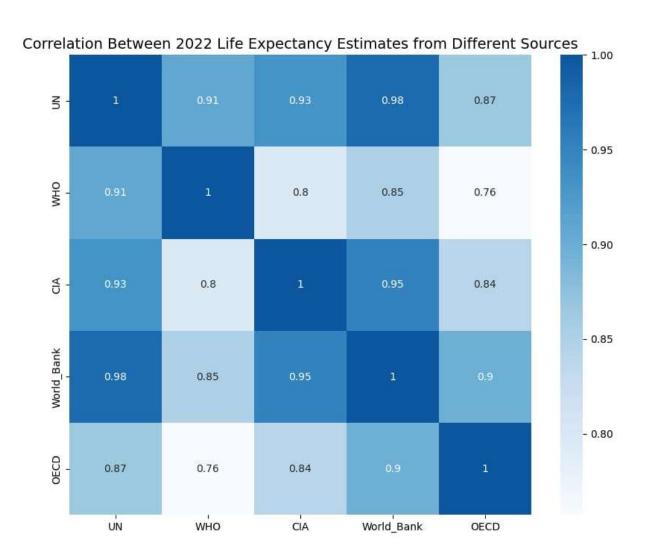






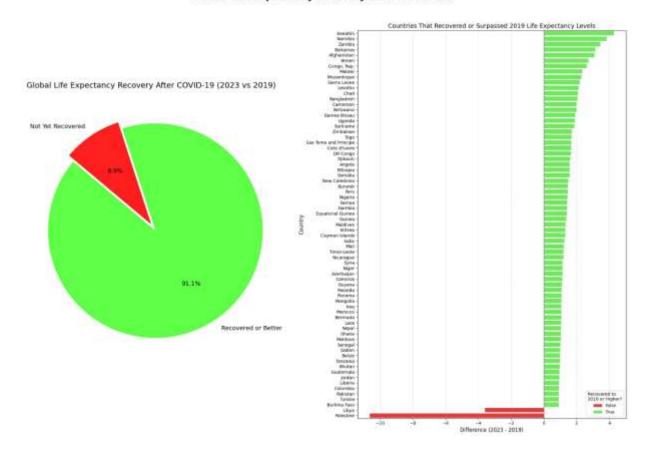
➤ In this plot we show number of countries that has more than one word in each file after implement Regular Expression on columns that have names of countries.

3. The comparisons and relationships between sources

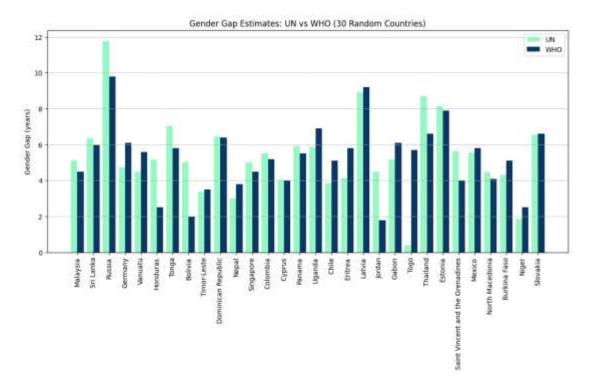


➤ The heatmap visually represents the degree of consistency among sources of life expectancy estimates for 2022 by examining their correlation coefficients.

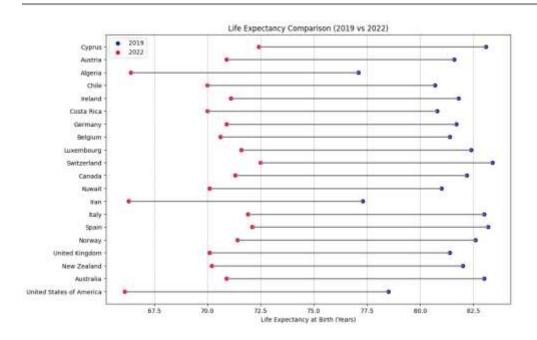
Global Life Expectancy Recovery After COVID-19



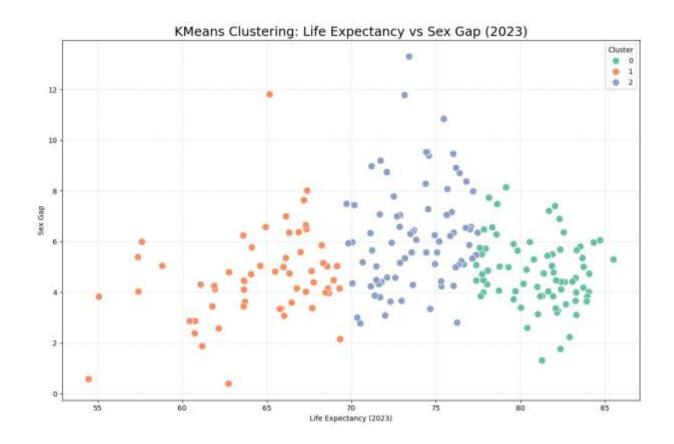
- ➤ The pie chart illustrates the percentage of countries that have recovered or surpassed their pre-pandemic life expectancy compared to those that have not yet recovered.
- ➤ The bar chart highlights the change in life expectancy for each country between 2019 and 2023, distinguishing between countries that have recovered and those that have not.



➤ The bar chart compares gender gap estimates in life expectancy between United Nations and World Health Organization data for 30 randomly selected countries.



➤ This chart displays the top 20 countries most affected by the COVID-19 pandemic in terms of decline in life expectancy at birth, based on a comparison between 2019 and 2022 data



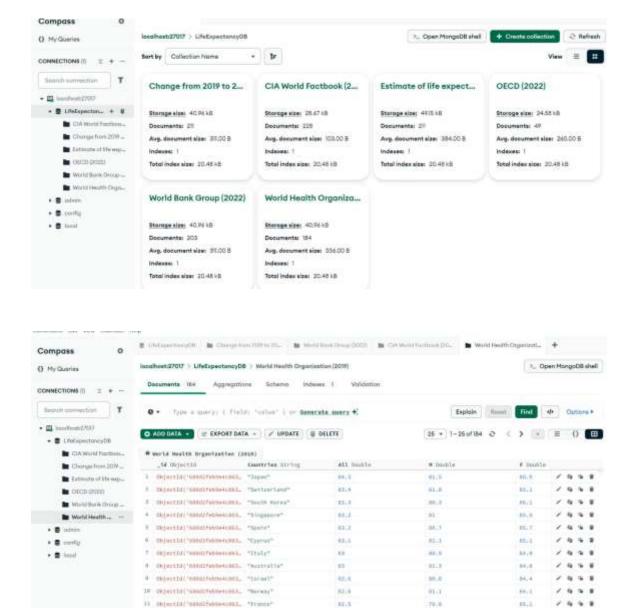
| | Cluster 0 | Cluster 1 | Cluster 2 |
|-------|------------------|-----------|------------------|
| Anti | gua and Barbuda | Libya | Hong Kong, China |
| | Sri Lanka | Yemen | Japan |
| | Argentina | Botswana | South Korea |
| | Ecuador | Laos | French Polynesia |
| | Guam | Senegal | Andorra |
| | | | |
| | Guyana | NaN | NaN |
| | Turkmenistan | NaN | NaN |
| | Greenland | NaN | NaN |
| | Philippines | NaN | NaN |
| ао То | ome and Principe | NaN | NaN |

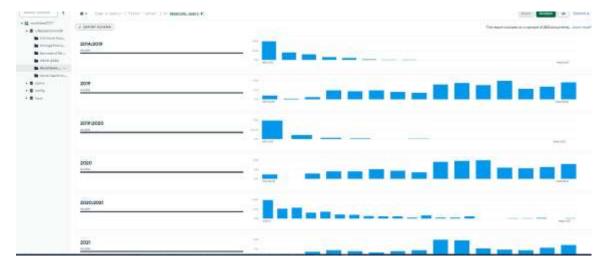
80 rows × 3 columns

➤ The chart is used to classify countries into three groups based on their 2023 life expectancy, recovery from COVID-19, and gender gap, understanding the relationship between life expectancy, the gender gap, and COVID-19 recovery levels.

4. Data Storage (MONGODB)

➤ We stored data from various sources about life expectancy information in a MongoDB database for easier access and analysis.

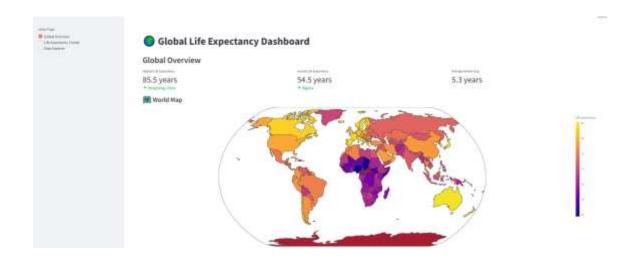


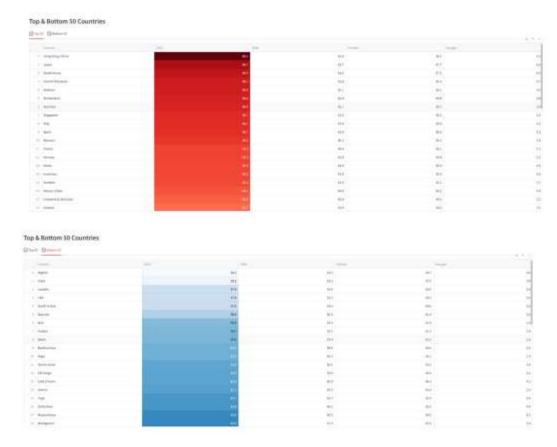


➤ This schema provides an overview of how life expectancy data is distributed for each year or period. It reveals any irregularities or missing values across the dataset.

5. The Streamlit Phase

- ➤ The dashboard contains three interfaces: Global Overview, Gender Analysis, and Data Explorer. Each interface offers a unique perspective on life expectancy.
- 1. Global Overview It displays:
- a) Summary Statistics about life expectancy and gender gap.
- b) Interactive world map showing life expectancy by country.





c) The tables present the top 50 and bottom 50 countries ranked by life expectancy in 2023, along with the gender gap in these countries

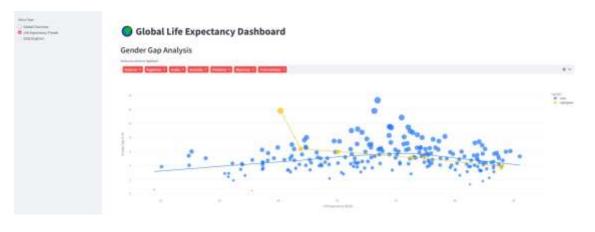




d) This part applies KMeans clustering to group countries based on life expectancy in 2023 and the gender gap and display it in split tables.

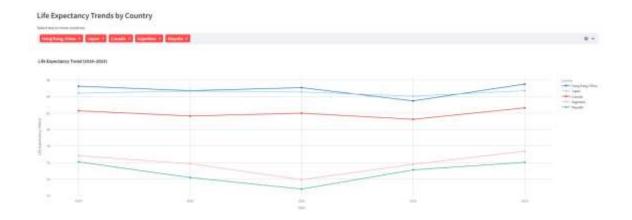
2- Life Expectancy Trends It displays:

a) Scatter plot: displays the relationship between life expectancy in 2023 and the gender gap, highlighting specific countries and a trend line to show the general pattern of the relationship between life expectancy and the gender gap.





b) The tables display the top 30 countries with the largest gender gaps in life expectancy, and the 30 countries with the smallest gender gaps.



c) This plot displays the trend of life expectancy for selected countries over the years 2019 to 2023.



d) This chart compares the changes in life expectancy at birth between 2019 and 2022 for different countries, showing how COVID-19 may have impacted life expectancy.

- 3- Data Explorer: it explore raw datasets from various sources. It desplay:
- a) A selection menu to choose a dataset.
- b) Basic information like number of rows & columns, column names and some statistics.
- c) Finnaly, represent full data with interactive way



Global Life Expectancy Dashboard -

