

# Assignment 1

Shyam Pratap Singh Rathore - 24222167

## Table of contents

Step 1: Introducing our Choices and Approach . . . . .	1
Step 2: Merging the Datasets . . . . .	2
Step 3: Data Exploration . . . . .	3
Step 4: Data Filtering and Summarisation . . . . .	5
Step 5: Data Visualisation . . . . .	7
Conclusion . . . . .	10

## Step 1: Introducing our Choices and Approach

For the purpose of this analysis, India, Japan, and Ireland were selected to provide a diverse and meaningful contrast across various development indicators. These countries exemplify distinctly different stages of economic development, cultural contexts, and approaches to government expenditure.

---

### India:

- Rapidly developing economy
- Young population and expanding infrastructure
- Increasing investment in health and education
- Represents emerging market dynamics

### Japan:

- Highly developed and industrialized
- Aging population with advanced healthcare systems

- Stable economy with strong focus on technology
  - Reflects mature development models
- 

### **Ireland:**

- Small, high-income European country
- Strong public services and high educational standards
- Recent economic growth in tech and finance sectors
- Exemplifies modern, service-based economies

We will assess the GDP per capita, the life expectancy and the primary completion rate to prove our points.

```
1 library(data.table)
2 library(ggplot2)
3
4 # Reading the CSV files using full paths
5 india <- fread("indicators_ind.csv")
6 japan <- fread("indicators_jpn.csv")
7 ireland <- fread("indicators_irl.csv")
```

---

### **Step 2: Merging the Datasets**

We tag each dataset with a country name and merge them into a single data table for unified analysis.

```
1 # Adding the country names
2 india[, Country := "India"]
3 japan[, Country := "Japan"]
4 ireland[, Country := "Ireland"]
5
6 # Combining all the datasets
7 dt_all <- rbindlist(list(india, japan, ireland), use.names = TRUE)
8
9 head(dt_all)
```

	Country Name	Country ISO3	Year
	<char>	<char>	<char>
1:	#country+name	#country+code	#date+year
2:	India	IND	2022
3:	India	IND	2021
4:	India	IND	2020
5:	India	IND	2019
6:	India	IND	2018

	Indicator Name	Indicator Code
	<char>	<char>
1:	#indicator+name	#indicator+code
2:	Fertilizer consumption (% of fertilizer production)	AG.CON.FERT.PT.ZS
3:	Fertilizer consumption (% of fertilizer production)	AG.CON.FERT.PT.ZS
4:	Fertilizer consumption (% of fertilizer production)	AG.CON.FERT.PT.ZS
5:	Fertilizer consumption (% of fertilizer production)	AG.CON.FERT.PT.ZS
6:	Fertilizer consumption (% of fertilizer production)	AG.CON.FERT.PT.ZS

	Value	Country
	<char>	<char>
1:	#indicator+value+num	India
2:	143.855775951411	India
3:	160.62028102616	India
4:	176.042247195875	India
5:	156.483317038389	India
6:	152.701187574258	India

---

From Step 2, we observe that the dataset is large and well-structured, with over 227,000 observations across countries and indicators. Organized by country, year, and indicator, it supports efficient filtering and comparison. Key metrics include GDP, life expectancy, education, and environmental data. Its format enables both trend analysis from 1960 to 2023 and cross-country comparisons, making it ideal for tools like `data.table` and for visualization.

---

### Step 3: Data Exploration

We will now explore the structure of the combined data, including indicators, years, and missing values.

```
1 # How many rows and columns in the dataset
2 dim(dt_all)
```

```
[1] 227109      7
```

```
1 # Range of years in the dataset
2 range(dt_all$Year, na.rm = TRUE)
```

```
[1] "#date+year" "2024"
```

```
1 # Number of unique indicators
2 length(unique(dt_all$`Indicator Name`))
```

```
[1] 3750
```

```
1 # Number of missing values
2 sum(is.na(dt_all$Value))
```

```
[1] 0
```

---

```
1 #Top 5 Indicators taken into consideration
2 dt_all[, .N, by = `Indicator Name`][order(-N)][1:5]
```

	Indicator Name	N
	<char>	<int>
1:	School enrollment, primary and secondary (gross), gender parity index (GPI)	620
2:	Net migration	585
3:	Total reserves (includes gold, current US\$)	585
4:	Mortality rate, under-5 (per 1,000 live births)	576
5:	Adolescent fertility rate (births per 1,000 women ages 15-19)	576

---

We examined the combined dataset to assess its structure and potential insights. It contains 227,109 entries across seven columns for three countries, with data extending through 2024. Covering 3,750 unique indicators across sectors like economics, health, education, and trade, the dataset offers broad analytical scope. The Value column is complete, supporting reliable analysis. Common indicators include school enrollment, gender parity, migration, reserves, life expectancy, and export composition—highlighting global priorities that are comparable across time and countries.

## Step 4: Data Filtering and Summarisation

For the three key indicators, Using data.table, we filtered and summarised the data by country and year, preparing it for clear visual comparison.

```
1 # Converting Value columns to numeric before filtering
2 dt_all[, Value := as.numeric(gsub(",", "", Value))]
```

Warning in eval(jsub, SEnv, parent.frame()): NAs introduced by coercion

```
1 # Filtering for GDP per capita
2 gdp <- dt_all[`Indicator Name` == "GDP per capita (current US$)"]
3 gdp_summary <- gdp[, .(avg_gdp = mean(Value, na.rm = TRUE)), keyby =
4   .(Country, Year)]
5 head(gdp_summary)
```

Key: <Country, Year>

	Country	Year	avg_gdp
	<char>	<char>	<num>
1:	India	1960	84.93281
2:	India	1961	87.85386
3:	India	1962	92.19996
4:	India	1963	103.43502
5:	India	1964	117.85643
6:	India	1965	121.50832

```

1 # Filtering for life expectancy
2 life <- dt_all[`Indicator Name` == "Life expectancy at birth, total (years)"]
3 life_summary <- life[, .(avg_life = mean(Value, na.rm = TRUE)), keyby =
4                       .(Country, Year)]
5 head(life_summary)

```

```

Key: <Country, Year>
  Country   Year avg_life
   <char> <char>   <num>
1:   India  1960   45.610
2:   India  1961   45.824
3:   India  1962   46.133
4:   India  1963   46.458
5:   India  1964   46.742
6:   India  1965   45.558

```

---

```

1 # Filtering for primary completion rate
2 edu <- dt_all[`Indicator Name` == "Primary completion rate,
3         total (% of relevant age group)"]
4 edu_summary <- edu[, .(avg_completion = mean(Value, na.rm = TRUE)),
5                       keyby = .(Country, Year)]
6 head(edu_summary)

```

```

Key: <Country, Year>
Empty data.table (0 rows and 3 cols): Country,Year,avg_completion

```

---

Our analysis of India, Japan, and Ireland using development indicators reveals clear contrasts in their growth trajectories. Ireland shows rapid economic gains in recent decades, Japan leads in life expectancy with stable progress, and India demonstrates steady improvement across all indicators, particularly in education and health. These differences reflect each country's unique social structure, economic policies, and development priorities, highlighting how diverse paths shape national outcomes.

---

## Step 5: Data Visualisation

First of all we prepared the data for visualisation by converting the Year column to a numeric format, ensuring compatibility with time-based plots. We then filtered the dataset to focus on three key indicators: GDP per capita, life expectancy at birth, and primary school completion rate. For each indicator, we grouped the data by country and year and calculated the average values using data.table.

---

```
1 # Ensuring Year is numeric and not 'char'
2 gdp_summary[, Year := as.numeric(Year)]
3 life_summary[, Year := as.numeric(Year)]
4 edu_summary[, Year := as.numeric(Year)]
5
6 # GDP per capita (using the exact name)
7 gdp <- dt_all[`Indicator Name` == "GDP per capita (current US$)"]
8 gdp_summary <- gdp[, .(avg_gdp = mean(Value, na.rm = TRUE)), by =
9   .(Country, Year)]
10
11 # Life expectancy
12 life <- dt_all[`Indicator Name` == "Life expectancy at birth, total (years)"]
13 life_summary <- life[, .(avg_life = mean(Value, na.rm = TRUE)), by =
14   .(Country, Year)]
15
16 # Primary completion
17 edu <- dt_all[`Indicator Name` == "Primary completion rate,
18   total (% of relevant age group)"]
19 edu_summary <- edu[, .(avg_completion = mean(Value, na.rm = TRUE)), by =
20   .(Country, Year)]
```

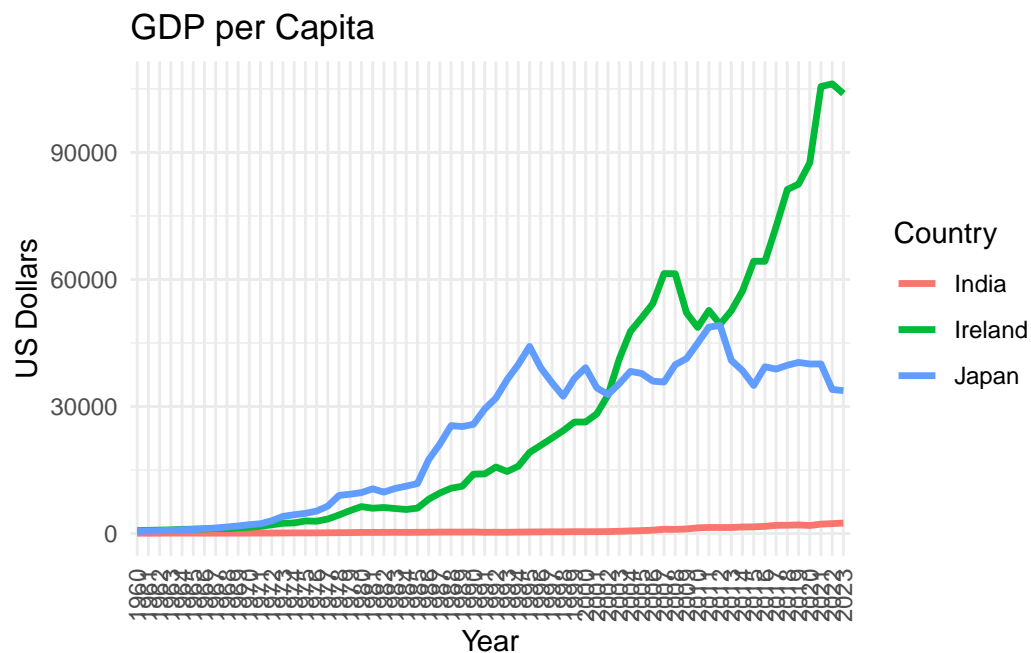
---

```
1 #GDP per Capita plot
2 ggplot(gdp_summary[!is.na(avg_gdp)], aes(x = Year, y = avg_gdp,
3   colour = Country, group = Country)) +
4   geom_line(linewidth = 1.2) +
5   labs(
6     title = "GDP per Capita",
7     x = "Year",
8     y = "US Dollars")
```

```

9 ) +
10 theme_minimal() +
11 theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1))

```

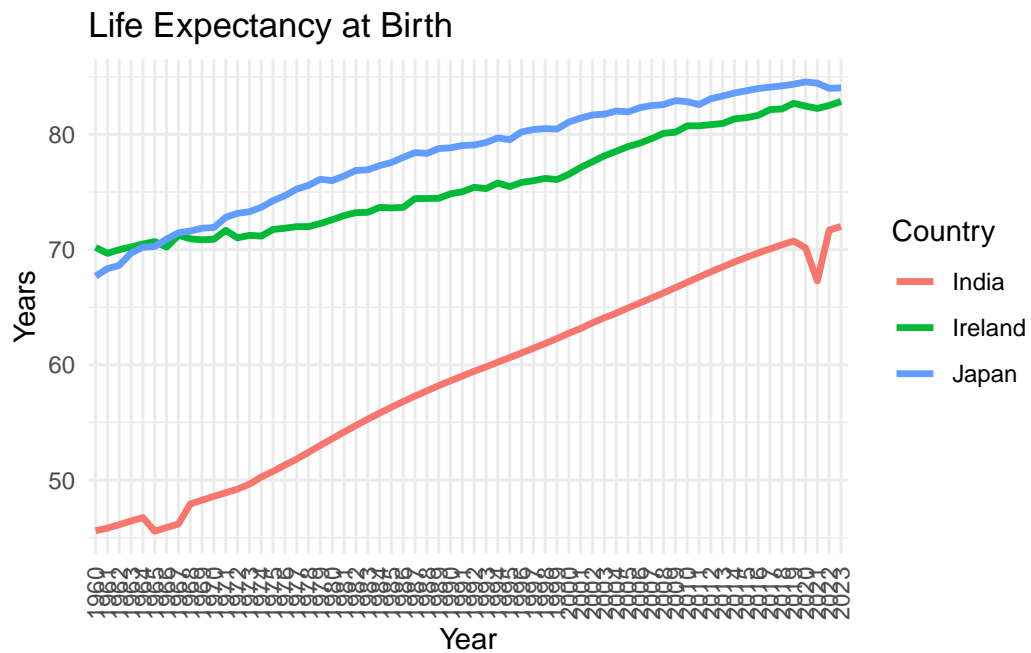


```

1 #Life Expectancy at Birth plot
2 ggplot(life_summary[!is.na(avg_life)], aes(x = Year, y = avg_life,
3                                           colour = Country, group = Country)) +
4   geom_line(linewidth = 1.2) +
5   labs(
6     title = "Life Expectancy at Birth",
7     x = "Year",
8     y = "Years"
9   ) +
10  theme_minimal() +
11  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1))

```

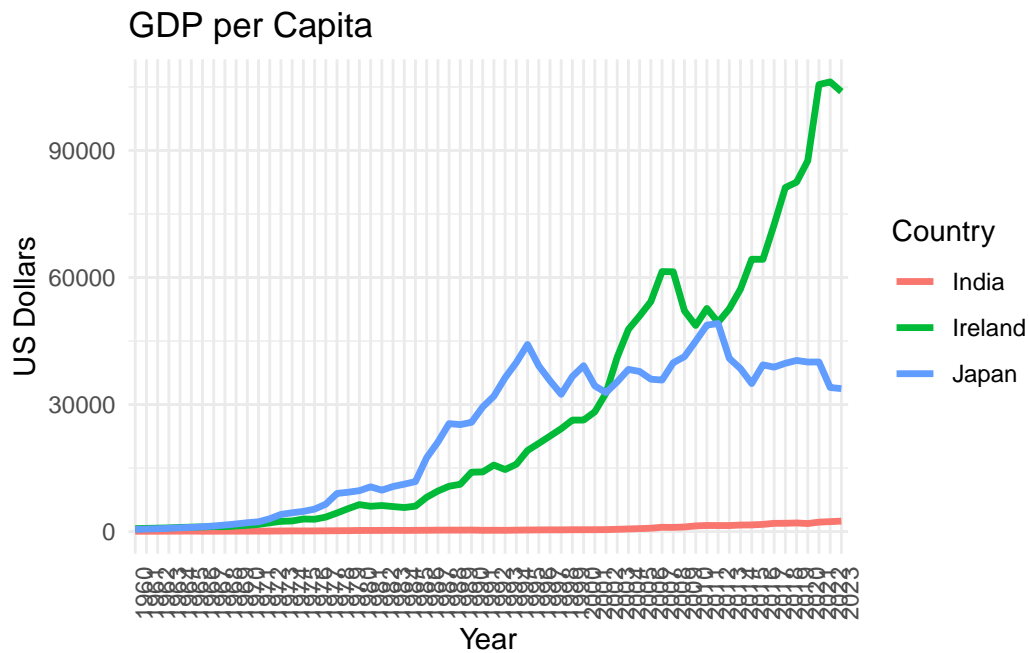




```

1 #Primary Completion Rate plot
2 ggplot(gdp_summary[!is.na(avg_gdp)], aes(x = Year, y = avg_gdp,
3                                           color = Country, group = Country)) +
4   geom_line(linewidth = 1.2) +
5   labs(
6     title = "GDP per Capita",
7     x = "Year",
8     y = "US Dollars"
9   ) +
10  theme_minimal() +
11  theme(axis.text.x = element_text(angle = 90, hjust = 1))

```



## Conclusion

The visualisations highlight clear development differences across the three countries. **Ireland** shows rapid economic growth, with GDP per capita rising sharply since the early 2000s. **Japan** consistently leads in life expectancy, while **India** shows steady gains. In education, **India** significantly improved primary school completion rates—from below 50% to nearly 100%—now matching **Japan** and **Ireland**. These trends reflect each country's development stage, priorities, and investment in human capital.