# Exercises: Relational data for country-level statistics
## YSC2210 - DAVis with R

Michael T. Gastner

## Introduction

In an earlier exercise, we created a plot similar to figure 1, which was made by the Gapminder foundation (Gapminder, 2016). The plot shows GDP per capita (x-axis), life expectancy (y-axis) and population (size) by country. In this exercise, we take a closer look at publicly available data for these variables.
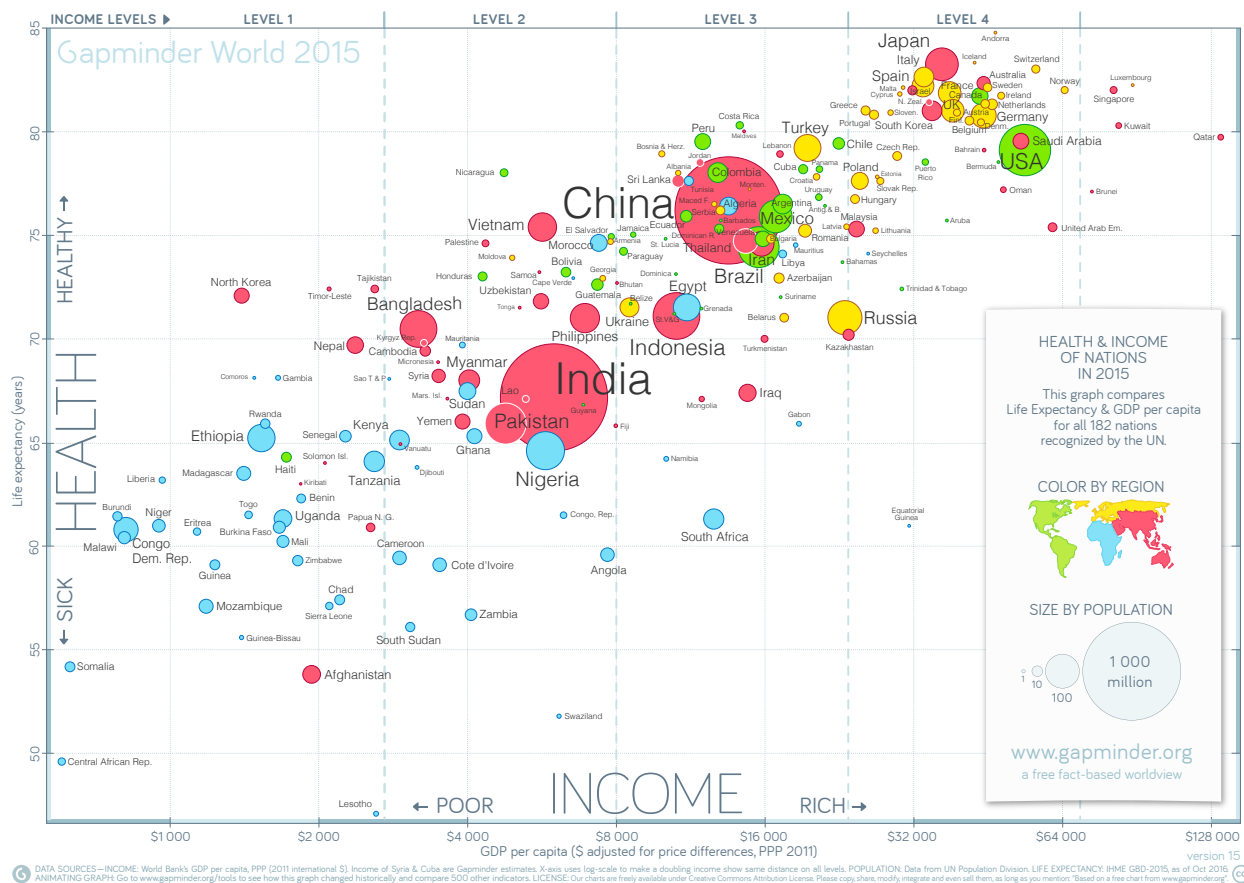


Figure 1: Image from Gapminder (2016)

## Learning objectives

We will practice our data wrangling skills. We will work with different kinds of joins and other functions from the **dplyr** package.

1

# Data

We need three data sets from the World Bank:

- GDP per capita (current US$): https://data.worldbank.org/indicator/NY.GDP.PCAP.KD
- Life expectancy at birth, total (years): https://data.worldbank.org/indicator/SP.DYN.LE00.IN
- Population: https://data.worldbank.org/indicator/SP.POP.TOTL

Please click on the EXCEL download button and save the downloaded XLS files in your project directory.[1]

At the end of this exercise, we compare these data with data from R's **gapminder** package.

# Tasks

(1) Import the World Bank data for GDP per capita, life expectancy and population.

(2) Is the column with three-letter country codes (second column from the left) the same in all three spreadsheets?

(3) Merge the three spreadsheets into a single tibble `wb` (for *W*orld *B*ank) with columns for:

- country name.
- country code.
- year.
- GDP per capita.
- life expectancy.
- population.

(4) Some rows in the World Bank spreadsheets do not represent a (single) country, for example 'East Asia & Pacific (excluding high income)'. We want to remove the corresponding rows from `wb`. We are going to automate this task by using the tibble `codelist` in the **countrycode** package.

The purpose of the **countrycode** package is to simplify the task of merging country-level data in different data bases. The same country often appears under a variety of names in official documents. For example, 'United States of America', 'U.S.A.' and 'US' all refer to the same country. The recommended practice when joining country-level data in different data bases is to use a standardised set of codes that uniquely identify each country. One option is to work with ISO 3166-1 alpha-3 codes.[2] These codes are in the column `iso3c` of `codelist`.

Perform an anti-join to find out which three-letter country codes in the World Bank spreadsheets do not have a matching code in `codelist`. What are the corresponding 'country names'? Do the results make sense?

(5) Use a **dplyr** 'join' function to remove all rows from `wb` that do not match any country code in `codelist`.

(6) A country can only be added to the scatter plot shown in figure 1 if all of the following three pieces of information about the country are known:

- GDP per capita.
- life expectancy.
- population.

Summarise the number of countries per year that cannot be plotted on the basis of the World Bank data. Here are the first few rows of a tibble that shows the number of missing countries for each year.

```
head(missing_values)
```

---

[1]We work with XLS files because there is a minor formatting issue with the CSV files provided by the World Bank.

[2]For background information, see https://www.iso.org/iso-3166-country-codes.html.

```
## # A tibble: 6 x 2
##    year  na_countries
##    <chr>        <int>
## 1 1960           130
## 2 1961           126
## 3 1962           126
## 4 1963           126
## 5 1964           126
## 6 1965           121
```

(7) Plot the number of missing countries per year. Comment on the result.

(8) R's **gapminder** package contains a tibble `gapminder_unfiltered` that is an alternative source of information about GDP, life expectancy and population. Create a tibble named `gap` that appends a column with three-letter country codes to `gapminder_unfiltered`. Let us remove the Netherlands Antilles, which do not have an officially assigned ISO 3166-1 alpha-3 code and do not appear in the World Bank data, presumably because these data are aggregated with those of the Netherlands. Here is the code snippet for this task:

```
gap <-
  gapminder_unfiltered |>
  filter(country != "Netherlands Antilles") |>
  mutate(country_code = countrycode(country, "country.name", "iso3c"))
```

After running this command, are there countries in `gap` without a country code? Are there countries that share the same country code?

(9) Which countries are in `gap`, but do not appear in `wb`? Which countries are in `wb`, but do not appear in `gap`?

(10) Let us compare the GDP data in `wb` and `gap` for the year 2007. Remove all unrelated rows and columns. Merge the information from `wb` and `gap` into a tibble `wb_gap` such that only those countries are included that appear in both tibbles.

(11) Append a column to `wb_gap` that shows the percentage difference of Gapminder's GDP estimate compared to the World Bank estimate. For example, if the World Bank's estimate is $5000 and Gapminder's estimate is $2500, the percentage difference is -50%.

(12) For which five countries is the percentage difference largest? For which five countries is it smallest (i.e. most strongly negative).

# References

Gapminder (2016). Updated Gapminder World Poster 2015! URL: https://www.gapminder.org/downloads/updated-gapminder-world-poster-2015/. Accessed on 2020-11-26.