

Module 2 Exercise

February 5, 2018

Instructions

For the remaining modules, you will be working directly in RStudio.

1 Create an R-Script File

- In RStudio, create a new R Script file. Save it as “Exercise1” in a new project folder.
- Install the packages you will be using for analysis (if not already installed):
 - tidyverse
 - knitr
- Set the working directory to your project folder.

2 Import the data

For this exercise, you will be working with data from the “World Development Indicators” data set produced by the World Bank.

- Download the dataset from <https://data.worldbank.org/data-catalog/world-development-indicators>.
 - You will need to unzip the download and move the “WDIData.csv” file into a suitable project folder. This will be your data set.
- Import the “WDIData.csv” data set (with a suitable name for the data)
- Using pipes, convert the data frame into a tibble.

3 Filter the data set for the desired indicators

“WDIData.csv” is a pretty large data set. For the exercise, you will be working with indicators related to poverty and inequality. Reduce the size of the data set in memory by filtering for only the indicator codes listed below:

- First create a vector containing the indicator codes you want to keep.
- Rename the indicator code column to just “indicator” to make it easier to work with.
- Use `filter()` with the `%in%` operator to selected the desired observations.

Indicator	Short Description
NY.GDP.PCAP.KD	GDP per capita (constant 2010 US\$)
SP.POP.TOTL	Population, total
SL.TLF.CACT.FM.ZS	Ratio of female to male labor force participation rate (%) (modeled ILO estimate)
SE.SEC.CUAT.UP.ZS	Educational attainment, at least completed upper secondary, population 25+, total (cum. %)
SL.UEM.NEET.ZS	Share of youth not in education, employment or training, total (% of youth population)
SL.UEM.BASC.ZS	Unemployment with basic education (% of total labor force with basic education)
SI.POV.UMIC	Poverty headcount ratio at \$5.50 a day (2011 PPP) (% of population)
SI.DST.FRST.20	Income share held by lowest 20%
SI.DST.02ND.20	Income share held by second 20%
SI.DST.03RD.20	Income share held by third 20%
SI.DST.10TH.10	Income share held by highest 10%

4 Tidy the data set

- a. First, get a sense of the data structure.
- b. Drop the following columns: “v63”, “Indicator Name”
- c. In order to easily work with the data, you will first need to tidy it.
 - Use the **gather** and **spread** functions so that the data adheres to tidy data principles.
 - You may need to think about how to specify the columns to be gathered since they appear as numbers.

5 Rename and change class of columns Columns

- To make the data truly presentable, rename the indicator columns to a short descriptive variable names (be sure not to include any spaces).
- Change the class of indicators to numeric.
- Rename all other variables to one-word names as well.

6 Further restrict the sample

We are specifically interested in inequality and poverty measures for “high-income” countries, with at least 4 million people, over the last ten years. If you look carefully at the dataset, you will notice that there are multiple issues that need to be addressed:

- a. There are a number of non-country regions that are present which need to be removed.
 - The list of county codes that need to be filtered out are supplied in `regions.csv`.
 - Import this file and filter for country codes not in `regions`.
 - **Tip:** Use the following code to get `regions` into a format where you can use it with an `%in%` filter:

```
regions <- import("regions.csv")
regions <- as.character(regions$CountryCode)
```

- Drop the `countrycode` column once you have removed all regional observations.
- b. The sample should be restricted to countries with a population greater than 4 million, GDP per capita greater than \$12,000, and year greater than or equal to

2007.

7 Create new variables

Now, create some new variables that we might be interested in: - log of GDP per capita - The absolute number of people in poverty (at \$5.50 a day) - Income share held by the bottom 60% of the population - Ratio of income held by top 10% to bottom 60%

8 Create and view summary statistics

- a. By country, create average values (with option `na.rm=TRUE`) for:
 - Ratio of income held by top 10% to bottom 60%
 - GDP per capita
 - Ratio of female to male labor force participation

- Share of population having completely at least completed upper secondary schooling
 - Share of youth not in education, employment or training, total (% of youth population)
 - Unemployment with basic education (% of total labor force with basic education)
- b. Then use `percent_rank` to create a table of percentile ranks for each of the indicators in (a) by country.
- c. Keep only the percentile rank columns and sort by per capita GDP percentile rank.

9 Display basic descriptive statistics for all variables in the general data set.