

# Warming Up to R

---

## Ilana Berlin

### Problem One

Re-open the scooby data set in R and glance through the rows and columns. Imagine you've been asked to use this set to answer the question, "do women who capture the villain get to unmask them as often as men who capture the villain?" What might you want to do in order to answer this? What data challenges could you foresee? Write a short paragraph.

The first thing that comes to mind when asked this question is to compare the sum of villains captured and unmasked by Daphne and Velma to the sum of the villains captured and unmasked by Fred, Shaggy, and Scooby. There may be null values in episodes where villains were not caught by the main members of Mystery Inc. or were not captured at all. There may also be nulls if a villain was not unmasked.

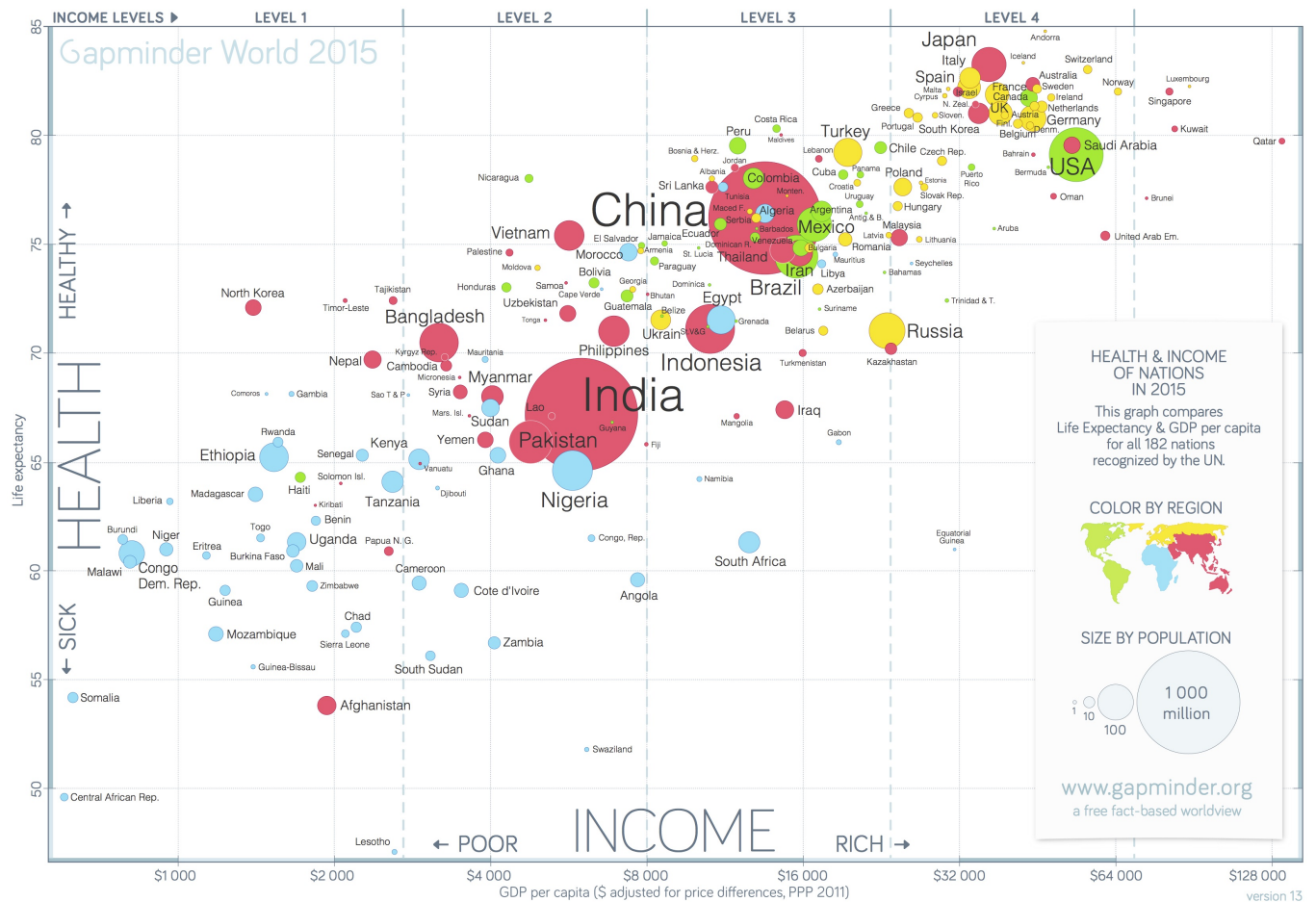
### Problem 2

Find and download the ggplot2 cheat sheet from the Posit (the company that makes RStudio) website. According to this, what are four possible (and equivalent) ways of plotting one categorical and one quantitative variable?

A bar chart (`geom_col()`) a boxplot (`geom_boxplot()`) a swarm plot (`geom_dotplot`) and a violin plot (`geom_violin`).

### Problem 3

The following plot represents data from the famous gapminder data set.



DATA SOURCES—INCOME: World Bank's GDP per capita, PPP (2011 international \$), with a few additions by Gapminder. X-axis uses log-scale to make a doubling incomes show same distance on all levels. POPULATION: Numbers from UN Population Division. LIFE EXPECTANCY: IHME GBD-2015, as of Oct 2016. INTERACTIVE GRAPH: [www.gapminder.org/tools](http://www.gapminder.org/tools) which lets you animate historic data for hundreds of indicators. LICENSE: Our charts are freely available under Creative Commons Attribution License. Please copy, share, modify, integrate and even sell them, as long as you mention: "Based on a free chart from [www.gapminder.org](http://www.gapminder.org)." 

There are four primary aesthetics here. Identify them by specifying both the variable and the way it's being represented. Also identify at least 3 non-data aspects of this plot.

The primary aesthetics used are  $x$  = GDP per capita,  $y$  = life expectancy, color = continent, and size = population. Non-data aspects include the arrows indicating rich and poor, the arrows indicating sick and healthy, the income levels, and the names of countries.