

Visualisation with ggplot2

2031 - Statistical Computing

Sam Clifford

2019-11-15

Introduction

About this practical session

In the lecture session we introduced visualisation with the histogram, x - y plots and other scatter plot techniques, and touched on Tufte's principles of graphical excellence.

This prac will investigate the visual display of data and what makes a good and a bad graph.

- Assumed skills
 - Writing R code into a script file
 - Identifying things that are visually pleasing
- Learning objectives
 - Identifying things that are informative
 - Being able to critique a graph
 - Understanding why and how data is encoded and decoded visually
 - Understanding the subjectivity of what is aesthetically pleasing
- Professional skills
 - Creating high quality graphics

Group formation

Organise yourselves into groups of 2-3 students to collaboratively solve the following exercises.

A reminder of expectations in the prac:

- Keep a record of the work being completed with a well-commented R script
- Allow everyone a chance to participate in the learning activities, keeping disruption of other students to a minimum while still allowing for fruitful discussion
- All opinions are valued provided they do not harm others
- Everyone is expected to do the work, learning seldom occurs solely by watching someone else do work

Activity 1 - Building an attempt at a plot

We will be looking at the gapminder data set as found in the gapminder package (Bryan 2017). This data has been collected from countries around the world and contains data on life expectancy, population and GDP per capita for 142 countries from 1952 to 2007.

Exercise: Copy and paste the code below to produce a plot showing how the relationship between GDP, life expectancy and population vary over time and continent. If you can't install the gapminder package, you can download the data from Moodle and load it with `read_csv()` from the readr package (loaded when tidyverse is loaded).

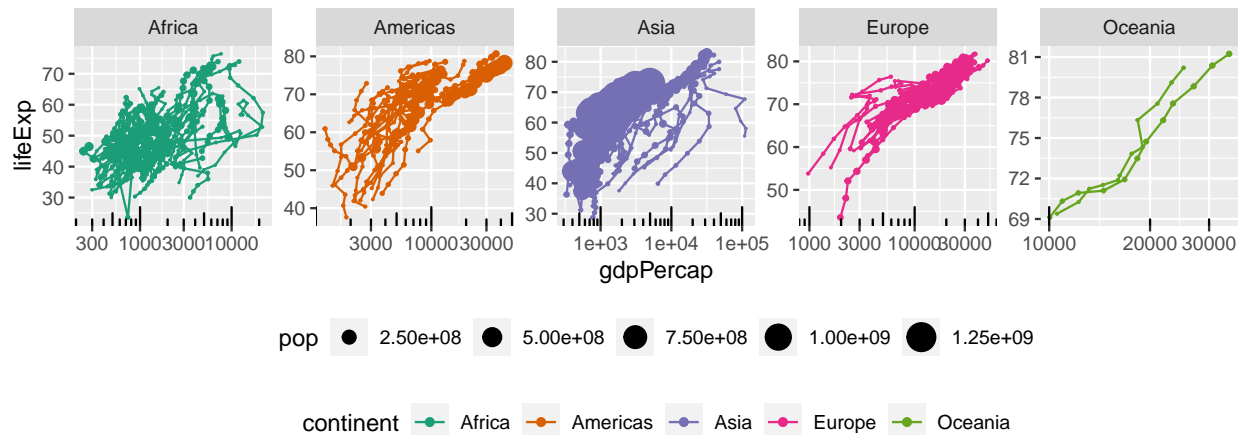
```
library(gapminder)
library(tidyverse)
data(gapminder)

ggplot(data = gapminder,
```

```

aes(x = gdpPerCap, y = lifeExp)) +
geom_path(aes(group = country, color = continent)) +
geom_point(aes(color = continent, size = pop)) +
scale_color_brewer(palette = 'Dark2') + scale_x_log10() +
annotation_logticks(sides = 'b') +
facet_wrap(~ continent, scales = 'free', nrow = 1) +
scale_size_area() +
theme(legend.position = 'bottom',
      legend.box = 'vertical')

```



Exercise: Discuss, within your group, what you think is good and bad about this plot. Does it conform to Tufte's principles of graphical excellence? Is it easy to interpret? Does it show the relationship we are interested in? List *three* important improvements that are needed for this graph to be useful.

Exercise: As a group, discuss what you think each line of code in the above block does. You may wish to answer as comments in your code (everything after a # is a comment) or in a separate document.

Activity 2 – Making a better graph

Based on the ideas discussed, build a graph which your group believes better shows the relationship between life expectancy and GDP. Think first about what story you want your plot to tell; are you interested in trends over space and/or time? Are you interested in a particular continent or even just one country?

You may choose to either modify the code given above or create your own graph from scratch. Make sure your code is written in your script file with appropriate comments.

Some things you may wish to consider:

- fixing up the axis labels
- a relevant title
- a different theme
- different plotting geometries
- different aesthetic options for colour, shape, etc.

You may wish to sketch the graph by hand before attempting to write the R code to generate it. This will help you and your group come to an agreement about the plot you want to make and will help the tutors understand what you're aiming for when you ask them for help.

If you get stuck, look at the [ggplot2](#) documentation or ask a tutor.

Exercise: Make a plot, save it to your computer and write comments in your code or standalone document that outline what the changes you made were and why.

Activity 3 – Making a worse graph

Make a new graph as in the previous activity but make it as bad as possible while still attempting to honestly show the information (i.e. don't add things to the plot which can't be derived from the variables in the plot).

Consider the principles of graphical excellence and how can we go against them to make a truly terrible plot. Think about what was bad about the plot provided earlier. Consider abusing the ability to map graphical options (e.g. color, fill, line type, point size) to our variables of interest.

Exercise: Make a plot, save it to your computer and write comments in your code or standalone document that outline what the changes you made were and why.

Activity 4 – Group discussion

Have participants present their best (and/or worst) graph from the last activities. What did they identify as good and bad and how has each group attempted to present the relationship?

Tidy up

Make sure you save your R script, and anything else you have produced and ensure everyone in your group has a copy. Email your worst graph to [Dr Sam Clifford](#). Leave the room in a better condition than you found it.

Further reading

A lot of the key ideas in data visualisation arose with Tufte (1983), and are summarised by Pantoliano (2012). Some of the history of data visualisation is summarised well by Friendly (2005) and Friendly (2006). [Tuftes website](#) is well worth exploring, particularly the discussion on how the visual presentation of information could have helped avert the *Challenger* disaster (Tufte 1997). For some more guidance on using ggplot2 for data visualisation, check Chapter 3 of Wickham and Grolemund (2017), the RStudio cheatsheets (RStudio 2012), and Chang (2017).

References

- Bryan, Jennifer. 2017. *Gapminder: Data from Gapminder*. <https://CRAN.R-project.org/package=gapminder>.
- Chang, Winston. 2017. *R Graphics Cookbook: Practical Recipes for Visualizing Data*. 2nd ed. O'Reilly Media. <http://www.cookbook-r.com/Graphs/>.
- Friendly, M. 2005. "Milestones in the History of Data Visualization: A Case Study in Statistical Historiography." In *Classification: The Ubiquitous Challenge*, edited by C. Weihs and W. Gaul, 34–52. New York: Springer. <http://www.math.yorku.ca/SCS/Papers/gfkl.pdf>.
- . 2006. "A Brief History of Data Visualization." In *Handbook of Computational Statistics: Data Visualization*, edited by C. Chen, W. Härdle, and A. Unwin. Vol. III. Heidelberg: Springer-Verlag. <http://www.datavis.ca/papers/hbook.pdf>.
- Pantoliano, Mike. 2012. "Data Visualization Principles: Lessons from Tufte." 2012. <https://moz.com/blog/data-visualization-principles-lessons-from-tufte>.
- RStudio. 2012. "RStudio Cheat Sheets." 2012. <https://www.rstudio.com/resources/cheatsheets/>.
- Tufte, Edward R. 1983. *The Visual Display of Quantitative Information*. Graphics Press.
- . 1997. "Visual and Statistical Thinking." In *Visual Explanations: Images and Quantities, Evidence and Narrative*. Graphics Press. https://www.edwardtufte.com/tufte/books_textb.
- Wickham, Hadley, and Garrett Grolemund. 2017. *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data*. O'Reilly Media. <http://r4ds.had.co.nz>.