

Introduction to data summarizing and visualization

Jeff Oliver

09 August, 2019

The R programming language provides many tools for data analysis and visualization, but all the options can be daunting. This lesson provides an introduction to wrangling data in R and using graphics to tell a story with data.

Learning objectives

1. Understand the difference between files and R objects
2. Modify data for proper data hygiene
3. Summarize information from raw data
4. Visualize data to convey information

[DESCRIPTION OR MOTIVATION; 2-4 sentences that would be used for an announcement]

Getting started

The tools: R and RStudio

For this lesson, we will use the R programming language in the RStudio environment. RStudio provides a convenient interface for working with files and packages in R. If you have not done so already, install R and RStudio; details can be found on the installation page.

Preparing our workplace

Key to successful programming is organization. In RStudio, we use Projects to organize our work (a “Project” is really just a fancy name for a folder that contains all our files). For this lesson, we’ll create a new project through the File menu (<style “font-family=‘Trebuchet MS’, ‘Helvetica’, ‘sans-serif’”>File > New Project). In the first dialog, select “New Directory”, and select “New Project” in the second dialog. Next you’ll be prompted to provide a directory name. This will be the name of our project, so we should give it an informative name. For this lesson, we will be using data from [DATA DESCRIPTION], for the directory name, enter ””. We need also to tell RStudio where to put the lesson on our computer; for this lesson, we will place the folder on our Desktop, so it is easy to find. In your own work, you may find it better to place project folders in your Documents folder.

The last thing we need to do to set up our workspace is to use file organization that reinforces best practices. In general, there should be a one-way flow of information: we take information from *data* and write code to produce *output*. We want to avoid any output from messing up our data, so we create separate folders for each. We want to create two folders, one for our data and one for any output, which may include results of statistical analyses or data visualization. In the R console,

```
dir.create("data")
dir.create("output")
```

[Get data]

Data in R

Data *outside* R

[Open file in Excel]

[Load in data]

[QA/QC]

head

Cleaning up

Missing data

take out some NAs

Fixing data

replace some values?

Summarizing data

Using `summary`

Summary statistics for groups

use `mean` and `sd` for `subset(s)` of data

Visualizing data

The `ggplot2` package

`install.packages()`

First rule of plots

Draw it first

Code it second

Simple scatterplot, no colors

Telling the story

Add group colors for story

Additional resources

- Official ggplot documentation
 - A handy cheatsheet for ggplot
 - A PDF version of this lesson
-

[Back to learn-r main page](#)

Questions? e-mail me at jcoliver@email.arizona.edu.