# Agenda

- About R and RStudio
- Hands-on:
  - variables
  - logical expressions
  - values, vectors, and data frames
  - R Studio projects
  - reading in data
  - exploring data
  - data wrangling: cleaning and reshaping
  - data visualization
  - data analysis
  - functions
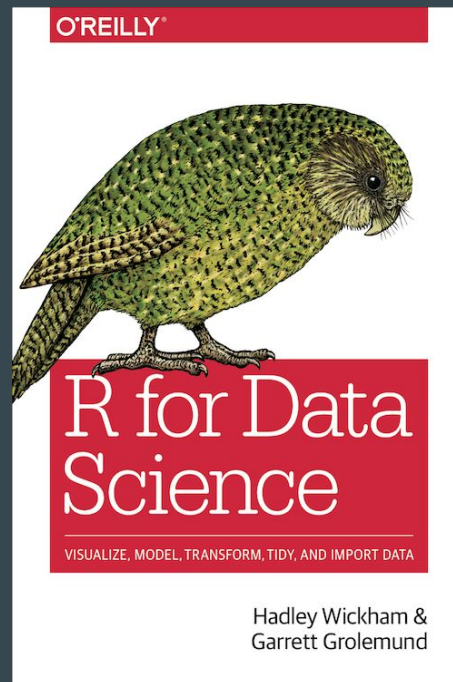- Resources for further learning

# Acknowledgments

A WALK ON THE R SIDE

# Workshop Housekeeping

Ask questions!

Respect every question and person asking the question

Help each other out!

If something is confusing in the workshop,
it probably needs improvement; let us know.
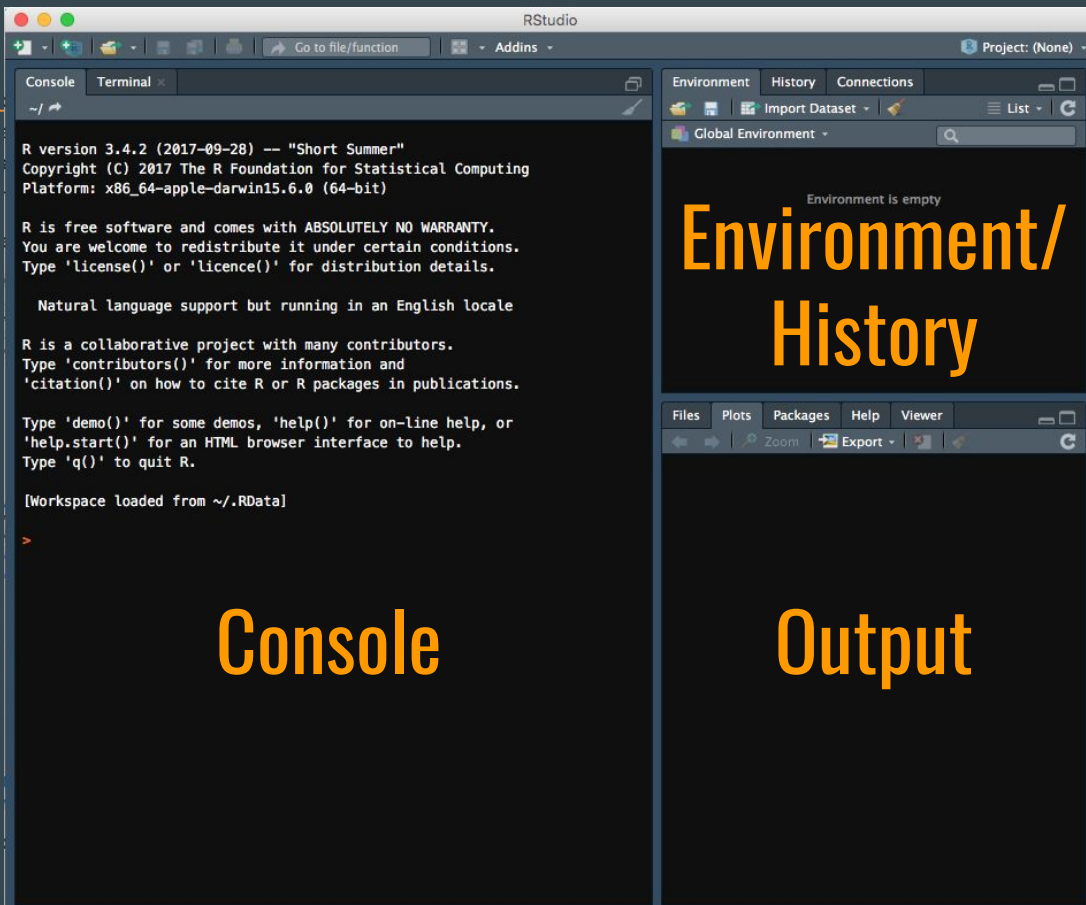
Stay as long as you like

# About R

- Open source
- For statistical computing and graphics
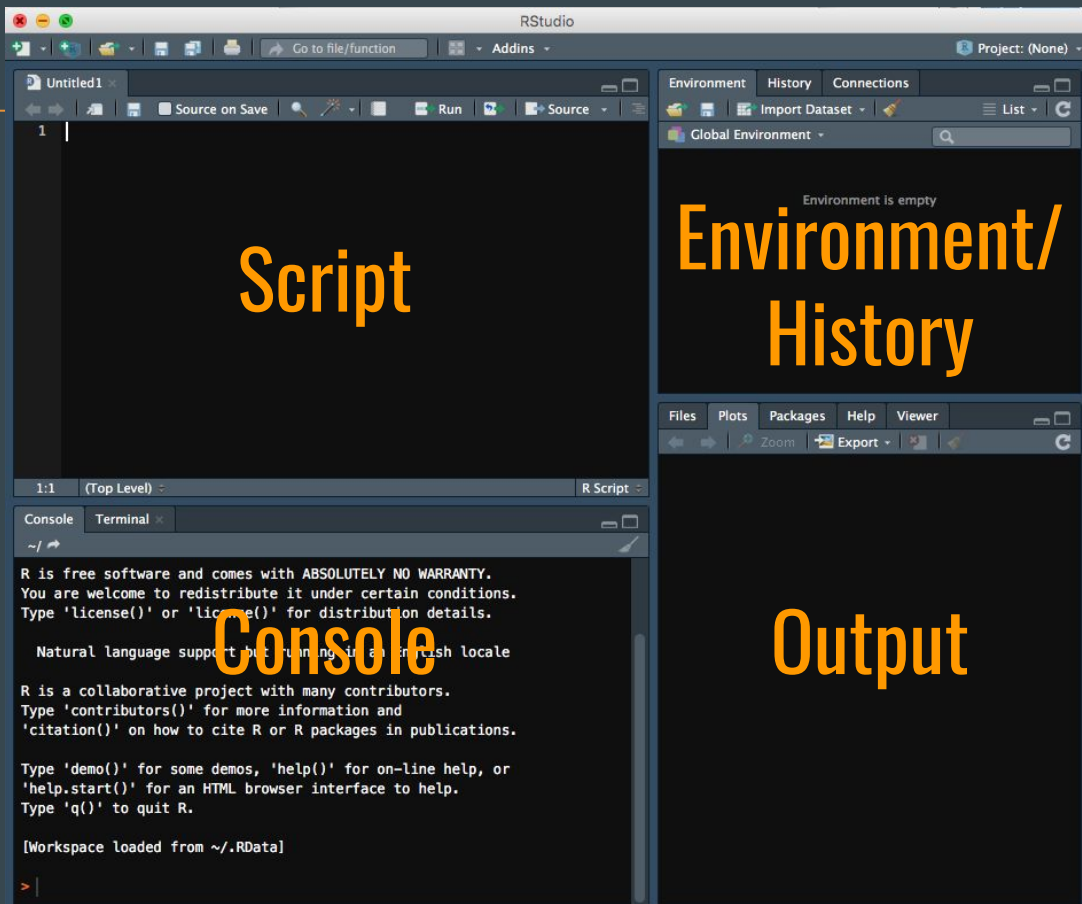- CRAN
  - R packages
  - R events
  - R journal
  - ...

# Other Ways To Use R
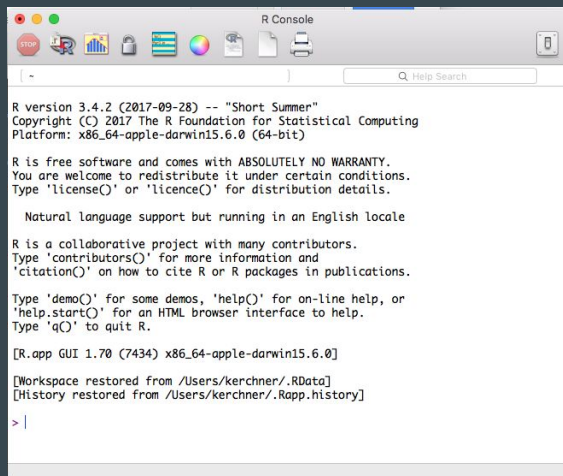
Plain R console

Jupyter Notebook (e.g. in Anaconda)

Let's tRy it!

# Variables

- Try using R as a "calculator" in the Console
  - Try some mathematical functions, too
- Create some variables
  - variable naming
  - `<-`  for assigning values to variables  (Option - on Mac, Alt - on Win)
  - numeric, character, logical
  - Watch the Environment pane!
  - `class()`
  - Coercion w/ `as.integer, as.character, as.logical, as...`

# Logical Expressions

- Operators include:

  ==, <, >, ! (not), & (and), | (or), etc.

# Vectors

# Vectors

- A vector is
  - A sequence of data elements (components) all of the same type.
- Create vectors with c()

Let's pause to explore some useful tabs in RStudio

# Data Frames

# Data Frames

- A data.frame stores a data table
- Comprised of vectors of equal length. <u>Vectors become columns.</u>
- Columns and rows can have names.
- `tibble` (from the tibble package) has some advantages over `data.frame`

# To summarize...

| Value | Vector | | Data Frame | | |
|:---:|:---:|:---:|:---:|:---:|:---:|
| | | | time | temp | boiling |
| 10.2 | 1 | 10.2 | 1 | 51 | 10.2 | FALSE |
| | 2 | 11.3 | 2 | 58 | 11.3 | FALSE |
| | 3 | 11.5 | 3 | 63 | 11.5 | FALSE |
| | 4 | 12.0 | 4 | 70 | 12.0 | TRUE |

A brief word on
`list` and `matrix`

Projects in RStudio

# Projects in RStudio

Recommendations:

- Use [Github for] **version control**!
- Create folders to keep things organized

It's time to import some data!

# Data Importing

- Prepare data as "tidy"
  - rectangular
  - one table per file

- Formats:  CSV, TSV, Excel, Fixed-Width, JSON... and with the right packages:  Stata, SPSS, SAS... (using `foreign`)

- A word about "big data"

# R Packages

A WALK ON THE R SIDE

# Installing and loading R packages

- install.packages('mypackage')
- library('mypackage')  -- or check the box on the Packages tab in RStudio

# Tidyverse Core Packages

## tidyverse.org

- ggplot2 - graphics
- dplyr - data manipulation
- tidyr - tidying data
- readr - reading in data
- tibble - modern data frame
- purrr

# Other often-used R packages

Basic stats functions, like ANOVA ➤ MASS

Mapping ➤ tmap, tmaptools, leaflet

Analyzing 2D and 3D shapes ➤ geomorph

Genomic data ➤ bioconductor

Cluster analyses ➤ cluster

Time series data ➤ forecast

Text mining ➤ qdap, sentimentr, tidytext

Interactive web visualizations ➤ shiny

# Exploring Data

- head, tail
- subsetting
- slicing and dicing

Data Wrangling

# Data Transformation using the dplyr package

- filter()
- arrange()
- select()

- mutate()
- summarize()
- group_by()
- ...

You will want to use a "pipe":  %>%

(shortcut: control-shift-M)

# Data Tidying

- gather()
- spread()
- separate()
- unite()

# Joining

"Merges" tables together

- left_join()
- right_join()
- ...

# Analyzing Data

# Data Visualization

# Data Visualization

- 3 main packages:
  - "base R"
  - lattice
  - ggplot2

# Functions

# Some Handy R Links

# Tutorials

- Software Carpentry:
  - http://swcarpentry.github.io/r-novice-inflammation
  - http://swcarpentry.github.io/r-novice-gapminder
- Data Carpentry:
  - http://datacarpentry.github.io/R-ecology-lesson/
  - http://www.datacarpentry.org/R-genomics/
- Lynda.com  lynda.it.gwu.edu - 3 video courses (~12 hours)
- r-tutor.com/r-introduction

  r-tutor.com/elementary-statistics
- R for Data Science http://r4ds.had.co.nz

# Classes at GW that teach R

- PSC 2012
  Visualizing and Modeling Politics
  Prof. Eric Lawrence
  currently scheduled next for Fall 2018

- Others?

# Reference Links

- [r-project.org](r-project.org)
- R search engine: [rseek.org](rseek.org)
- [rstudio.com](rstudio.com)
  - Cheat Sheets
- [stackoverflow.com](stackoverflow.com)

# Thanks!

- Dan Kerchner            [kerchner@gwu.edu](mailto:kerchner@gwu.edu)
- Dr. Kes Schroer         [schroerk@gwu.edu](mailto:schroerk@gwu.edu)
- Vishwesh Haldevanekar   [vishwesh_s_h@gwu.edu](mailto:vishwesh_s_h@gwu.edu)

These slides:  go.gwu.edu/gwlibrworkshop

Statistics Appointments with Vishwesh:

calendly.com/vishwesh_s_h