

# SPEAKER\_NOTES.md - R for Reproducible Scientific Analysis

Speaker Notes for the 2017-01-11 Software Carpentry R for Reproducible Scientific Analysis lesson

**TYPE ALL EXAMPLES AS YOU GO. THIS KEEPS THE SPEED SANE, AND ALLOWS YOU TO EXPLAIN EVERY STEP.**

START SLIDES WITH `reveal-md slides.md --theme=white`

## Introduction to R and RStudio

---

SLIDE (Why `R` / `RStudio` ?)

- Talk around slide

SLIDE ("But I already know Excel!")

SLIDE ( `R` / `RStudio` presentation)

- Live presentation section
- Everyone start up `RStudio`

### Summarise windows

- Four (maybe three) subwindows:
  - Interactive `R` console
  - Editor (may be missing on startup - will appear when files are opened)
  - Environment/History
  - Files/Plots/Packages/Help

## Create a working directory with version control

- **We're following practices of project management**
  - We'll create a project directory, with `Git` version control
  - Helps ensure data integrity
  - Makes sharing code easier (lab-mates, publication)
  - Easier to recover after a Christmas break
- **Create the new directory LIVE**

- `File->New Project`
- `New Directory`
- `Empty Project`
- Enter sensible name, e.g. `swc-r_reproducible`
- Check box for `Create a git repository`
- `Create project`



Red sticky for a question or issue



Green sticky if complete

- Describe contents of new folder
  - `.gitignore`
  - `.Rproj`

**SLIDE** (Best practices)

- Talk around slide

## Create directory structure

**SLIDE** (Creating files/directories)

- **LIVE PRESENTATION**
- **Create subdirectory for data**
  - In `Files` tab, create `data` subdirectory
- **Create new `R` script**
  - `File -> New File -> R script`
  - save in working directory with sensible name, e.g. `swc-script.R`



Red sticky for a question or issue



Green sticky if complete

- **LIVE VERSION CONTROL EXAMPLE**
- **Show Git tab on right**
- **Stage files**
  - Three files shown (including `.gitignore` and the new script file)
  - Yellow status markers mean they're not in the repository
  - Click check-boxes to stage them

- Note that **we don't version disposable output**

- **Commit files**

- Click `Commit`
- Describe new dialogue window
- Show contents/changes to files
- Add commit message ("initialise repository") - good commit messages are short and imperative
- Commit
- Show commit summary
- Exit



Red sticky for a question or issue



Green sticky if complete

### **SLIDE** (Challenge 1)

Run through challenge (5min?) - hint about editing `.gitignore`

- Right-click link on presentation and download to `data`
- Create `graphs` subdirectory in `Files` tab
- Edit `.gitignore` to add `graphs/` folder and save
- Stage `.gitignore` in Git tab
- Commit in Git tab, and add appropriate commit message
- Demo History window for Git

### **SLIDE** (`R` as a calculator)

## **Interacting with `R`**

- **Two ways**

- Type commands in the console
- Use the script editor and save the script

- **Console**

- Output shown here
- Good for experimentation
- Commands 'forgotten' when you close a session

- **Script**

- Keeps record of what you did

- Easier to reproduce and share

## Working at the console

- `R` shows a `>` if it is expecting input

```
1 > 1 + 100
2 [1] 101
```

- `R` shows `+` if it's waiting for completion ( `Esc` to exit)

```
1 > 1 +
2 +
```

## Working from script file

- Can write same commands in the script file ( `1 + 100` )
  - Use `Run` to execute
  - Use `Ctrl-Enter` to execute
  - Output appears in the console
  - Show `#` comments - good practice to comment
  - More examples (order of precedence):

```
1 > 3 + 5 * 2
2 [1] 13
3 > (3 + 5) * 2
4 [1] 16
```

- Show `Source` operation: add the following lines to script:

```
1 # Using R as a calculator script demo
2 1 + 100
3 3 + 5 * 2
4 (3 + 5) * 2
```

- Run script

```
1 > # Using R as a calculator script demo
2 > 1 + 100
3 [1] 101
4 > 3 + 5 * 2
5 [1] 13
6 > (3 + 5) * 2
7 [1] 16
```

- More examples
  - scientific notation

```
1 > 1/40
2 [1] 0.025
3 > 2/10000
4 [1] 2e-04
5 > 5e3
6 [1] 5000
```

## Mathematical functions

- General format: `fn(arg)`
  - autocompletion - example: `factorial(6)`

```
1 > sin(1)
2 [1] 0.841471
3 > log(1)
4 [1] 0
5 > log10(10)
6 [1] 1
7 > exp(0.5)
8 [1] 1.648721
```

## Comparisons

- Return `TRUE` / `FALSE` logical values

```
1 > 1 == 1
2 [1] TRUE
3 > 1 == 2
4 [1] FALSE
5 > 1 != 2
6 [1] TRUE
7 > 1 < 2
8 [1] TRUE
9 > 1 > 2
10 [1] FALSE
11 > 1 <= 2
12 [1] TRUE
13 > 1 >= 2
14 [1] FALSE
```

- Computer representation of numbers is approximate: important for comparisons
  - Any physicists/computer scientists in the room?
  - Numbers may not be equal, but be 'the same'
  - Use `all.equal` instead of `==`

```
1 > all.equal(pi-1e-7, pi)
2 [1] "Mean relative difference: 3.183099e-08"
3 > all.equal(pi-1e-8, pi)
4 [1] TRUE
5 > pi-1e-8 == pi
6 [1] FALSE
```

## Variables and assignment

- **Variables hold values**, just like in Python
- Two ways to assign variables: `<-` and `=`
  - The `<-` form is more widely used
  - Consistency more important than choice

```
1 > x <- 1/40
2 > x
3 [1] 0.025
4 > x = 1/40
5 > x
6 [1] 0.025
```

- **Look at the Environment tab** automatic updates

```
1 > x <- 100
```

- Variables can be used as arguments to functions

```
1 > log(x)
2 [1] 4.60517
3 > sqrt(x)
4 [1] 10
```

- Variables can be used to reassign values to themselves

```
1 > x
2 [1] 100
3 > x <- x + 1
4 > x
5 [1] 101
```

### SLIDE (Good variable names)

- Talk around slide

### SLIDE (MCQ1)

- Pose question

## Package management

### SLIDE (Package Management)

- See what packages are installed with `installed.packages()`
  - **demo this one**
- Add a new package using `install.packages("packagename")`
  - **demo this one with** `install.packages("ggplot2")`
- Update packages with `update.packages()`
  - **demo this one**
- You can remove a package with `remove.packages("packagename")`
- To make a package available for use, use `library(packagename)`
  - **demo**
  - Note that there are no quotes, this time

```
1 > ggplot()
2 Error: could not find function "ggplot"
3 > library(ggplot2)
4 Warning message:
5 package 'ggplot2' was built under R version 3.2.3
6 > ggplot()
7 Warning message:
8 In max(vapply(evald, length, integer(1))) :
9   no non-missing arguments to max; returning -Inf
```

### SLIDE (Challenge 2)

Solution:

```
1 install.packages("plyr")
2 install.packages("gapminder")
3 install.packages("dplyr")
4 install.packages("tidyr")
```

## Getting help for functions

### SLIDE (Functions, and getting help)

- Talk around slide

- Demo: `round(3.14159)` :
  - argument: `3.14159`
  - value: `3`

```
1 > round(3.14159)
2 [1] 3
```

## SLIDE (Getting help for functions)

- **Carrying on with** `round()` **from last slide**
- What other arguments can `round()` take?
  - Use `args(fname)`

```
1 > args(round)
2 function (x, digits = 0)
3 NULL
```

- Can use the `digits` argument by naming it, or not (but order matters)

```
1 > round(3.14159, digits=2)
2 [1] 3.14
3 > round(3.14159, 2)
4 [1] 3.14
```

- **Best practice:** always use the argument name
  - clearer to others
  - if function changes, order may change
  - difficult to remember the purpose of each argument, if not explicit
- **What does a function do?**
  - Use `?fname` or `help(fname)` to get the complete help text
  - Demo: `?round` - go through main points
- **What package is my function in?**
  - (i.e. I can't find it, and don't know what to install)
  - Demo: `??melt` - show that we need `reshape2`
- **Is there a function that does X?**
  - e.g. you know the name of a test, such as Kolmogorov-Smirnov
  - Demo: `help.search("smirnov")`, `?ks.test`



## SLIDE (Where can I get more help?)

- Talk around slide

## SLIDE (Asking the right questions)

- Talk around slide
- For `dput()` example use `dput(head(iris))`
- Demo `sessionInfo()`

# Functions

---

## SLIDE (Functions)

## SLIDE (Learning objectives)

- Talk around slide
- **Why functions?**
  - You've already seen the power of functions, for encapsulating complex analyses into simple commands
  - Functions work similarly in `R` as they do in the shell/Python

## SLIDE (What is a function?)

- Talk around slide

## Defining a function

## SLIDE (Defining a function)

- Talk around slide
- **Create a new `R` script file to hold functions**
  - `File -> New File -> R Script`
  - `File -> Save -> functions-lesson.R`
  - Check what's happened in Git tab
- **Write new function in script**
  - Describe parts of function:
  - *prototype* with inputs

- code block/body
- indentation (readability)
- addition, and return statements
- function scope, internal variables (readability)
- assignment of function to variable
- comments (readability)

```

1 # Returns sum of two inputs
2 my_sum <- function(a, b) {
3   the_sum <- a + b
4   return(the_sum)
5 }
6 # Converts fahrenheit to Kelvin
7 fahr_to_kelvin <- function(temp) {
8   kelvin <- ((temp - 32) * (5 / 9)) + 273.15
9   return(kelvin)
10 }

```

#### • Run the functions

- `source` the script
- tab-completion works!
- boiling and freezing points

```

1 > fahr_to_kelvin(32)
2 [1] 273.15
3 > fahr_to_kelvin(212)
4 [1] 373.15

```

#### SLIDE (Challenge 1)

Solution:

```

1 kelvin_to_celsius <- function(temp) {
2   celsius <- temp - 273.15
3   return(celsius)
4 }

```

#### SLIDE (Challenge 2)

Solution:

```

1 fahr_to_celsius <- function(temp) {
2   kelvin <- fahr_to_kelvin(temp)
3   celsius <- kelvin_to_celsius(kelvin)
4   return(celsius)
5 }

```

## INSERTED EXAMPLE

- Just as in Python, we can use `for` loops to apply a function to several values
- Avoids repetition

```
1 for (i in 32:100) {  
2   print(fahr_to_celsius(i))  
3 }
```

- Can also apply functions to vectors

```
1 fahr_to_celsius(32:100)
```

- Also `if` and `if/else` statements, as in Python:

```
1 if (5 > 1) {  
2   print("condition is true")  
3 }
```

```
1 if (5 < 1) {  
2   print("condition is true")  
3 } else {  
4   print("condition is false")  
5 }
```

- **COMMIT TO LOCAL GIT REPO**

## SLIDE (Testing functions)

- Talk around slide
- **Known good values**
  - water freezes at 32F/0C, boils at 212F/100C

```
1 > fahr_to_celsius(32)  
2 [1] 0  
3 > fahr_to_celsius(212)  
4 [1] 100
```

- **Known bad values**
  - All values are fair game on Fahrenheit/Celsius, but can't go below 0K

```
1 > kelvin_to_celsius(-10)  
2 [1] -283.15
```

- **We'd need to modify this for real use!**

**SLIDE** (Not the best approach...)