

Introduction to R: R Basics

Session 1, Part A

Nick Graetz¹

¹ University of Pennsylvania, Population Studies Center

9/4/2020

IN THIS LECTURE

1. What is R?
2. RStudio interface
3. Packages
4. R as calculator
5. Anatomy of a function
6. Help files
7. R scripts

WHAT IS R?

- ▶ R is a language for statistical computing and graphics
- ▶ Originally developed in 1992 by Robert Gentleman and Ross Ihaka based on the programming language S
- ▶ The core of the R language is maintained by the R Core Team
- ▶ A (very) large number of packages which add additional functionality are maintained by other contributors

WHY USE R?

- ▶ R can do many useful things
 - ▶ Flexible data management
 - ▶ Powerful statistical capabilities, particularly for modeling
 - ▶ Extensive graphics capabilities

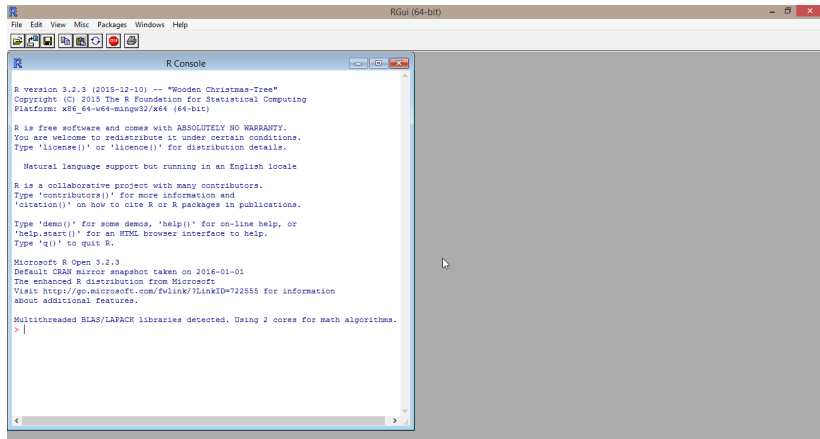
WHY USE R?

- ▶ R can do many useful things
 - ▶ Flexible data management
 - ▶ Powerful statistical capabilities, particularly for modeling
 - ▶ Extensive graphics capabilities
- ▶ R is free software
 - ▶ You don't have to pay for it (and you can share it with anyone)
 - ▶ You can use and modify it as you see fit

WHY USE R?

- ▶ R can do many useful things
 - ▶ Flexible data management
 - ▶ Powerful statistical capabilities, particularly for modeling
 - ▶ Extensive graphics capabilities
- ▶ R is free software
 - ▶ You don't have to pay for it (and you can share it with anyone)
 - ▶ You can use and modify it as you see fit
- ▶ R has a large (and enthusiastic) user base
 - ▶ This makes finding help relatively straightforward
 - ▶ New methods are often implemented in R very quickly

R (GUI) INTERFACE



WHAT IS RSTUDIO?

- ▶ “Integrated development environment”

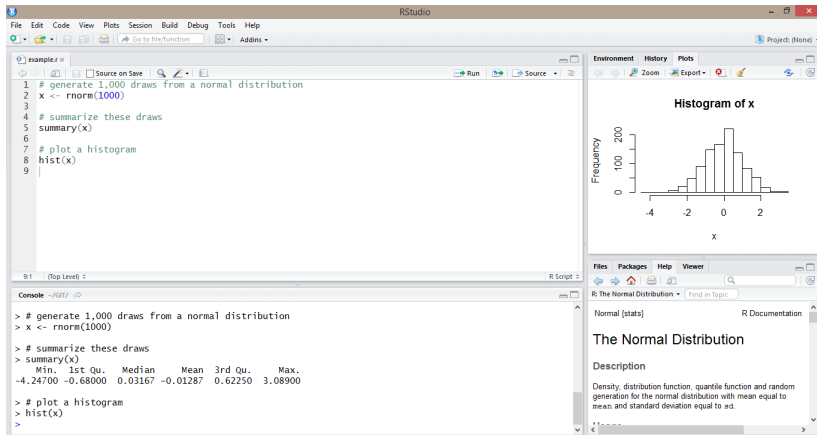
WHAT IS RSTUDIO?

- ▶ “Integrated development environment”
- ▶ Convenient interface for R which incorporates a number of useful features for developing code
 - ▶ syntax highlighting
 - ▶ code completion
 - ▶ code navigation
 - ▶ debugging tools
 - ▶ etc.

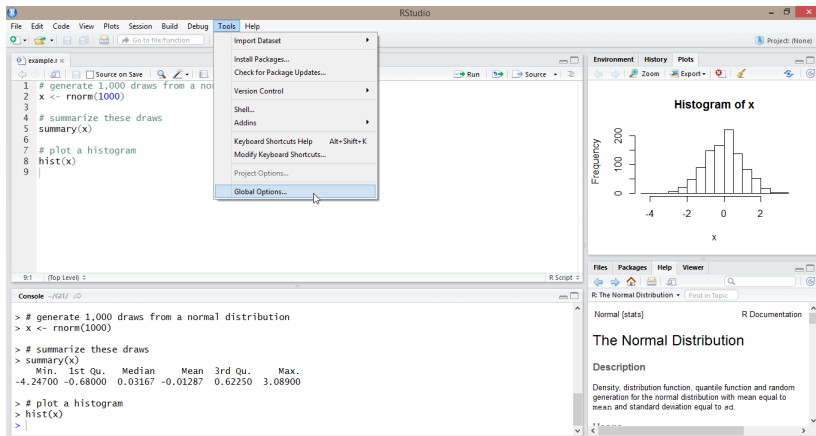
WHAT IS RSTUDIO?

- ▶ “Integrated development environment”
- ▶ Convenient interface for R which incorporates a number of useful features for developing code
 - ▶ syntax highlighting
 - ▶ code completion
 - ▶ code navigation
 - ▶ debugging tools
 - ▶ etc.
- ▶ Also provides integration with other useful tools
 - ▶ Shiny (for developing web apps)
 - ▶ R Markdown (for authoring documents and slides)
 - ▶ Git/Subversion (for version control)

RSTUDIO INTERFACE



RSTUDIO INTERFACE



The screenshot displays the RStudio application window. The 'Tools' menu is open, showing options such as 'Import Dataset', 'Install Packages...', 'Check for Package Updates...', 'Version Control', 'Shell...', 'Addins', 'Keyboard Shortcuts Help', 'Modify Keyboard Shortcuts...', 'Project Options...', and 'Global Options...'. The 'Global Options...' option is highlighted by the mouse cursor.

The background shows the R script editor with the following code:

```
1 # generate 1,000 draws from a normal distribution
2 x <- rnorm(1000)
3
4 # summarize these draws
5 summary(x)
6
7 # plot a histogram
8 hist(x)
9
```

The Console window at the bottom shows the output of the executed code:

```
> # generate 1,000 draws from a normal distribution
> x <- rnorm(1000)
> # summarize these draws
> summary(x)
   Min.  1st Qu.  Median    Mean 3rd Qu.    Max.
-4.24700 -0.68000  0.03167 -0.01287  0.62250  3.08900
> # plot a histogram
> hist(x)
>
```

The Environment window on the right displays a histogram titled 'Histogram of x'. The x-axis is labeled 'x' and ranges from -4 to 2. The y-axis is labeled 'Frequency' and ranges from 0 to 200. The histogram shows a bell-shaped distribution centered around 0.

The Files window at the bottom right shows the 'R Documentation' for 'The Normal Distribution'. The description states: 'Density, distribution function, quantile function and random generation for the normal distribution with mean equal to mean and standard deviation equal to sd.'

RSTUDIO INTERFACE

The screenshot displays the RStudio application window. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Tools, and Help. The main workspace is divided into four panes: Source (left), Environment/History/Plots (top right), Files/Packages/Help/Viewer (bottom right), and Console (bottom left).

The **Options** dialog box is open, showing the **General** tab. The R version is [Default] [64-bit] C:\Program Files\Microsoft\MSO\R-3.2.5. The default working directory is set to ~/GIT. The following options are checked: Re-use idle sessions for project links, Restore most recently opened project at startup, Restore previously open source documents at startup, and Restore .RData into workspace at startup. The 'Save workspace to .RData on exit' dropdown is set to 'Never'. Other options include 'Always save history', 'Remove duplicate entries in history', 'Show .Last.value in environment listing', 'Use debug error handler only when my code contains errors', 'Automatically expand tracebacks in error inspector', and 'Automatically notify me of updates to RStudio'.

The **Source** pane shows a script with the following code:

```
1 # generate 1,000 draws from a normal d
2 x <- rnorm(1000)
3
4 # summarize these draws
5 summary(x)
6
7 # plot a histogram
8 hist(x)
9
```

The **Console** pane shows the output of the code:

```
> # generate 1,000 draws from a normal dis
> x <- rnorm(1000)
>
> # summarize these draws
> summary(x)
   Min.  1st Qu.  Median    Mean 3rd Qu.
-4.24700 -0.68000  0.03167 -0.01287  0.622
> # plot a histogram
> hist(x)
>
```

The **Plots** pane displays a histogram titled "Histogram of x". The x-axis is labeled "X" and ranges from -4 to 2. The y-axis is labeled "Frequency" and ranges from 0 to 200. The histogram shows a distribution of data points.

The **Files/Packages/Help/Viewer** pane shows the "R The Normal Distribution" page, which includes a description of the normal distribution and links to the R documentation.

RSTUDIO INTERFACE

The screenshot displays the RStudio application window. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Tools, and Help. The toolbar below the menu contains icons for file operations and development tools. The main workspace is divided into four panes:

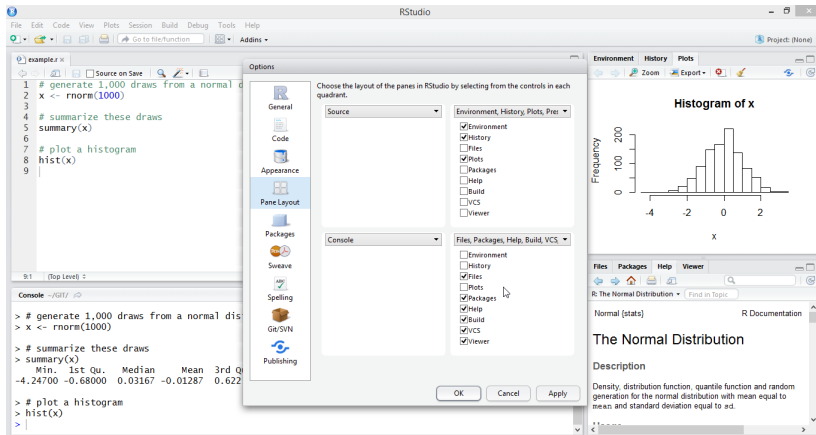
- Source Editor:** Contains an R script with the following code:

```
1 # generate 1,000 draws from a normal distribution
2 x <- rnorm(1000)
3
4 # summarize these draws
5 summary(x)
6
7 # plot a histogram
8 hist(x)
9
```
- Console:** Shows the output of the executed code:

```
> # generate 1,000 draws from a normal distribution
> x <- rnorm(1000)
> # summarize these draws
> summary(x)
      Min.   1st Qu.   Median     Mean   3rd Qu.
-4.24700  -0.68000   0.03167  -0.01287   0.62200
> # plot a histogram
> hist(x)
>
```
- Options Dialog:** The 'Appearance' tab is selected. It shows settings for Zoom (100%), Editor font (Lucida Console), Font size (11), and a list of editor themes. 'Eclipse' is currently selected.
- Plots Pane:** Displays a histogram titled 'Histogram of x'. The x-axis is labeled 'x' and ranges from -4 to 2. The y-axis is labeled 'Frequency' and ranges from 0 to 200.

The bottom right pane shows the 'R Documentation' for 'The Normal Distribution', including a description of the density, distribution function, quantile function, and random generation for the normal distribution with mean equal to μ and standard deviation equal to σ^2 .

RSTUDIO INTERFACE



PACKAGES

Most basic R functionality is part of `base` and is loaded automatically when you start R. Additional functionality can be added through packages.

Most basic R functionality is part of `base` and is loaded automatically when you start R. Additional functionality can be added through packages.

The first time you use a package, it needs to be installed:

```
> install.packages ("ggplot2")
```

After that, you just need to load the package using the `library()` command whenever you start a new instance of R:

```
> library (ggplot2)
```

R can be used as a calculator by just typing in the console.

All of the basic arithmetic operators (+, -, *, /, ^) do what you would expect them to do, following normal order of operations conventions:

```
> 230 + 97  
[1] 327  
> 500/20  
[1] 25
```

Parentheses can be used to alter the order of operations:

```
> 300/20^1/2
```

```
[1] 7.5
```

```
> (300/20)^(1/2)
```

```
[1] 3.872983
```

R AS CALCULATOR: QUICK EXERCISE

1. How many seconds are in September?
2. What is 80 degrees Fahrenheit in degrees Celsius?
3. How much longer is 1 mile than 1600 meters (in feet)?

R AS CALCULATOR: QUICK EXERCISE

1. How many seconds are there in September?

```
> 30 * 24 * 60 * 60  
[1] 2592000
```

2. What is 80 degrees Fahrenheit in degrees Celsius?

```
> (80 - 32) * (5/9)  
[1] 26.66667
```

3. How much longer is 1 mile than 1600 meters (in feet)?

```
> 5280 - 1600 * 3.28084  
[1] 30.656
```

FUNCTIONS

R functions are used to transform input into output in some way.

For example...

```
> log(10)
[1] 2.302585
```

```
> exp(3)
[1] 20.08554
```

```
> sqrt(80)
[1] 8.944272
```

FUNCTIONS: ANATOMY

```
> log(x = 300, base = 10)
[1] 2.477121
```

1. Function name: **log()**
2. Argument name(s): **x, base**
3. Argument value(s): **300, 10**
4. Output: **2.4771213**

FUNCTIONS: ARGUMENT ORDER

Arguments can be specified in any order *if they are named*:

```
> log(x = 300, base = 10)
[1] 2.477121
```

```
> log(base = 10, x = 300)
[1] 2.477121
```


FUNCTIONS: ARGUMENT NAMES

Arguments don't need to be named, but then *there is only one correct order*:

```
> log(x = 300, base = 10)
```

```
[1] 2.477121
```

```
> log(base = 10, x = 300)
```

```
[1] 2.477121
```

```
> log(300, 10)
```

```
[1] 2.477121
```

```
> log(10, 300)
```

```
[1] 0.4036944
```

FUNCTIONS: DEFAULTS

Some (but not all) arguments have defaults and don't need to be specified, assuming you are happy with the default:

```
> log(x = 300)
[1] 5.703782
```

```
> log(base = 10)
```

```
Error in eval(expr, envir, enclos): argument "x" is missing, w
```

FUNCTIONS: COMBINING

Functions can be combined or nested with other functions and operators:

```
> exp(log(10) + log(10))  
[1] 100
```

```
> log(x = (4 * 10) / 7, base = 10)  
[1] 0.756962
```

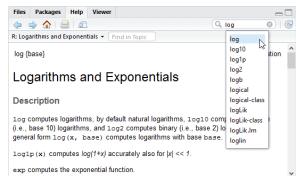
HELP FILES

Every function has a help file.

You can access a help file from the console:

```
> help(log)
```

or from the help tab in RStudio:



`log {base}`

R Documentation

Logarithms and Exponentials

Description

`log` computes logarithms, by default natural logarithms, `log10` computes common (i.e., base 10) logarithms, and `log2` computes binary (i.e., base 2) logarithms. The general form `log(x, base)` computes logarithms with base `base`.

`log1p(x)` computes $\log(1+x)$ accurately also for $|x| \ll 1$.

`exp` computes the exponential function.

`expm1(x)` computes $\exp(x) - 1$ accurately also for $|x| \ll 1$.

Usage

```
log(x, base = exp(1))  
logb(x, base = exp(1))  
log10(x)  
log2(x)
```

```
log1p(x)
```

```
exp(x)  
expm1(x)
```

Arguments

`x`
a numeric or complex vector.

`base`
a positive or complex number: the base with respect to which logarithms are computed. Defaults to $e = \exp(1)$.

Details

All except `logb` are generic functions: methods can be defined for them individually or via the [Math](#) group generic.

`log10` and `log2` are only convenience wrappers, but logs to bases 10 and 2 (whether computed *via* `log` or the wrappers) will be computed more efficiently and accurately where supported by the OS. Methods can be set for them individually (and otherwise methods for `log` will be used).

`logb` is a wrapper for `log` for compatibility with S. If (S3 or S4) methods are set for `log` they will be dispatched. Do not set S4 methods on `logb` itself.

All except `log` are [primitive](#) functions.

Value

A vector of the same length as `x` containing the transformed values. `log(0)` gives `-Inf`, and `log(x)` for negative values of `x` is `NaN`. `exp(-Inf)` is 0.

For complex inputs to the `log` functions, the value is a complex number with imaginary part in the range $[-\pi, \pi]$: which end of the range is used might be platform-specific.

References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole. (for `log`, `log10` and `exp`.)

Chambers, J. M. (1998) *Programming with Data. A Guide to the S Language*. Springer. (for `logb`.)

See Also

[Trig](#), [sqrt](#), [Arithmetic](#).

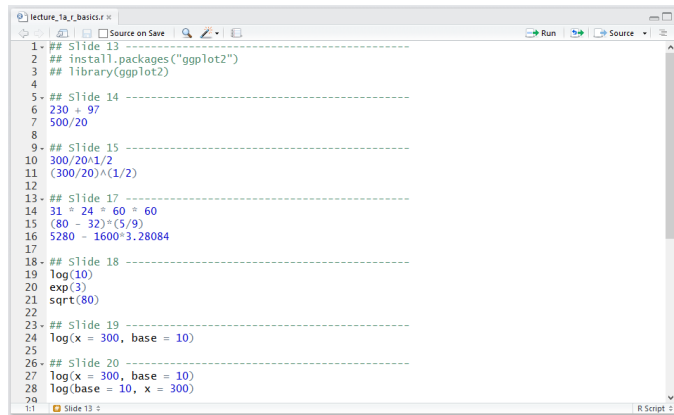
Examples

```
log(exp(3))  
log10(1e7) # = 7
```

```
x <- 10^(1+2*1:9)  
cbind(x, log(1+x), log1p(x), exp(x)-1, expm1(x))
```


R SCRIPTS

An R script is a text file (`.r` extension) with a series of R commands and (ideally) some useful commentary.



```
lecture_1a_r_basics.r x
Source on Save
Run Source
1- ## Slide 13 -----
2  ## install.packages("ggplot2")
3  ## library(ggplot2)
4
5- ## Slide 14 -----
6  230 + 97
7  500/20
8
9- ## Slide 15 -----
10 300/20^1/2
11 (300/20)^(1/2)
12
13- ## Slide 17 -----
14 31 * 24 * 60 * 60
15 (80 - 32)*(5/9)
16 5280 - 1600*3.28084
17
18- ## Slide 18 -----
19 log(10)
20 exp(3)
21 sqrt(80)
22
23- ## Slide 19 -----
24 log(x = 300, base = 10)
25
26- ## Slide 20 -----
27 log(x = 300, base = 10)
28 log(base = 10, x = 300)
29
1:1 Slide 13 ▾ R Script ▾
```

WHY USE A SCRIPT?

Typing in the console is fine for quick calculations or experimentation with a command, but a script provides...

- ▶ a full record of all commands required to carry out an analysis
- ▶ a convenient mechanism for repeating an analysis without needing to retype everything (no need to reinvent the wheel)
- ▶ a starting point for writing new code
- ▶ a vehicle for providing context and commentary for your code

WHY USE A SCRIPT?

Any analysis you do should be saved as a script!

Without a script...

- ▶ you will forget what you've done
- ▶ you will forget why you did it
- ▶ no one else will ever know what you did or why you did it
- ▶ you will have to do things over again for no reason

RUNNING A SCRIPT

If your script is open in RStudio, you can run the whole thing using `ctrl + shift + enter` or just a single line (or highlighted block) using `ctrl + enter`.

Or you can run a script from the command line using the `source()` function:

```
> source(file="J:/temp/bootcamp_r_training/lectures/lecture_1a_r_basics.R")
```

COMMENTING A SCRIPT

R will ignore any line in a script that starts with #, so you can use this to add comments to your code:

```
> # add 1-5
> 1 + 2 + 3 + 4 + 5
[1] 15
>
> # find the natural log of 10
> log(10)
[1] 2.302585
```

COMMENTING A SCRIPT

Use comments to:

- ▶ Label blocks of code. This will help you navigate your code later
- ▶ Explain why you're doing something (if it's not self-evident)
- ▶ Write yourself (and other users) notes about particularly tricky lines of code

COMMENTING A SCRIPT

Use comments to:

- ▶ Label blocks of code. This will help you navigate your code later
- ▶ Explain why you're doing something (if it's not self-evident)
- ▶ Write yourself (and other users) notes about particularly tricky lines of code

You want to provide enough information so that your future self, or someone else, can quickly understand the structure and purpose of your code at a later date.

However, it is possible to provide too much information, making your code more cumbersome (e.g., writing out what each line of code does).

It's also good practice to use '#' to provide some sort of header at the top of your code:

```
#####  
## Author:      John Doe  
##  
## Description: A short description of what this code does  
##              and any important context for why.  
##  
## Output:      A list of files that are output by this  
##              code.  
##  
## Notes:       Anything someone should know when running  
##              this code.  
#####
```