

Preliminaries in R

What is R?

What is R?

- A programming language good for data analysis / statistics
- A base package of some software with many available user-created add-ons (packages)
- Free and open-source software
- A scripting language

R is currently popular in a number of fields, including:

- Statistics / Biostatistics
- Machine learning
- Data journalism
- Ecology
- Financial engineering
- Bioinformatics

What is R?

R is a programming language popular for statistical computing.

"The best thing about R is that it was developed by statisticians. The worst thing about R is that... it was developed by statisticians."

-Bo Cowgill, Google, at the Bay Area R Users Group

Other programming languages popular for statistical computing include:

- SAS
- SPSS
- Matlab
- Julia
- Python

Free and open-source software

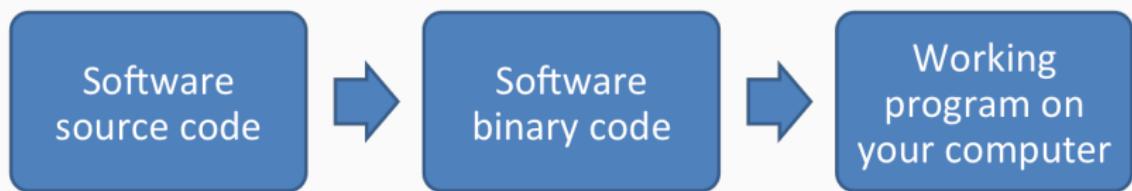


Source: Ruben Rodriguez. licensed under the [Creative Commons Attribution 4.0 International license](#).

[Stallman story]

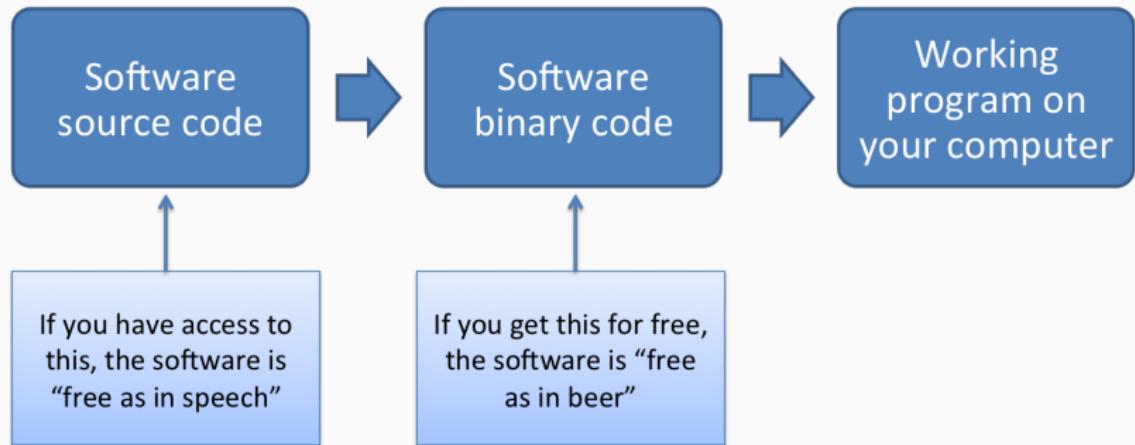
Free and open-source software

How a lot of software is created:



Free and open-source software

A basic sketch of how software can be “free”:

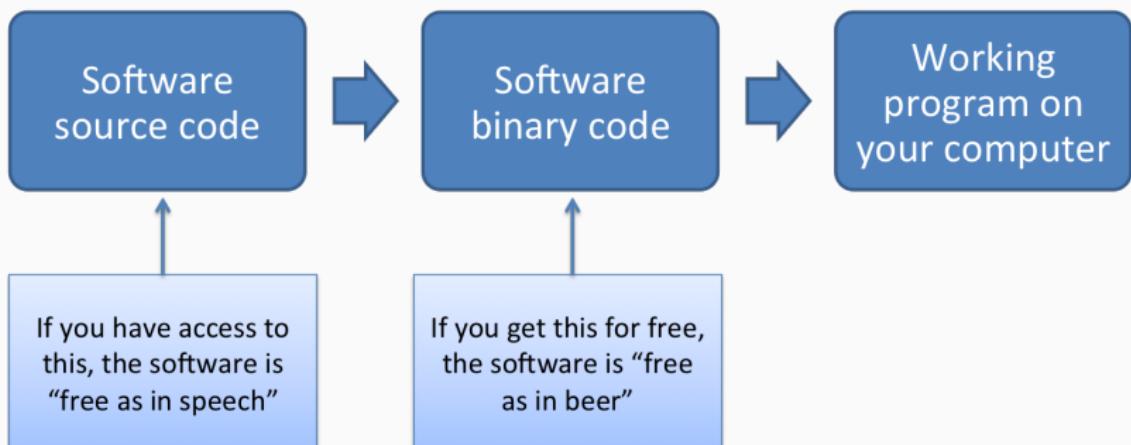


- **Gratis:** Free as in beer
- **Libre:** Free as in speech

Free and open-source software

With open-source software (free as in speech), you can:

- Check out the code to figure out how the software works
- Share the code (and software) with other people
- Make any changes you want to the code



Free and open-source software

"Despite its name, open-source software is less vulnerable to hacking than the secret, black box systems like those being used in polling places now. That's because anyone can see how open-source systems operate. Bugs can be spotted and remedied, deterring those who would attempt attacks. This makes them much more secure than closed-source models like Microsoft's, which only Microsoft employees can get into to fix."

Woolsey and Fox. *To Protect Voting, Use Open-Source Software*. New York Times. August 3, 2017.

Free and open-source software

Funding agencies are starting to ask for grant proposals to develop open-source tools. For example, a recent call from the NIH asks for:

“Open-source, generalizable, and scalable bioinformatics tools”

NIH RFA-RM-17-012: “Metabolomics Data Analysis and Interpretation Tools (U01)”

Scripting languages

- Point-and-click software
- Scripting language
- Compiled language
- Assembly language
- Machine code

Scripting languages

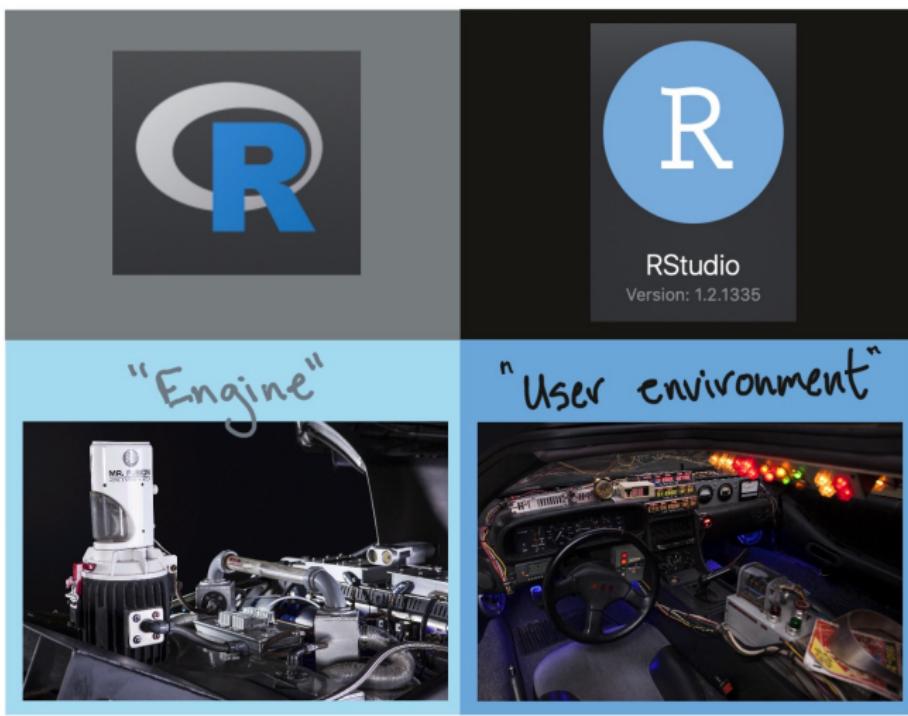
Scripting languages

RStudio vs. R

In this class, I'm having you download both R and RStudio. It's helpful for you to know the difference between the two.

R provides the engine, while **RStudio** provides a nice place to work while using that engine (the leather interior, say).

RStudio vs. R



RStudio vs. R

[Engine versus interior]

RStudio vs. R



When you open **R**, it does have a user interface, it's just not as nice as RStudio's.

RStudio vs. R

RStudio (the software) is an integrated development environment (IDE) for R. You download it separately from R, but it's a "nicer" way to work in R.

This IDE includes:

- An interface with "panes" for key tasks you'll be doing (e.g., one pane with the R console, one for scripts, one to view graphs)
- Code highlighting
- Version control (git) and interface with GitHub
- Tools for Shiny web app development
- Tools for R package development

RStudio vs. R

RStudio (the company) is a leader in the R community. Currently, the company:

- Develops and freely provides the RStudio IDE
- Provides excellent resources for learning and using R (e.g., cheatsheets, free online books)
- Is producing some of the most-used R packages
- Employs some of the top people in R development
- Is a key member of The R Consortium (others include Microsoft, IBM, and Google)

Setting up

If do not already have them, you will need to download and install both R and RStudio.

- Go to <https://cran.r-project.org> and download the latest version of R for your system. Install.
- Go to the [RStudio download page](#) and download the latest version of RStudio for your system. Install.
- Defaults should be fine for everything.

How to “talk” to R

How to “talk” to R

1. Open an **R session**
2. At the **prompt** in the **console**, enter an **R expression**
3. Read R's “response” (the **output**)
4. Repeat 2 and 3
5. Close the R session

Opening an R session

An **R session** is an instance of you using R.

Double-click on the icon for “RStudio” on your computer. When RStudio opens, you will be in a “fresh” R session, unless you restore a saved session (which I strongly recommend against).

This means that, once you open RStudio, you will need to “set up” your session, including loading any packages you need (which we’ll talk about later) and reading in any data (which we’ll also talk about).

The prompt in the console

In RStudio, there screen is divided into several “panes”. We’ll start with the pane called “Console”.

The **console** lets you “talk” to R. This is where you type an **expression** at the **prompt** (the caret symbol, “>”). You press the “Return” key to submit this message to R.

This is called . . .

How R might respond

Once you press “Return”, R will respond in one of three ways:

1. R does whatever you asked it to do with the expression and prints the output (if any) of doing that, as well as a new prompt so you can ask it something new
2. R doesn't think you've finished asking you something, and instead of giving you a new prompt (“>”) it gives you a “+”. This means that R is still listening, waiting for you to finish asking it something.
3. R tries to do what you asked it to, but it can't. It gives you an **error message**, as well as a new prompt so you can try again or ask it something new.

R expressions and function calls

To “talk” with R, you need to know how to give it a complete **expression**.

Most expressions you’ll want to give R will be some combination of two elements:

1. **Function calls**
2. **Object assignments**

We’ll go through both these pieces and also look at how you can combine them together for some expressions.

Function calls

In general, function calls in R take the following structure:

```
## Generic code
function_name(formal_argument_1 = named_argument_1,
              formal_argument_2 = named_argument_2,
              [etc.])
```

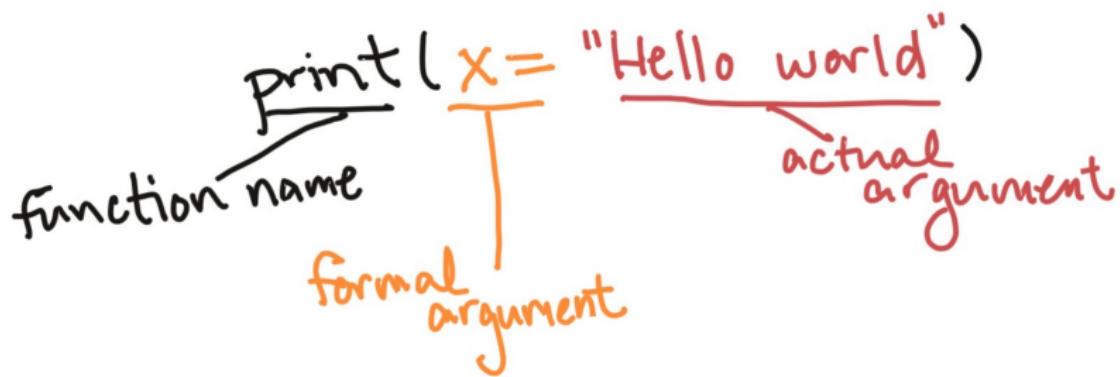
Function calls

Here is an example of this structure:

```
print(x = "Hello world")
```

```
## [1] "Hello world"
```

Function calls



In this example, we're **calling** a function with the **name** `print`. It has one **argument**, with a **formal argument** of `x`, which in this call we've provided the **named argument** "Hello world".

Function calls

The **arguments** are how you customize the call to an R function.

For example, you can use change the named argument value to print different messages with the `print` function:

```
print(x = "Hello world")
```

```
## [1] "Hello world"
```

```
print(x = "Hi Fort Collins")
```

```
## [1] "Hi Fort Collins"
```

Function calls

Some functions do not require any arguments. For example, the `getRversion` function will print out the version of R you are using.

```
getRversion()
```

```
## [1] '3.6.1'
```

Function calls

Some functions will accept multiple arguments. For example, the `print` function allows you to specify whether the output should include quotation marks, using the `quote` formal argument:

```
print(x = "Hello world", quote = TRUE)
```

```
## [1] "Hello world"
```

```
print(x = "Hello world", quote = FALSE)
```

```
## [1] Hello world
```

Function calls

Arguments can be **required** or **optional**.

For a required argument, if you don't provide a value for the argument when you call the function, R will respond with an error. For example, `x` is a **required argument** for the `print` function, so if you try to call the function without it, you'll get an error:

```
print()
```

```
Error in print.default() : argument "x" is
missing, with no default
```

Function calls

For an **optional argument** on the other hand, R knows a **default value** for that argument, so if you don't give it a value for that argument, it will just use the default value for that argument.

For example, for the print function, the quote argument has the default value TRUE. So if you don't specify a value for that argument, R will assume it should use quote = TRUE. That's why the following two calls give the same result:

```
print(x = "Hello world", print = TRUE)
```

```
## [1] "Hello world"
```

```
print(x = "Hello world")
```

```
## [1] "Hello world"
```

Function helpfiles

Often, you'll want to find out more about a function, including:

- Examples of how to use the function
- Which arguments you can include for the function
- Which arguments are required versus optional
- What the default values are for optional arguments.

You can find out all this information in the function's **helpfile**, which you can access using the function ?.

Function helpfiles

For example, the `mean` function will let you calculate the mean (average) of a group of numbers. To find out more about this function, at the console type:

```
?mean
```

This will open a helpfile in the “Help” pane in RStudio.

Function helpfiles

Helpfile for "mean"



mean (base)

R Documentation

Arithmetic Mean

Description

Generic function for the (trimmed) arithmetic mean.

Usage

mean(x, ...)

#> Default method:

mean(x, trim = 0, na.rm = FALSE, ...)

Arguments

- x An `object`. Currently there are methods for numeric/logical vectors and `date`, `date-time` and `time-interval` objects. Complex vectors are allowed for `trim = 0`, only.
- trim the fraction 0 to 0.5 of observations to be trimmed from each end of `x` before the mean is computed. Values of `trim` outside that range are taken as the nearest endpoint.
- na.rm a logical value indicating whether `NA` values should be stripped before the computation proceeds.
- ... further arguments passed to or from other methods.

Value

If `trim` is zero (the default), the arithmetic mean of the values in `x` is computed, as a numeric or complex vector of length one. If `x` is not logical (coerced to numeric), numeric (including integer) or complex, `NA_real_` is returned, with a warning.

If `trim` is non-zero, a symmetrically trimmed mean is computed with a fraction of `trim` observations deleted from each end before the mean is computed.

References

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

See Also

`weighted.mean`, `mean.POSIXct`, `colMeans` for row and column means.

Examples

```
x <- c(0.10, 50)
m1 <- mean(x)
c(m1, mean(x, trim = 0.10))
```

Required and optional arguments

Descriptions of arguments

What the function outputs

Examples of using the function

The helpfile includes sections giving the function's **usage**, **arguments**, **value**, and **examples**.

Function helpfiles

Helpfile for "mean"

↓

mean(x, trim = 0, na.rm = FALSE, ...)

Arithmetic Mean

Description

Generic function for the (trimmed) arithmetic mean.

Usage

mean(x, ...)

At Default S3 method:

mean(x, trim = 0, na.rm = FALSE, ...)

Arguments

x – An object. Currently there are methods for numeric/logical vectors and `date`, `date-time` and `time.interval` objects. Complex vectors are allowed for `mean` only.
trim – the fraction (0 to 0.5) of observations to be stripped from each end of `x` before the mean is computed. Values of `trim` outside that range are taken as 0 or against endpoint.
na.rm – a logical value indicating whether NA values should be stripped before the computation proceeds.
... – further arguments passed to or from other methods.

Value

If `trim` is zero (the default), the arithmetic mean of the values in `x` is computed, as a numeric or complex vector of length one. If `x` is not logical (converted to numeric), it must be (including integer) or complex. `NA` and `NaN` is returned, with a warning.
If `trim` is non-zero, a symmetrically trimmed mean is computed with a fraction of `trim` observations deleted from each end before the mean is computed.

References

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

See Also

`weighted.mean`, `mean.POSIXct`, `colMeans` for row and column means.

Examples

```
x <- c(0, 10, 50)
x0 <- mean(x)
c(x0, mean(x, trim = 0.1))
```

What the function outputs →

Required argument

default values

optional arguments

↑ Examples of using the function

↑ Required and optional arguments

↑ Descriptions of arguments

You can figure out which arguments are **required** and which are **optional** in the Usage section of the helpfile.

Objects, object names, and assignments

Function calls will usually produce something called an **object**.

If you just call a function, as we've been doing, then R will respond by printing out that object.

However, we'll often want to use that object some more. For example, we might want to use it as an argument later in our "conversation" with R, when we call another function later.

If you want to re-use the results of a function call later, you can **assign** that **object** to an **object name**.

R's MVP: The *gets arrow*

The *gets arrow*, `<-`, is R's assignment operator. It takes whatever you've created on the right hand side of the `<-` and saves it as an object with the name you put on the left hand side of the `<-` :

```
## Note: Generic code-- this will not work  
[object name] <- [object]
```

R's MVP: The *gets arrow*

For example, if I just type "Hello world", R will print it back to me, but won't save it anywhere for me to use later:

```
"Hello world"
```

```
## [1] "Hello world"
```

R's MVP: The *gets arrow*

However, if I assign it to an object, I can “refer” to that object in a later function call.

For example, the code below assigns the **object** "Hello world" the **object name** message. Later, I can just refer to this object using the name message, for example in a function call to the print function:

```
message <- "Hello world"  
print(x = message)  
  
## [1] "Hello world"
```

History of <-



Photo by Marcin Wichary, User:AlanM1 - Derived (cropped) from, CC BY 2.0,
<https://commons.wikimedia.org/w/index.php?curid=20744606>

R's MVP: The *gets arrow*

When you enter an **assignment expression** like this at the R console, if everything goes right, then R will “respond” by giving you a new prompt, without any kind of message.

However, there are three ways you can check to make sure that the object was assigned to the object name:

1. Enter the object’s name at the prompt and press return. The default if you do this is for R to “respond” by calling the `print` function with that object as the `x` argument.
2. Call the `ls` function (which doesn’t require any arguments). This will list all the object names that have been assigned in the current R session.
3. Look in the “Environment” pane in RStudio. This also lists all the object names that have been assigned in the current R session.

R's MVP: The *gets arrow*

Here's an example of the first two strategies:

1. Enter the object's name at the prompt and press return:

```
message
```

```
## [1] "Hello world"
```

2. Call the `ls` function:

```
ls()
```

```
## [1] "message"
```

"Environment" pane

"Environment" tab

The screenshot shows the RStudio interface with the "Environment" tab selected. The top menu bar includes Environment, History, Connections, Build, and Git. Below the menu is a toolbar with icons for file operations like Open, Save, and Import Dataset, and a "Global Environment" dropdown. A search bar and a "List" button are also present. The main pane is titled "Values" and displays a single entry: "message" with the value "'Hello world'". Handwritten annotations in red point to various parts of the interface:

- A bracket on the left points to the list of objects in the pane, labeled "list of assigned objects".
- An arrow points to the "Clear all objects" button in the toolbar, labeled "clear all objects".
- An arrow points to the "object name" part of the first row, labeled "object name".
- An arrow points to the "some of the content of the object" part of the first row, labeled "some of the content of the object".

Object names

There are some rules for the names you can use for an object name:

- Use only letters, numbers, and underscores [and dots?]
- Don't start with anything but a letter

From [Hadley Wickham's R style guide](#) [link still current?]

- Use lower case for variable names (`fedex_number`, not `FedExNumber`)
- Use an underscore as a separator (`fedex_number`, not `fedex.number` or `fedexNumber`)
- Avoid using names that are already defined in R (e.g., don't name an object `mean`, because a `mean` function exists)

Assigning objects to object names

If you try to assign an object to a name that doesn't follow the "hard" rules, you'll get an error.

For example, all of these expressions will give you an error:

```
1message <- "Hello world"  
_message <- "Hello world"  
message! <- "Hello world"
```

“Composing” to combine function calls

What if you want to “compose” a call from more than one function call?

One way to do it is to assign the output from the first function call to a name and then use that name for the next call.

For example:

```
message <- paste("Hello", "world")
print(x = message)

## [1] "Hello world"
```

“Composing” to combine function calls

You can also “nest” one function call inside another function call. For example:

```
print(x = paste("Hello", "world"))
```

```
## [1] "Hello world"
```

Just like with math, the order that the functions are evaluated moves from the inner set of parentheses to the outer one.

Using R scripts

The console can be great for quick functions to explore the data. However, for most data analysis work you'll want to use a script, so you can save all the code you used for the analysis. This improves the *reproducibility* of your analysis.

R scripts

An **R script** is a file where you can write down and save R code.

When you write, run, and save your R code in a script rather than running it one line at a time in the console, you can easily go back and re-do exactly what you did again later. You can also share the script for someone else to use, or run it on a different computer.

R scripts

RStudio has one pane that shows any R scripts you have open. If you'd like to create new R scripts, you can do that in RStudio with the following steps:

- Open a new script file in RStudio: File -> New File -> R Script.
- I recommend that you make an "R" folder in all of the R projects that you create and save all your script files in that folder.
- Save scripts using the extension .R

R scripts

Running code in R scripts line-by-line:

- To run code from an R script file in RStudio, you can use the Run button (or Command-R).
- This will run whatever's on your cursor line or whatever's highlighted.

Sourcing an R script (i.e., running all the code saved in the script):

- To run the whole script, you can also use the source function with the filename.
- You can also use the “Source” button on the script pane.

R scripts

The screenshot shows the RStudio interface with several tabs at the top: "rcise.Rmd", "InCourseExercises_Week2.Rmd", "Chapter2.Rmd", and "Untitled28*". Below the tabs is a toolbar with icons for Save, Run, and Source. The main area displays R code:

```
1 ## Some example code to show a script file
2
3 one_to_ten <- 1:10
4
5 course_dates <- data.frame(session = c(1, 2, 3),
6                             topic = c("Basic R",
7                                   "Getting and Cleaning Data 1",
8                                   "Exploring Data 1"))
9
10 a <- 1:4 ; b <- rnorm(10)
11 |
```

Annotations with arrows point to specific parts of the code:

- A red arrow points to the "Save" button in the toolbar, labeled "Save" button.
- A red arrow points to the "Run" button in the toolbar, labeled "Run" button.
- A red arrow points to the "Source" button in the toolbar, labeled "Source" button.
- An arrow points from the text "Commented code line" to the first three blank lines of the code (lines 1-3).
- An arrow points from the text "One command across several lines" to the assignment statement on line 10.
- An arrow points from the text "Two commands on one line" to the semicolon-separated statements on line 10.

Comment characters

Sometimes, you'll want to include notes in your code. You can do this in all programming languages by using a **comment character** to start the line with your comment.

In R, the comment character is the hash symbol, #. R will skip any line that starts with # in a script.

```
# Don't print this.
```

```
"But print this"
```

```
## [1] "But print this"
```

Closing an R session

Do **not** save the history of your R session when you close RStudio.
Instead, get in the habit of writing your R code in reproducible formarts (R scripts, RMarkdown documents)

In-Course Exercises

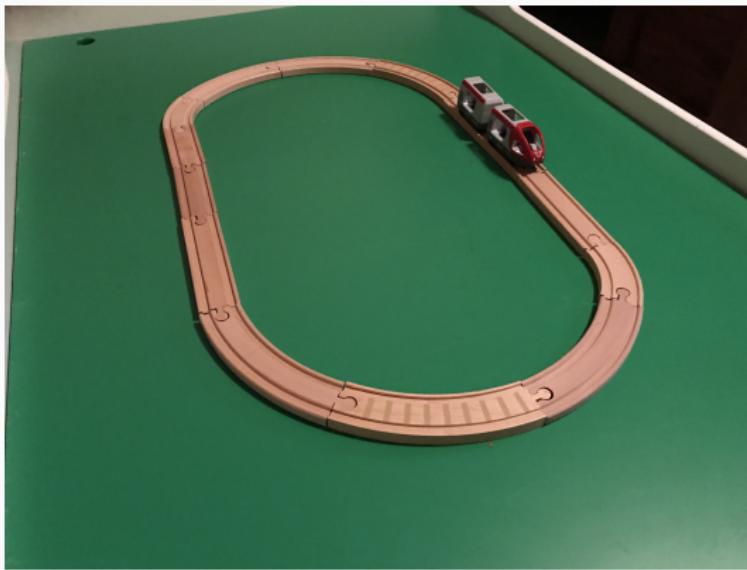
We'll take a break now to finish this week's in-course exercise. You can find all the instructions on the online coursebook (<https://geanders.github.io/RProgrammingForResearch/r-preliminaries.html>).

Get together with your group members and do Sections 1.8.4 through 1.8.8 of the In-Course Exercise for Week #1 ("R Preliminaries").

The “package” system

R packages

Your original download of R is only a starting point:



R packages

To take full advantage of R, you'll want to add on packages:



R packages

You can get packages to add-on to your version of R from:

- CRAN
- Bioconductor (specifically for bioinformatics-related packages)
- GitHub
- Your friends and collaborators
- Make them yourself

The Comprehensive R Archive Network (CRAN) is the primary source of R packages, with thousands of available packages. Each of the packages available on CRAN has a unique name and has passed some broad checks.



Dirk Eddelbuettel @eddelbuettel · 27 Jan 2017

Big congratulations to @gbwanderson whose new package 'hurricaneexposure' just became package 10,000 on CRAN !!

CRAN Package Updates @CRANberriesFeed

9999 packages on CRAN right now, so imagine dozens of R nerds hanging in suspense waiting for the package to make it 10k ...

2

35

93



Installing from CRAN

The main way you will install new packages will be with the `install.packages` function. By default, this installs packages from CRAN.

For example, to get the package “phonenumbers”, you could use:

```
install.packages("phonenumbers")
```



Installing from CRAN

Installing a package downloads it (usually from CRAN) to your computer.

You will need to be online to install a package, because it downloads the package from an online repository.

Once you install a package, you do not need to install it again. It's on your computer. (At least, until the package maintainer creates a new version.)

What's in a package?

How to use a package you've installed

1. Open an R session
2. Load the package
3. Use the package's vignette and helpfiles to figure out how to use the package (optional)
4. Use functions in the package

Loading an installed package

Once you have a package, you can load it to an R session using the `library()` function.

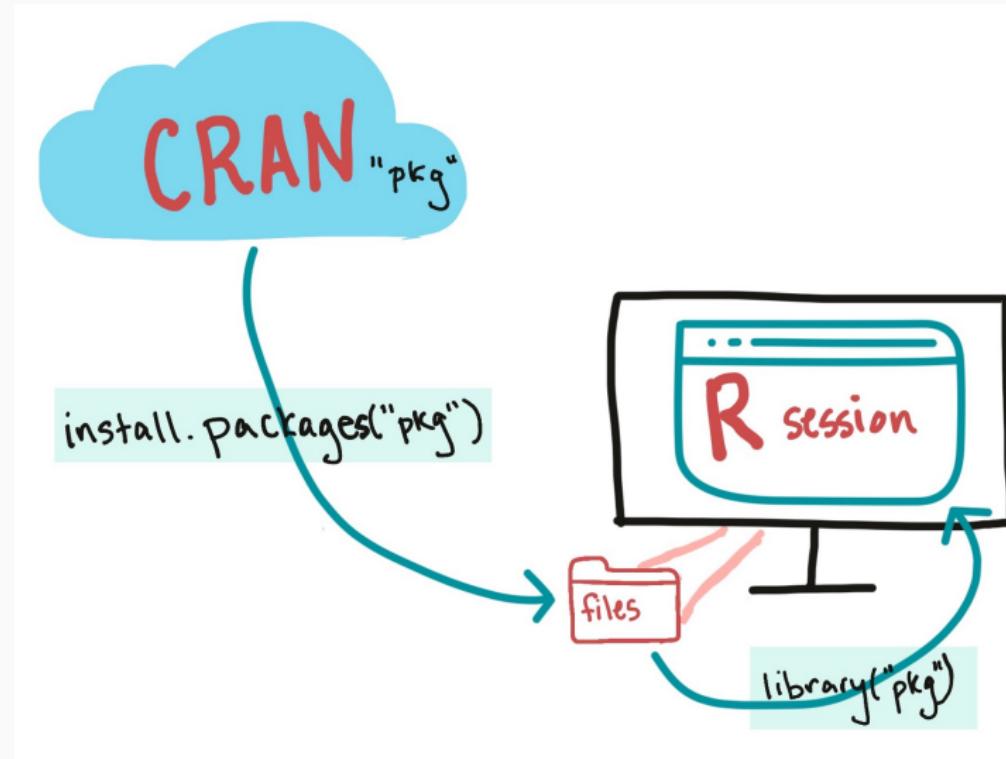
```
library("phonenumbers")
```

Once it's loaded, you can use all its functions.

```
message <- "HelloWorld"  
letterToNumber(message)
```

```
## [1] "4355696753"
```

Installing versus loading



Install a package (with `install.packages`) to get it onto your computer. **Load** it (with `library`) to get it into your R session.

Package vignettes

Many packages will come with a “vignette”, or a tutorial on how to use the package. These are very helpful tools for figuring out how to use a package.

To get a list of all the vignettes a package has, use the `vignette` function, specifying a package name with `package`:

```
vignette(package = "phonenumbers")
```

Package vignettes

To open a vignette once you know its name, you can also use the vignette function. For example, to open the vignette called "phonenumbers" for the phonenumbers package, run:

```
vignette("phonenumbers", package = "phonenumbers")
```

Many packages only have one vignette, with the same name as the package. In that case, you can open the vignette using a shorter call:

```
vignette("phonenumbers")
```

Function helpfiles

Packages include helpfiles for all the functions they intend for people to use.

You can access the helpfile for a function using ? followed by the function name. For example, to get the helpfile for letterToNumber, you can run:
?letterToNumber

Object classes

What are object classes?

...

One class of objects can be a subtype of a more general object class. We'll see an example of this today.

Today, we'll look at two key object classes for working with data:

1. Vectors
2. Dataframes

What are object classes?

For these two object classes (vectors and dataframes), we'll look at:

1. How to make a new object with that class
2. How to extract values from objects with that class

In later classes, we'll spend a lot of time learning how to do other things with objects from these two classes.

Vectors

Dataframes

Creating vector objects

Creating dataframe objects

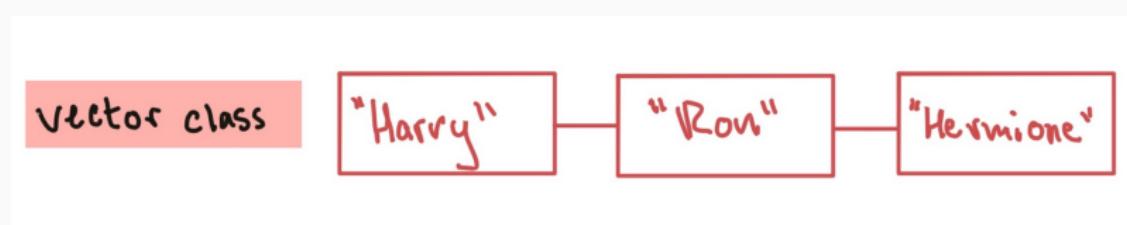
Extracting values from a vector

Extracting values from a dataframe

R's most basic object classes

The two most common types of objects for data in R that we will use in this class are **vectors** (1D) and **dataframes** (2D).

Vectors



A **vector** is an object class where the object is made of a string of values.
All those values must be of the same **class** (i.e., all numbers, all
characters, all dates)

Vectors

To create a vector, you can use the concatenate function, c:

```
c("Harry", "Ron", "Hermione")
```

```
## [1] "Harry"     "Ron"       "Hermione"
```

Vectors

If you want to use that object later, you can assign it to a name:

```
main_characters <- c("Harry", "Ron", "Hermione")  
print(x = main_characters)
```

```
## [1] "Harry"      "Ron"        "Hermione"
```

Vectors

A handwritten diagram illustrating the structure of the assignment expression `main_characters ← c("Harry", "Ron", "Hermione")`. The diagram uses arrows and labels to identify parts of the code:

- assignment operator**: Points to the `←` symbol.
- main_characters**: Points to the variable name `main_characters`.
- object / name**: Points to the variable name `main_characters`.
- c**: Points to the function name `c`.
- function calls**: Points to the function call `c("Harry", "Ron", "Hermione")`.
- arguments**: Points to the arguments `"Harry", "Ron", "Hermione"`.

This **assignment expression** follows the structure we covered earlier for function calls and expressions.

Vectors

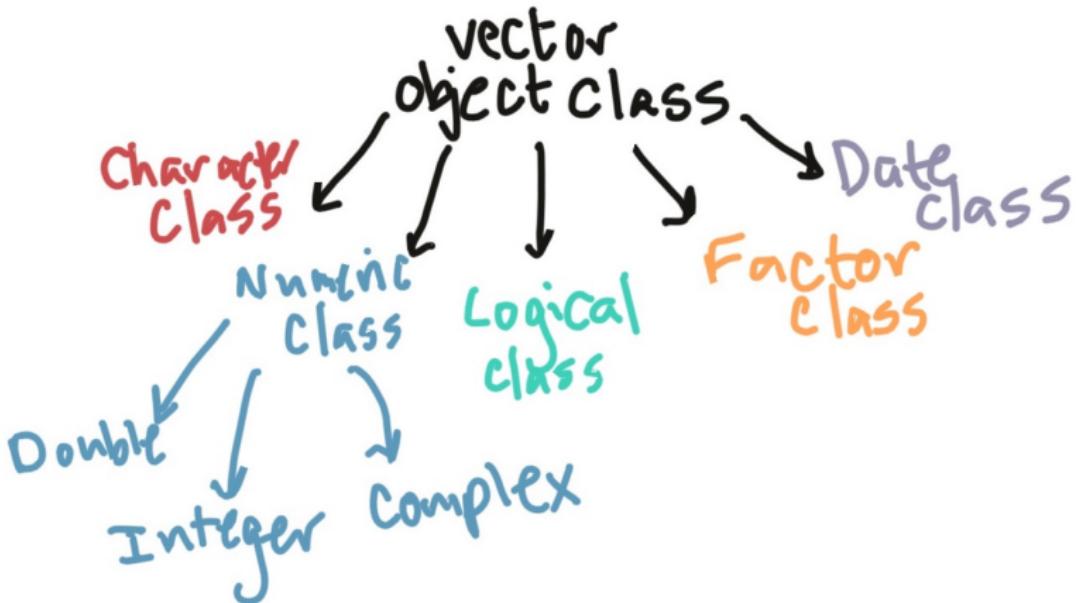
The *length* of the vector is how many values it has in it. The `main_characters` vector includes three values (“Harry”, “Ron”, and “Hermione”), so it has a length of 3.

You can use the `length` function to figure out how long a vector is:

```
length(x = main_characters)
```

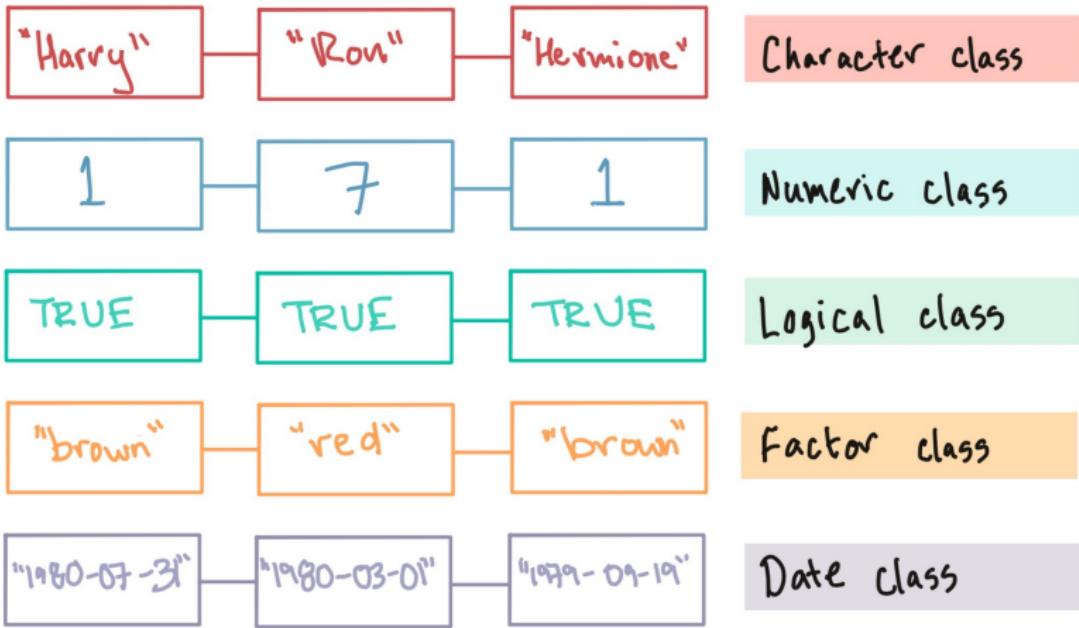
```
## [1] 3
```

Vectors



Vectors can have several different **classes**, including numeric, character, factor, date, and logical. Some of these classes are further divided into more specific classes.

Vectors



Vectors

Typically, when you use the `c` function, you'll be creating a numeric, character, or logical vector. More complex classes (like factors and dates) require a bit more work to create "from scratch".

- For character vectors, use quotation marks around each element.

```
main_characters <- c("Harry", "Ron", "Hermione")
```

- For numeric, don't use quotation marks.

```
n_kids <- c(1, 7, 1)
```

Vectors

You can use the `class` function to figure out the class of a vector.

```
class(main_characters)
```

```
## [1] "character"
```

```
class(n_kids)
```

```
## [1] "numeric"
```

Vectors

If you create a vector with a mix of classes, R will default the whole vector to the most generic class:

```
mixed_classes <- c(1, 3, "five")  
mixed_classes
```

```
## [1] "1"     "3"     "five"
```

[Order of generic]

Extracting values from vectors

[Tidyverse alternative to indexing?]

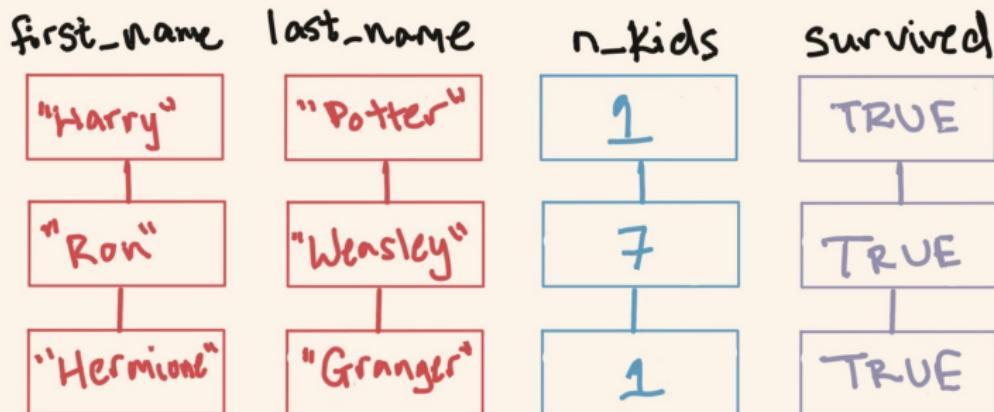
Vectors

You can pull out certain values by using indexing (`[...]`) to identify the locations you want to get:

```
main_characters[2] # Get the second value  
  
## [1] "Ron"  
  
main_characters[c(1, 3)] # Get first and third values  
  
## [1] "Harry"      "Hermione"  
  
main_characters[1:3] # Get the first three values  
  
## [1] "Harry"      "Ron"        "Hermione"
```

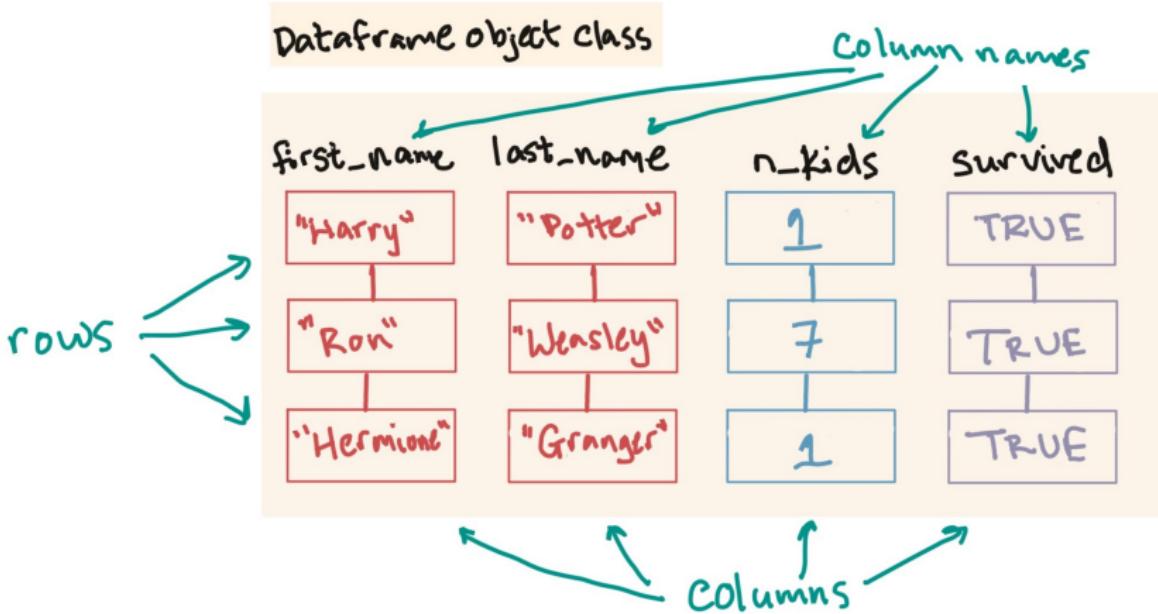
Dataframes

Dataframe object class



A **dataframe** combines one or more vectors of the same length stuck together side-by-side.

Dataframes



A dataframe contains **rows** and **columns**, and each column has a **column name**.

Dataframes

You can create dataframes using the `tibble` function from the `tibble` package. We'll look at how to do that.

However, most often you will create a dataframe by reading in data from a file using something like `read_csv` from the `readr` package—most datasets will not be short enough that you want to enter them by hand in R.

Dataframes

The format for using tibble is:

```
## Generic code
library("tibble")
[name of object] <- tibble([1st column name] =
                           [1st column content],
                           [2nd column name] =
                           [2nd column content],
                           etc.)
```

Dataframes

```
library("tibble")
hp_data <- tibble(first_name = c("Harry", "Ron", "Hermione"),
                   last_name = c("Potter", "Weasley", "Granger"),
                   n_kids = c(1, 7, 1),
                   survived = c(TRUE, TRUE, TRUE))
hp_data

## # A tibble: 3 x 4
##   first_name last_name n_kids survived
##   <chr>      <chr>     <dbl>  <lgl>
## 1 Harry       Potter      1  TRUE
## 2 Ron         Weasley     7  TRUE
## 3 Hermione    Granger     1  TRUE
```

Dataframes

You can also create dataframes by joining together vector objects you previously created, as long as they have the same length and line up:

```
hp_data <- tibble(first_name = main_characters,  
                  last_name = c("Potter", "Weasley", "Granger"),  
                  n_kids = n_kids,  
                  survived = c(TRUE, TRUE, TRUE))
```

```
hp_data
```

```
## # A tibble: 3 x 4  
##   first_name last_name  n_kids  survived  
##   <chr>      <chr>     <dbl>   <lgl>  
## 1 Harry      Potter      1 TRUE  
## 2 Ron        Weasley     7 TRUE  
## 3 Hermione   Granger     1 TRUE
```

Extracting values from dataframes

The `dplyr` package has two functions for extracting data from dataframes by position: `slice` to extract rows based on their row position and `select` to extract columns based on their column position.

Extracting values from dataframes

For example, if you wanted to get the first two rows of the `hp_data` dataframe, you could run:

```
library("dplyr")

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
## 
##     filter, lag

## The following objects are masked from 'package:base':
## 
##     intersect, setdiff, setequal, union

slice(.data = hp_data, c(1:2))

## # A tibble: 2 x 4
##   first_name last_name n_kids survived
```

Extracting values from dataframes

If you wanted to get the first and fourth columns, you could run:

```
select(.data = hp_data, c(1, 4))
```

```
## # A tibble: 3 x 2
##   first_name survived
##   <chr>      <lgl>
## 1 Harry      TRUE
## 2 Ron        TRUE
## 3 Hermione   TRUE
```

Extracting values from dataframes

You can compose calls from both functions. For example, you could extract the values in the first and fourth columns of the first two rows with:

```
select(slice(.data = hp_data, c(1:2)), c(1, 4))
```

```
## # A tibble: 2 x 2
##   first_name survived
##   <chr>     <lgl>
## 1 Harry     TRUE
## 2 Ron       TRUE
```

Dataframes

Usually, instead of creating a dataframe from vectors, you'll read one in from data on an outside file, for example using `read_csv` from the `readr` package.

For example, to read in a dataset from a csv file called "daily_show_guests.csv":

```
library("readr")
daily_show <- read_csv("daily_show_guests.csv",
                      skip = 4)
```

Dataframes

This has read data in from the external file into a dataframe object in my R session:

```
daily_show[1:4, 1:4]
```

```
## # A tibble: 4 x 4
##   YEAR GoogleKnowlege_Occupation Show      Group
##   <dbl> <chr>                  <chr>    <chr>
## 1 1999 actor                  1/11/99  Acting
## 2 1999 Comedian               1/12/99  Comedy
## 3 1999 television actress   1/13/99  Acting
## 4 1999 film actress          1/14/99  Acting
```

Dataframes

You can use the functions `dim`, `nrow`, and `ncol` to figure out the dimensions (number of rows and columns) of a dataframe:

```
dim(daily_show)
```

```
## [1] 2693      5
```

```
nrow(daily_show)
```

```
## [1] 2693
```

```
ncol(daily_show)
```

```
## [1] 5
```

Dataframes

Base R also has some useful functions for quickly exploring dataframes:

- `str`: Show the structure of an R object, including a dataframe
- `summary`: Give summaries of each column of a dataframe.

For example, to get summaries of the first two columns of the Daily Show data, you can run:

```
summary(daily_show[ , c(1, 2)])
```

```
##          YEAR      GoogleKnowlege_Occupation
##  Min.    :1999   Length:2693
##  1st Qu.:2003   Class  :character
##  Median  :2007   Mode   :character
##  Mean    :2007
##  3rd Qu.:2011
##  Max.    :2015
```

In-Course Exercises

Here are a few tips for doing these in-course exercises:

- Sometimes, I will ask you to try out things that we haven't covered yet in class. In this case, I'll give you detailed instructions. You can also use Google, R help files, and any other resources you want to help you figure out how to do something for an in-course exercises. If you're really stumped, I'm happy to help.
- For each section of the exercise, I've included some example R code to show one possible solution. Try to not look at this until you've tried the exercise yourself without it. If you move quickly to this example R code, you will not learn R very quickly.
- Some groups tend to move more quickly than others. I often include "If you have extra time" sections for groups who are moving more quickly.
- Use Google!

In-Course Exercises

The #1 rule for this class:

Don't use spaces in file pathnames!

This includes the names of files and the names of directories. Keep this in mind during today's exercise. (We'll talk more about why not next week.)

In-Course Exercises

We'll take a break now to get started on this week's in-course exercise. You can find all the instructions on the online coursebook (<https://geanders.github.io/RProgrammingForResearch/r-preliminaries.html>).

Get together with your group members and do Sections 1.8.1 through 1.8.4 of the In-Course Exercise for Week #1 ("R Preliminaries").

PROBABLY

As of Wednesday, November 7th, 2012, Nate Silver is probably a witch.

His unusually accurate predictions are, thus far, explained by his use of validated statistical methods. His disregard of momentum, gut feelings, and the interpretations of people paid to promote certain viewpoints is not the result of supernatural assistance.

While we on the Is Nate Silver a Witch editorial board are strict rationalists, Mr Silver's performance has been uncanny enough to raise small but significant doubts as to whether his methodology is entirely of this world. We are following the situation closely.

Brought to you by epistemology and @vruba.



For the first time in our history, the winners of the White House Turkey Pardon were chosen through a highly competitive online vote. And once again, Nate Silver completely nailed it. The guy is amazing.

— Barack Obama —

AZ QUOTES

For the story, see <http://www.politico.com/blogs/media/2012/11/obama-cracks-nate-silver-joke-at-turkey-pardoning-150132>.

Some odds and ends

Missing values

In R, NA is used to represent some types of missing values in a vector. This value can show up in numerical or character vectors (or in vectors of some other classes):

```
c(1, 4, NA)
```

```
## [1] 1 4 NA
```

```
c("Jane Doe", NA)
```

```
## [1] "Jane Doe" NA
```

Missing values

Geeky license plate earns hacker \$12,000 in parking tickets

A California man's vanity license plate backfires spectacularly.

JONATHAN M. GITLIN - 8/13/2019, 11:25 AM



[Enlarge](#) / This can be a bad idea, evidently.

For the full story: <https://arstechnica.com/cars/2019/08/wiseguy-changes-license-plate-to-null-gets-12k-in-parking-tickets/>

The \$ operator

We've talked about how you can combine vectors of the same length to create a dataframe.

To go the other direction (pull a column from a dataframe), you can use the \$ operator.

For example, say you have the following dataset and want to pull the color column as a vector:

```
example_df <- data.frame(color = c("red", "blue"),
                           value = c(1, 2))

example_df
```

```
##   color value
## 1   red     1
## 2  blue     2
```

The \$ operator

You can pull the color column as a vector using the name of the dataframe, the dollar sign, and then the name of the column:

```
example_df$color
```

```
## [1] red blue
## Levels: blue red
class(example_df$color)
```

```
## [1] "factor"
```

(Note: You can use tab completion in RStudio after you put in example_df\$.)

paste and paste0

If you want to paste together several character strings to make a length-one character vector, you can use the `paste` function to do that:

```
paste("abra", "ca", "dabra")
```

```
## [1] "abra ca dabra"
```

By default, spaces are used to separate each original character string in the final string.

paste and paste0

If you want to remove these spaces, you can use the `sep` argument in the `paste` function:

```
paste("abra", "ca", "dabra", sep = "")
```

```
## [1] "abracadabra"
```

A short-cut function is `paste0`, which is identical to running `paste` with the argument `sep = ""`:

```
paste0("abra", "ca", "dabra")
```

```
## [1] "abracadabra"
```

<- vs. =

You can make assignments using either `<-` or `=`, and you'll see both when you're reading other people's code.

However, R gurus advise using `<-` in your own code.

The arrow shows the *direction* of assignment. Also, as you move to doing more complex things, problems might crop up if you use `=` for assignment in a few specific cases.

<- vs. =

For now, though, it will be helpful for you to know that these two calls do the same thing:

```
one_to_ten <- 1:10
```

```
one_to_ten
```

```
## [1] 1 2 3 4 5 6 7 8 9 10
```

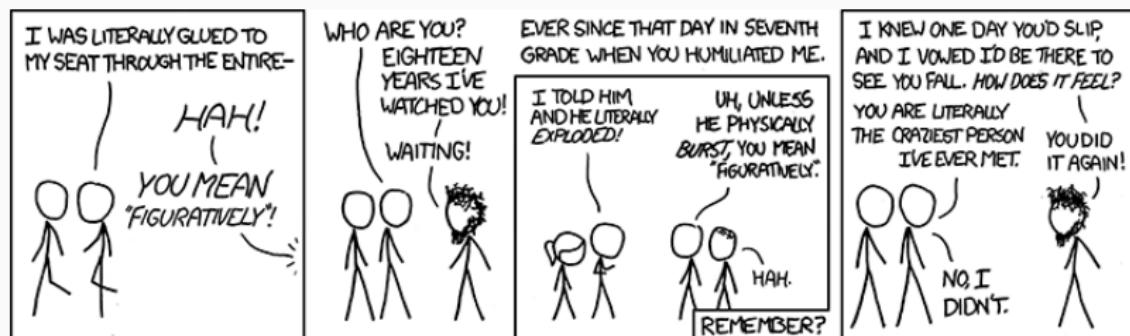
```
one_to_ten = 1:10
```

```
one_to_ten
```

```
## [1] 1 2 3 4 5 6 7 8 9 10
```

When to use quotations

- Use quotation marks if you **literally** mean that specific character string
- Don't use quotation marks if you want to reference an R object by its name
- Never use quotation marks on the left-hand side of <- or =.



Source: xkcd.com

When to use quotations

- Use quotation marks if you **literally** mean that specific character string

```
c("Harry", "Ron")
```

```
## [1] "Harry" "Ron"
```

- Never use quotation marks on the left-hand side of <- or =.

```
Harry <- c("Good at Quidditch")
```

```
Ron <- c("Good at Wizard Chess")
```

- Don't use quotation marks if you want to reference an R object by its name

```
c(Harry, Ron)
```

```
## [1] "Good at Quidditch"      "Good at Wizard Chess"
```

Alternative package::function notation

The library function is the most common way you'll access functions from R packages, but it's not the only way to do this. There is another type of notation that will allow you to use any external function from any package you have installed on your computer: the package::function notation.

You can call functions by specifying the package name, a double colon, and then the function name you want to use from that package. For example:

```
phonenumbers::letterToNumber(message)
```

```
## [1] "4355696753"
```

Alternative package::function notation

Search path

- base R
- phonenumbers package
- decoder ring package

Alternative package::function notation

The package::function notation is not typically used because it substantially increases how much you have to type in your code.

However, there can be cases where a function name is ambiguous. For example, you might want to use functions from two different packages that have the same name. In this case, using the package::function notation makes it crystal clear which function you mean.

In practice, this problem is most likely to come up when you've loaded both `plyr` and `dplyr`, which share several function names and are both popular packages.

Alternative package::function notation

There is another useful trick that you can do with the package::function notation.

RStudio has tab complemention, which means that once you start typing an object name or function, if you press the Tab key, RStudio will give you suggestions for how you can finish the name.

If you want to scroll through the names of all the external functions in a package, you can do so by typing something like ?phonenumbers:: in the console (don't press Return) and then pressing the Tab key.

“Composing” to combine function calls

There's one additional way to compose several function calls into one R expression.

There's a special “pipe” function (`%>%`) in a package called “magrittr” that you can use to compose a large R expression from several smaller function calls. This will “pipe” the object that is created from one function call into the first argument for the next function call.

We'll explore this a **lot** more later in the course, but here's a taste of how that looks:

```
library(magrittr)

paste("Hello world") %>%
  print()

## [1] "Hello world"
```

“Composing” to combine function calls

Here's another example. Earlier, we used the following code to extract certain rows and columns from a dataframe:

```
select(slice(hp_data, c(1:2)), c(1, 4))
```

```
## # A tibble: 2 x 2
##   first_name survived
##   <chr>     <lgl>
## 1 Harry     TRUE
## 2 Ron       TRUE
```

You can use a pipe to express this much more cleanly:

```
slice(.data = hp_data, c(1:2)) %>%
  select(c(1, 4))
```

```
## # A tibble: 2 x 2
##   first_name survived
##   <chr>     <lgl>
## 1 Harry     TRUE
```