# Introduction to Statistical Programming with R

Aileen Brown Cuevas
Thomas Wiemann
*University of Chicago*

July 29, 2021

## Outline

1. Statistical Programming

2. Why use R

3. Getting started
   - ▶ Base R
   - ▶ R Studio
   - ▶ Jupyter Notebooks

4. Version Control
   - ▶ Overview
   - ▶ Git with GitHub
   - ▶ Sample Problem Set Repository
   - ▶ Best Practices
   - ▶ Optional Homework

## Statistical Programming

Statistical programming may be defined as the implementation of statistical procedures and statistical analyses on a computer.

To do so efficiently, the choice of the right tool (i.e., programming language) for the problem at hand is crucial. Four distinctions are particularly relevant:

▶ high-level languages (more human readable, less control) vs low-level languages (less human readable, more control)

▶ domain languages (optimized for specific tasks) vs general-purpose languages (suitable for large variety of tasks)

▶ availability of packages

▶ open-source/free vs proprietary/commercial

For statistical programming, this leaves us with C/C++ for low-level high-performance code, and R, Python, and Julia for higher level code.

# Why use R

Pros:

- ▶ It's free, open source, and available on every major platform.
- ▶ Widely used and sought-after language (e.g., see IEEE Spectrum)
- ▶ Large and incredibly helpful community (e.g., see Chadley Wickham)
- ▶ Excellent availability of data manipulation and visualization packages (e.g., `dplyr`, `ggplot2`).
- ▶ Possible to connect to low-level programming languages such as C++ (e.g., `Rcpp`).
- ▶ *Thousands* of packages designed and maintained by statisticians. Lot's of cutting-edge methods implemented in R.

Cons:

- ▶ Thousands of packages designed and maintained by *statisticians*. We're not professional programmers – code often messy.
- ▶ Easy to write, difficult to perfect. Naive code can turn out slow.
- ▶ Can be very memory intensive.
- ▶ Advanced numerical optimization methods largely missing.

# Programming Pluralism



**Victor Chernozhukov**
@VC31415

Recommendations for future students in econometrics and statistics. For programming, invest in both Python and R. For math, probability or stochastic processes, or both (I'd recommend engineering dept over math dept courses; with the exception: math.uci.edu/~rvershyn /pape....)

11:35 PM · Mar 20, 2021 · Twitter Web App

**104** Retweets  **7** Quote Tweets  **519** Likes

**Julius P** @EconJulius · Mar 21
Replying to @VC31415
Why not stata instead of R? Is this solely a preference thing?

2

1

**Victor Chernozhukov** @VC31415 · Mar 21
Stata is a respectful, high quality professional grade software for empirical work using highly standardized methods. I don't personally use it, but it is widely used.

3

*Source:* `https://twitter.com/VC31415/status/1373402880280047617`

Takeaway: use the most convenient tool for the project at hand.

# Getting started: Base R

Download & install R: https://cran.rstudio.com/



Rough! Better use a integrated development environment (IDE).

## Getting started: IDEs

Most popular IDE for R is R Studio:

- ▶ Designed specifically for R.

- ▶ Project management and integrated version control.

- ▶ Debugging support.
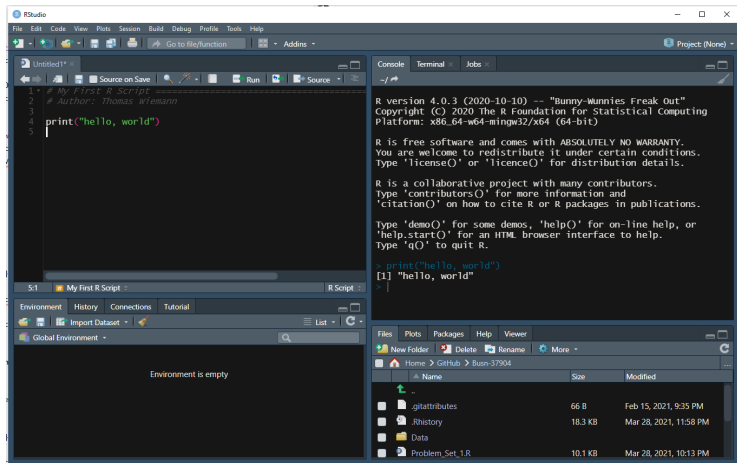
- ▶ Integrated R Markdown.

A sometimes convenient alternative to R Markdown is a Jupyter Notebook, especially if you're familiar with Jupyter Notebooks through working with Python or Julia. (Note: Ju-Py-[t]-R)

Personally, I do 99% of R coding in R Studio and 1% in Jupyter Notebooks.

# Getting started: R Studio

Download & install R Studio:
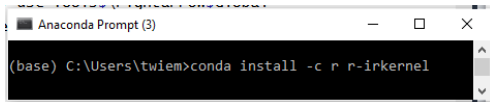`https://rstudio.com/products/rstudio/download`



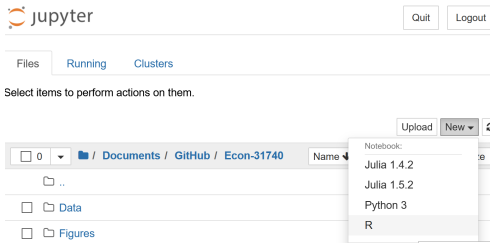To change visual appearance, use Tools→Global Options→Appearance and change the editor theme (here: *Chaos*).

# Getting started: Jupyter Notebooks for R

- Download & install Anaconda:
  `https://www.anaconda.com/products/individual`
- Install R kernel via Anaconda prompt



- Open Jupyter Notebook and select the R kernel

# Getting started: Jupyter Notebooks for R (Contd.)

Jupyter Notebooks are a browser-based application and so will open in Chrome, Firefox, etc.



They are particularly useful for sharing output scripts that include brief sections of code and lots of explanation.

## Version Control: Overview

A key feature of any editor is the $\mathrm{UNDO}$ operator (e.g., `ctrl + z`). Intuitively, a version control system is a sophisticated $\mathrm{UNDO}$ operator.
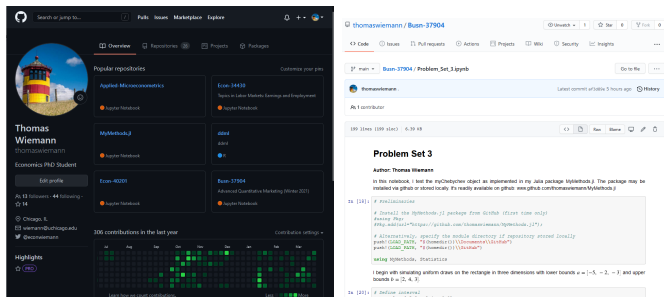
Working with version control allows you to

▶ track changes in your code;

▶ compare different versions of your code;

▶ rebuild your code as it existed at any prior point.

Invaluable for debugging and collaboration!

A particularly popular version control system is Git. It's free, available on all major operating systems, and widely supported by free hosting websites (e.g., GitHub). R Studio also has a Git integration.

## Version Control: Git with GitHub

GitHub is a free hosting website for code. Excellent for sharing code (like problem sets). Lots of software packages are hosted on GitHub.
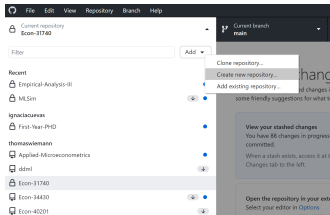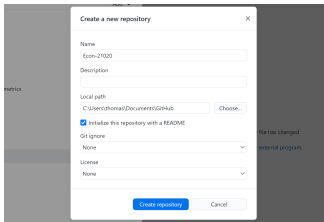


Installation:

1. Sign-up for GitHub: `www.github.com`
2. Download the GitHub desktop app: `https://desktop.github.com`
3. After installation, open GitHub app and sign-in with your account.

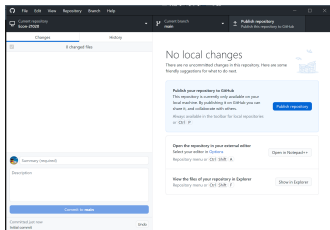We're now ready to start our first repository!
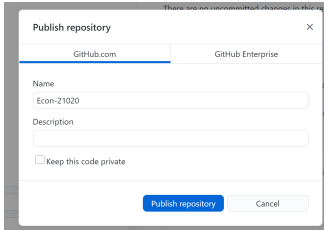
# Sample problem set repository



(a) File→Add→Create new...
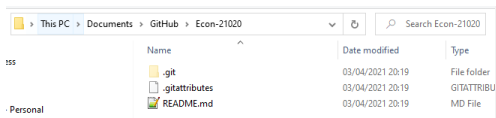
(b) Name appropriately

(c) "publish" repository

(d) Decide whether to keep private
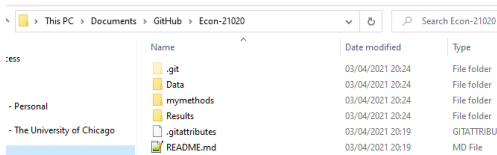
## Sample problem set repository (Contd.)

After initialization, the repository contains only the bare essentials.



Before adding code, it's useful to organize the repository for a clean workflow. For a problem set repository, it's often convenient to have separate folders for data, results, and complementary code (here: `mymethods`).
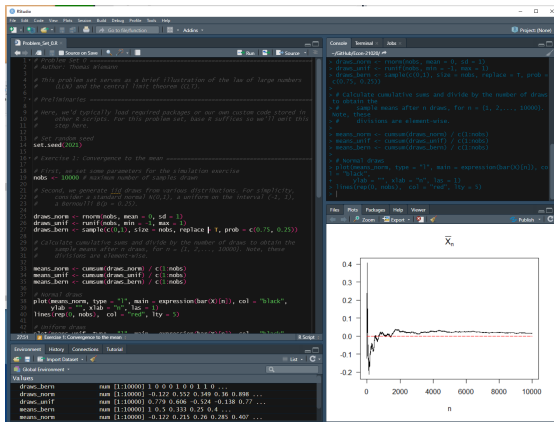
## Sample problem set repository (Contd.)

Let's now add our first R script. As an example for a problem set, Problem_Set_0.R provides a basic simulation exercises for the law of large numbers and the central limit theorem.



The file should be stored directly in the repository (i.e., it's filepath should look something like: "~/GitHub/Econ-21020/Problem_Set_0.R").

# Sample problem set repository (Contd.)



(a) Add brief description and commit

(b) "Push" to upload to GitHub

(c) Repository immediately updated...

(d) ... but all versions backed-up

## Sample problem set repository (Contd.)

Particularly useful: line-by-line comparisons between different versions of the repository.

## Best Practices

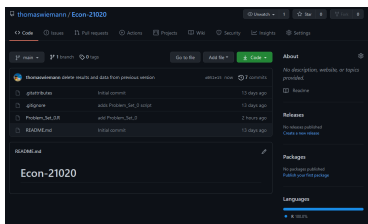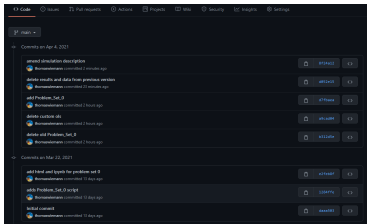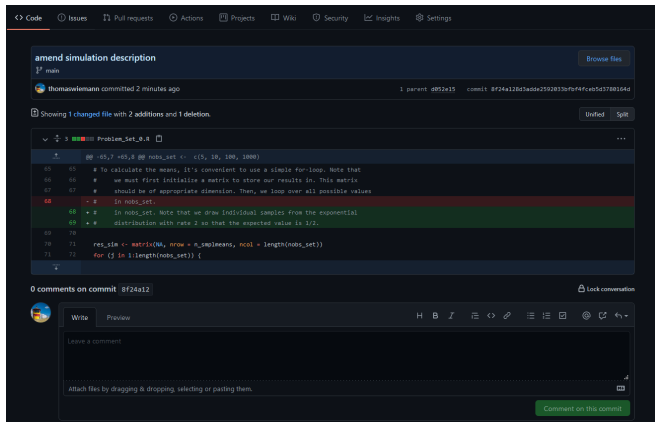Getting used to version control and GitHub can take some time. A few best practices may help you get started. (Note that these are not universal guidelines and depend on the project, but it's a good set of rules for homework repositories.)

- ▶ Commit recent changes after completing a logical step (e.g., finishing an exercise). This should be done relatively frequently.
- ▶ Always add a brief description to your commit.
- ▶ Upload to GitHub each time you stop working on your project.
- ▶ Avoid uploading large files. Small data sets are fine, large data sets may need to be stored locally.
- ▶ When collaborating, "push" and "pull" frequently.

As a student at UChicago, you qualify for GitHub pro. Comes with a few benefits, including unlimited collaborators on private repositories. Link: https://education.github.com/benefits?type=student.

## Optional Homework

Exercise: setup a problem set repository for this course.

To do so, complete the following steps:

1. Download and install R.
2. Download and install R Studio.
3. Create a GitHub account.
4. Download and install the GitHub desktop application.
5. Create a new **public** repository with the name Econ-21020 and publish it to GitHub.
6. Add a single R script called hello_world.R. The file should contain the following code: print('hello, world!').
7. Commit your changes to the repository and 'push' to GitHub.
8. Copy the link of the repository and include it in your next problem set. The link should look something like this (but of course with your username): www.github.com/thomaswiemann/Econ-21020.

All necessary instructions can be found in this slide deck.

# Introduction to Statistical Programming with R (Part. II)

AILEEN BROWN CUEVAS
THOMAS WIEMANN
*University of Chicago*

July 29, 2021

# Outline

## Data Types

The simplest kind of variables in R are *atomic types*. These do not have an internal structure. Examples are:

- ▶ integer: 1L, 2L, 42L, ...
- ▶ numeric: pi, 0.25, Inf, ...
- ▶ logical: TRUE, FALSE.
- ▶ character: "lorem ipsum", "a", "", ...

More complex data types are combinations of atomic types. It's useful to organize the most common data types based on their dimensionality and whether they are homogeneous (i.e., whether all of it's contents must be of the same type).

|     | Homogeneous   | Heterogeneous |
| --- | ------------- | ------------- |
| 1d  | Atomic vector | List          |
| 2d  | Matrix        | Data frame    |
| nd  | Array         |               |

## Data Types: Atomic vectors

Atomic vectors are usually created with c(), short for combine. Even when nesting, atomic vectors are always flat.

R does not have any scalar variables. Single numbers or characters are vectors of length 1. Vectors of length 0 also exist (and are a common cause for bugs).

Example: Atomic vectors

```
c(1L, c(2L, c(3L, 4L))) # integer
#> [1] 1 2 3 4

c(1, 2, 3, 4) # numeric
#> [1] 1 2 3 4

c(T, F) # logical
#> [1] TRUE FALSE

c("these are", "some strings") # character
#> [1] "these are" "some strings"
```

## Data Types: Factors

Another fundamental data type that is essentially a vector with attributes is `factor`. These are used to store categorical and ordinal variables.

Factor variables can easily be initialized from an existing vector by using `factor()`.

Example: Factors

```
factor(c("red", "blue", "green", "red", "green"))
#> [1] red blue green red green
#> Levels: blue green red

factor(c(1, 2, 1, 2, 3))
#> [1] 1 2 1 2 3
#> Levels: 1 2 3
```

## Data Types: Matrices

A two-dimensional collection of atomic vectors is a `matrix`. They are convenient for storing and manipulating data with linear algebra.

Matrices are easily initialized with `matrix()`. Note that R uses column-wise initialization of matrices. Vectors (or matrices) of appropriate dimension may also be combined via row- or column-concatenation using `rbind()` or `cbind()`.

Example: Matrices

```
matrix(c(1, 2, 3, 4), ncol = 2, nrow = 2)
#>      [, 1] [, 2]
#> [1, ]    1    3
#> [2, ]    2    4

cbind(c(1, 2), c(3, 4))
#>      [, 1] [, 2]
#> [1, ]    1    3
#> [2, ]    2    4
```

## Data Types: Arrays

The n-dimensional generalization of `matrix` is called `array`. They are much less common but it's good to be aware of them.

Use `array()` for initialization.

Example: Arrays

```
array(c(1:8), dim = c(2, 2, 2))
#> , , 1
#>      [, 1] [, 2]
#> [1, ]    1    3
#> [2, ]    2    4
#>
#> , , 2
#>      [, 1] [, 2]
#> [1, ]    5    7
#> [2, ]    6    8
```

## Data Types: Lists

A `list` is a 1-dimensional data type with elements of arbitrary type.
Note that this is their key difference to atomic vectors.

Use `list()` for initialization (instead of `c()`).

Example: Lists

```
list(1:3, "a", TRUE, 1L)
#> [[1]]
#> [1] 1 2 3
#>
#> [[2]]
#> [1] "a"
#>
#> [[3]]
#> [1] TRUE
#>
#> [[4]]
#> [1] 1
```

## Data Types: Data frames

Data frames are the most common way of storing data in R. Essentially, a data frame is a list of equal-length vectors. Given this construction, it shares some properties of both the matrix and the list, which can be very convenient for data analysis.

Use `data.frame()` for initialization. Multiple data frames (of appropriate dimension) may be concatenated via `rbind()` or `cbind()`.

Example: Data frames

```
data.frame(x = 2:4, y = c("a", "b", "c"))
#>   x y
#> 1 2 a
#> 2 3 b
#> 3 4 c
```

## Data Types: Which one to choose?

Sometimes the data gives you little choice in the data structure. More often than not, however, you have to make the decision. Doing so incorrectly can have noticeable consequences. Poor choices can make your code needlessly complicated, slow, and memory-intensive.

As a simple guideline for the choice between data frames and matrices, consider how linear algebra intensive your code is. If you're using linear algebra operations frequently, use a `matrix` instead of a data frame.

For advanced projects (outside the scope of this course), also consider other data structures provided by R packages (e.g., `data.table`, sparse matrices).

## Operators

Operators are convenient shortcuts to many of the key operations that frequently arise in data analysis. The key groups of operators are

▶ Assignment operators

▶ Arithmetic operators

▶ Relational operators

▶ Logical operators

▶ Subsetting operators

We'll discuss each briefly and conclude with a note on operator preference.

(There are also two linear algebra operators. We'll turn to them later.)

## Operators: Assignment

Assignment associates an internal object (e.g., a vector of integers) with a variable in your code. The key assignment operator in R is "<-". It's not the only one – but it's the only one you should use.

There are other assignment operators you may sometimes see in other's code. "=" is identical to "<-" and both deal with local assignment. The global assignment operator in R is "<<-".

Assignment: Assignment Operators

```
x <- c(1, 2)
x # print(x)
#> [1] 1 2

y <- cbind(x, x)
y # print(y)
#>      [, 1] [, 2]
#> [1, ]   1    1
#> [2, ]   2    2
```

## Operators: Arithmetic

The familiar arithmetic operators are implemented with their usual symbols. In particular, "+" is addition, "−" is subtraction, "∗" is multiplication, "/" is division, and "ˆ" is the exponent operator. There are two other useful operators: "%%" calculates the modulus and "%/%" calculates integer division.

Note that these operations are element-wise! R is vectorized.

Example: Arithmetic Operators

```
c(1, 2, 3, 4) + 1
#> [1] 2 3 4 5

c(1, 2, 3, 4) * c(1, 2, 3, 4)
#> [1] 1 4 9 16

c(3, 4, 5, 6) %% 3
#> [1] 0 1 2 0
```

Relational operators are very convenient for generation of logical variables. In R, the key relational operators are the familiar "greater than" and "smaller than" operators – ">" and "<" – as well as their weak extensions – ">=" and "<=". Equality is indicated by a double equality sign: "==". Not-equal is using the negation operator: "!=".

As before: these operations are element-wise!

Example: Relational Operators

```
x <- 1
y <- x + 1
c(y < x, y <= x, y != x)
#> [1] FALSE FALSE TRUE

(c(3, 4, 5, 6) %% 3) == 0
#> [1] TRUE FALSE FALSE TRUE
```

## Operators: Logical

Logical operators are convenient for aggregating logical variables. Here, "&" is element-wise *logical-and*, and "|" is element-wise *logical-or*. The negation operator is "!".

In addition to the element-wise logical operators, there are also non-element-wise alternatives – "&&" and "||". Both only evaluate the first logical variable. Be cautious when using them.

Example: Logical Operators

```
!(c(TRUE, FALSE) & TRUE)
#> [1] FALSE TRUE

x <- 1
y <- x + 1
(y < x - 1) | (y > x)
#> [1] TRUE

(c(1, 2, 3, 4) == c(1, 5, 5, 5)) && (x == 1)
#> [1] TRUE
```

## Operators: Subsetting

Subsetting allows you to retrieve and manipulate parts of your data structure. This is crucial for succinctly expressing complex operations.

R has three subsetting operators:
- ▶ [ for vectors, matrices, arrarys, and data frames
- ▶ [[ for lists
- ▶ $ for named lists and data frames

Example: Subsetting Operators

```
df <- data.frame(x = 2:4, y = c("a", "b", "c"))
df[, 1]
#> [1] 2 3 4

df[, 1] == df[, "x"] & df[, 1] == df$x
#> [1] TRUE TRUE TRUE
```

Precedence rules for operators can be complicated. Rather than memorize these rules and risk confusion, use parentheses to avoid ambiguity.

Example: Operator Precedence

```
TRUE || FALSE == FALSE || FALSE
#> [1] TRUE

(TRUE || FALSE) == (FALSE || FALSE)
#> [1] FALSE

TRUE || (FALSE == FALSE) || FALSE
#> [1] TRUE
```

## Functions

Functions are a fundamental building block of R: Anything that "happens" is the result of a function call.

Much like it's mathematical equivalent, a function takes some inputs (called *arguments*) and returns an output. In R, a function can only return a single object (e.g., a vector or a list).

Writing your own functions is – thankfully – incredibly easy. If you catch yourself copy-pasting lines of code, consider writing your own function and calling it instead.

Example: Custom Functions

```
myplus <- function(x, y) {
  res <- x + y
  return(res)
}#MYPLUS

myplus(1, 2) == (1 + 2)
#> [1] TRUE
```

## Functions

Base R implements many key functions for data analysis. Being familiar with them is important for efficient programming.

Some of the key function groups are:

▶ Distributional functions

▶ Statistical functions

▶ Logical functions

▶ Linear Algebra

▶ Importing and Exporting Data

We'll briefly go over each now.

## Functions: Distributions

The four key distributional functions we need are

▶ the probability density (mass) function,

▶ the cumulative density (mass) function,

▶ the quantile function,

▶ the sampling function for generating random draws.

R implements these for all standard distributions with the same pattern:

(d, p, q, r) * (binom, chisq, exp, f, norm, t, unif)

where d denotes the pdf (pmf), p denotes the CDF (CMF), q denotes the quantile function, and r denotes the sampling function.

binom denotes the binomial distribution, chisq denotes the $\chi^2$-distribution, exp denotes the exponential distribution, f denotes the F-distribution, norm denotes the normal distribution, t denotes the t-distribution, and unif denotes the uniform distribution.

(Note: There are many other distributions implemented in R!)

## Functions: Sample Statistics

Sample statistics are functions of the observed data. They are useful for summarizing key features of the data.

Particularly useful functions are:

▶ sum, mean, sd, var

▶ cov, cor

▶ quantile

▶ table

Example: Sample Statistics

```
x <- rexp(n = 100000, rate = 2) # E[X] = 1/2
mean(x)
#> [1] 0.501

x <- rnorm(n = 10000, mean = 0, sd = 1)
quantile(x = x, probs = c(0.025, 0.975))
#>   2.5%  97.5%
#> -1.967  1.947
```

# Functions: Sample Statistics (Contd.)

In the example below, we test the fairness of an unfair coin given a sample of $n = 100$ observations. Recall that a binomial random variable with 1 trial is simply a Bernoulli random variable.

Example: Test fairness of a coin flip

```
# Simulate data from Bernoulli s.t. E[X] = 0.6.
#      That is, the coin is unfair!
n <- 100
x <- rbinom(n = n, size = 1, prob = 0.6)

# Calculate test statistic for H_0: E[X] = 0.5.
z <- sqrt(n) * (mean(x) - 0.5) / sd(x)

# Conduct two sided test using the CLT.
2 * pnorm(q = -abs(z), mean = 0, sd = 1)
#> [1] 0.002
```

## Functions: Linear Algebra

Linear algebra is incredibly useful for econometrics and statistics – both for theory and for implementation.

The key functions are

- ▶ the transpose function: `t`
- ▶ the inverse function: `solve`
- ▶ matrix multiplication: `%*%`
- ▶ out product: `%o%`

Implementing linear regression from scratch has never been so easy!
Recall $\hat{\beta}_{ols} = \left(X^\top X\right)^{-1} X^\top y$.

Example: Linear Algebra

```
myols <- function(y, X) {
  beta <- solve(t(X) %*% X) %*% (t(X) %*% y)
  return(beta)
}#MYOLS
```

## Functions: Logical functions

The key logical functions are

- any
- all
- which

Two others that are sometimes convenient are which.min and which.max.

Example: Logical Functions

```
x <- c(TRUE, FALSE, TRUE)
any(x)
#> [1] TRUE

all(x)
#> [1] FALSE

which(x)
#> [1] 1 3
```

## Functions: Importing and Exporting Data

Here we discuss importing and exporting comma separated values-files (i.e., csv-files). This can easily be done with `read.csv` and `write.csv`.

There are many R packages for importing other file formats. See, e.g., `readstata13`, `haven`, etc.

Example: Importing and Exporting Data

```
# Simulate some data
n <- 100
df <- data.frame(x = rnorm(n, 0, 1),
                 y = rbinom(n, 1, 0.6))

# Export to .csv
write.csv(df, file = "my_simulated_data.csv",
          row.names = FALSE)

# Import from .csv
df_2 <- read.csv(file = "my_simulated_data.csv")
head(df_2$y) # prints first 5 entries
#> [1] 0 1 0 1 1
```

## Functions: Packages

Base R is great (and you should use it) but there are excellent extensions available in the form of packages.

Official packages are hosted on CRAN. These packages are well maintained and have to adhere to certain quality standards. Examples are ggplot2 for plotting or dyplr for data manipulation.

In-official packages (or development versions) are often hosted on GitHub. Using the package devtools, these can be directly installed as well.

After installation, packages need to be activated each time you want to use them. This is done via library.

Example: Installing Packages

```
# CRAN packages
install.packages(c("ggplot2", "dyplr", "devtools"))

# GitHub package
library(devtools)
install_github("jkcshea/l1svr") # needs devtools
```

## Documentation

Before calling a function, it's important you familiarize yourself with the syntax and default arguments. The perfect resource for this are the R help files. These can be conveniently accessed via the "?" operator.

The following is a screenshot from the output after calling ?rnorm.



Other helpful resources are Google, Stack Overflow, Stack Exchange, etc.

# Style-Guide

"*Good coding style is like using correct punctuation. You can manage without it, but it sure makes things easier to read.*" – Wickham (2019)

Adopting (and adapting) a so-called style guide is the better alternative to re-inventing punctuation from scratch. The following gives a few highlights from the `tidyverse` style guide (`https://style.tidyverse.org/`).

We focus on:
▶ Naming
▶ Spacing
▶ Comments
▶ Sections

## Style-Guide: Naming

Coming up with good variable names can be challenging!

▶ Variable and function names should be in lowercase.
▶ Use an underscore to separate words withing a name.
▶ Variable names should be nouns and function names should be verbs.
▶ Strive for names that are clear and concise.

Example: Variable names

```
# Good
sample_mean_x <- mean(x)
avg_x <- mean(x)

# Bad
SampleMean_1 <- mean(x)
avgx <- mean(x)
```

## Style-Guide: Spacing

▶ Place spaces around all infix operators (=, +, -, <-, etc.). The exception is ":", which should be used without spaces.

▶ Always put a space after a comma, and never before.

▶ Place spaces before left parentheses, except in a function call.

▶ Don't add space within parentheses.

Example: Spacing

```
# Good
x <- 1:10
if (1 %in% x) print("Wow!")
x <- c(x, x)

# Bad
x <- 1 : 10
if(1%in%x)print("Wow!")
x <- c ( x , x )
```

## Style-Guide: Comments

- ▶ Use # to add a comment.
- ▶ Comment frequently but not excessively.
- ▶ For multi-line comments, add *four* spaces at the beginning of each line after the first one.

Important: Limit your code to 80 characters per line. Many IDEs add a vertical line to help with this.

Example: Multi-line comment

```
# Commenting code is important. Be concise. When
#     necessary, use indentation to increase
#     readability.
```

Example: Simple comment

```
x <- 1
x <- x + 1 # This works well for brief comments
```

## Style-Guide: Sections

Use comments to break up a single R script into multiple sections. For example, each section dedicated to a single exercise of a problem set.

Using "=" does the job nicely:

Example: Sectioning an R Script

```
# Section name ===================================
# Adding a brief explanation here may be good.

x <- 1 # Some code
x <- x + 1

# Another section name ===========================
# Etc...
```

# Coding in R Studio

# Coding in R Studio (Contd.)

There are essentially three options for executing code in R Studio:

- ▶ Entering code directly in the R console. Note that this makes reproduction (practically) impossible.

- ▶ Selecting a single line in the R script and hitting enter. This runs the line of code and advances to the next line.

- ▶ Block-selecting multiple lines of code and hitting enter. This runs all selected lines in sequence.

Frequently run your code to check whether everything works as intended.

# References

Wickham, H. (2019). *Advanced R.* CRC press.