

## Cross References

John Doe

2025-04-11



# Contents

<b>1</b>	<b>About</b>	<b>5</b>
1.1	Usage . . . . .	5
1.2	Render book . . . . .	5
1.3	Preview book . . . . .	6
<b>2</b>	<b>Rstudio</b>	<b>7</b>
2.1	Dark Theme . . . . .	12
2.2	Update R . . . . .	14
2.3	Packages Management . . . . .	15
2.4	Using Git with RStudio . . . . .	22
2.5	Copilot . . . . .	23
2.6	Save R Workspace . . . . .	23
2.7	Options . . . . .	25
2.8	R Startup . . . . .	26
<b>3</b>	<b>Knit Rmd</b>	<b>29</b>
3.1	Chunk Options . . . . .	35
3.2	Print Verbatim R code chunks . . . . .	38
3.3	Rmd Basics . . . . .	41
3.4	Citing References . . . . .	43
3.5	Cross References . . . . .	43
3.6	Equations . . . . .	47
<b>4</b>	<b>Basic R</b>	<b>49</b>
4.1	Data Input & Output . . . . .	51
<b>5</b>	<b>Machine Learning</b>	<b>55</b>
5.1	Random Forest . . . . .	56
5.2	Neural Network . . . . .	60



# Chapter 1

## About

This is a *sample* book written in **Markdown**. You can use anything that Pandoc’s Markdown supports; for example, a math equation  $a^2 + b^2 = c^2$ .

### 1.1 Usage

Each **bookdown** chapter is an .Rmd file, and each .Rmd file can contain one (and only one) chapter. A chapter *must* start with a first-level heading: **# A good chapter**, and can contain one (and only one) first-level heading.

Use second-level and higher headings within chapters like: **## A short section** or **### An even shorter section**.

The **index.Rmd** file is required, and is also your first book chapter. It will be the homepage when you render the book.

### 1.2 Render book

You can render the HTML version of this example book without changing anything:

1. Find the **Build** pane in the RStudio IDE, and
2. Click on **Build Book**, then select your output format, or select “All formats” if you’d like to use multiple formats from the same book source files.

Or build the book from the R console:

```
bookdown::render_book()
```

To render this example to PDF as a `bookdown::pdf_book`, you'll need to install XeLaTeX. You are recommended to install TinyTeX (which includes XeLaTeX): <https://yihui.org/tinytex/>.

### 1.3 Preview book

As you work, you may start a local server to live preview this HTML book. This preview will update as you edit the book when you save individual .Rmd files. You can start the server in a work session by using the RStudio add-in “Preview book”, or from the R console:

```
bookdown::serve_book()
```

## Chapter 2

# Rstudio

### Rstudio shortcuts

Command Palette: `⌘ + P`, all shortcuts can be accessed via the Command Palette.

keyboard combination	function
<code>opt + _</code> <code>ESC</code> or <code>ctrl + C</code> <code>⌘ + M</code>	insert assignment operator <code>&lt;-</code> exit + prompt Add magrittr's pipe operator “%>%” After R4.1, you can set this too native pipe <code> &gt;</code>
<code>ctrl + [/]</code> <code>cmd + D</code> <code>cmd + 1</code> <code>cmd + 2</code> <code>ctrl + shift + S</code>	indent or unindent delete one row move cursor to console window move cursor to editor window add 80 hyphens <code>---</code> to signal a new chapter (Addin)
<code>ctrl + shift + =</code>	add 80 equals <code>===</code> to signal a new Chapter (Addin)
<code>shift + cmd + N</code> <code>cmd + ↑ / ↓</code>	new R script in console, get a list of command history
<code>shift + ↑ / ↓</code> <code>fn + F2</code>	select one line up/down <code>view()</code> an object, don't select the object
<code>cmd + shift + 1</code> <code>ctrl (+ shift) + tab</code>	activate X11() window next (last) tab in scriptor (this applies to all apps); hit <code>ctrl</code> first, then <code>shift</code> if necessary, last tab

**Source**

keyboard combination	function
cmd + return	Run current line/selection
opt + return	Run current line/selection (retain cursor position)

**Rmd related**

keyboard combination	function
cmd + shift + K	<b>Knit</b> rmd
cmd + opt + C	run current code chunk in Rmd
cmd + opt + I	insert code chunks in Rmd, i.e., <code>```{r}</code> and <code>```</code>

---

Q: How to print output in console rather than inline in Rmd?

A: Choose the gear in the editor toolbar and choose “Chunk Output in Console”.

---

Q: How to insert Emojis in Rmd?

A: There are several options:

- You can type directly a lot of Emojis, such as and . Try this first, if it doesn’t show properly, then try the following solutions.
  - If the emoji can show in the script, then you can use it directly.
- Using a html tag, e.g., `<span> </span>` will show like this

This seems to be the most straightforward solution to me.

Note that the emoji won’t display correctly in your Rmd file, but when you render the Rmd and deploy to html pages, the emoji will show properly.

- Using Hexadecimal code. (You need to look up the code somewhere, which is a hassle. )

We can add emojis to an HTML document by using their hexadecimal code. These code starts with `&#x` and ends with `;` to specify browser that these are hexadecimal codes. For example,

```
<p>Smily face <span>&#x1F600;</span> </p>
```



will give you

Smily face

Go to this site: <https://emojipedia.org/emoji/>

Grab the **codepoint** for the emoji you want (e.g., U+1F600 for grinning face)

Replace U+ with `&#x` so it becomes `&#x1F600`, and add a semicolon ; at the end.

Finally, enclose that into an html tag, e.g., `<span>`.

- With RStudio Visual mode. (You need to change mode back and forth. )

First change to the Visual mode. To insert an emoji, you can use either the **Insert** menu or the requisite markdown shortcut plus auto-complete:

I am personally NOT a fan of Visual Mode because it changes your source code silently ...

---

### Set working directory

```
# get the dir name of the current script
dir_folder <- dirname(rstudioapi::getSourceEditorContext())$path)
setwd(dir_folder) # set as working dir
```

RStudio projects are associated with R working directories. You can create an RStudio project:

- In a brand new directory
- In an existing directory where you already have R code and data
- By cloning a version control (Git or Subversion) repository

Why using R projects:

1. I don't need to use `setwd` at the start of each script, and if I move the base project folder it will still work.
2. I have a personal package with a custom project, which creates my folders just the way I like them. This makes it so that the basic locations for data, outputs and analysis is the same across my work.

Double-click on a `.Rproj` file to open a fresh instance of RStudio, with the working directory and file browser pointed at the project folder.

Q: What is an **R session**? And when do I use it?

A: Multiple concurrent sessions can be useful when you want to:

- Run multiple analyses in parallel
- Keep multiple sessions open indefinitely
- Participate in one or more shared projects

---

### Launch a new project-less RStudio session

```
# run in console
rstudioapi::terminalExecute("open -n /Applications/RStudio.app", show = FALSE)
```

-n Open a new instance of the application(s) even if one is already running.

`rstudioapi::terminalExecute(command, workingDir = NULL, env = character(), show = TRUE)` tells R to run the system command in quotes.

- `command` System command to be invoked, as a character string.
- `workingDir` Working directory for command
- `env` Vector of name=value strings to set environment variables
- `show` If FALSE, terminal won't be brought to front

The `rstudioapi` package provides an interface for interacting with the RStudio IDE with R code. Using `rstudioapi`, you can:

- Examine, manipulate, and save the contents of documents currently open in RStudio,
  - Create, open, or re-open RStudio projects,
  - Prompt the user with different kinds of dialogs (e.g. for selecting a file or folder, or requesting a password from the user),
  - Interact with RStudio terminals,
  - Interact with the R session associated with the current RStudio instance.
- 

### Set up Development Tools

<https://cran.r-project.org/bin/macosx/tools/>

- install Xcode command line tools

```
sudo xcode-select --install
```

- install GNU Fortran compiler

Using **Apple silicon** (aka arm64, aarch64, M1) Macs Fortran compiler

- Go to <https://www.xquartz.org/>, download the .dmg and run the installer.
- Verify that build tools are installed and available by opening an R console and running

```
install.packages("pkgbuild")
pkgbuild::check_build_tools()
```

---

### Insert Code Session

To insert a new code section you can use the **Code -> Insert Section** command. Alternatively, any comment line which includes at least four trailing dashes (-), equal signs (=), or pound signs (#) automatically creates a code section.

### Define your own shortcuts

<https://www.statworx.com/ch/blog/defining-your-own-shortcut-in-rstudio/>

<https://www.r-bloggers.com/2020/03/defining-your-own-shortcut-in-rstudio/>

Install the shortcut packages.

Add code session separators, --- or ==.

```
install.packages(
  # same path as above
  "~/Downloads/shoRtcut_0.1.0.tar.gz",
  # indicate it is a local file
  repos = NULL)
install.packages(
  # same path as above
  "~/Downloads/shoRtcut2_0.1.0.tar.gz",
  # indicate it is a local file
  repos = NULL)
```

Now go to Tools > Modify Keyboard Shortcuts and search for “dashes”. Here you can define the keyboard combination by clicking inside the empty Shortcut field and pressing the desired key-combination on your keyboard. Click Apply, and that’s it!

---

### Tips and Tricks

- To add comments to a function, you can type “Roxygen comment” into the Command Palette ( + +P) while the cursor is in a function and it will automatically add a template structure for writing a comment about your function.

Keyboard shortcut: R

- Snippets are a way to make a shortcut for inserting text based on a “code”.

To find the snippets and edit them, use the Palette (Cmd-Shift-P) and type “edit snippets”. There you will find some predefined snippets. You can also create your own.

For instance, when in an R script (or code chunk), typing “fun” followed by pressing Tab, a template for a function will be inserted that looks like:

```
name <- function(variables) {
}
```

You can just fill in the name of the function, then press **Tab** to move to the variables, change the name, then press **Tab** again to move to the function code area and write your function without moving your fingers from the keyboard.

- Show argument definitions as you type functions.

When you type an existing R function such as `round()`, not only does **tab** give you the options, but there's an explanation beneath each variable, telling you its role in the function:

## 2.1 Dark Theme

<https://community.rstudio.com/t/fvleature-req-word-background-highlight-color-in-find-and-spellcheck/18578/3>

<https://rstudio.github.io/rstudio-extensions/rstudio-theme-creation.html>

<https://docs.posit.co/ide/user/ide/guide/ui/appearance.html#creating-custom-themes-for-rstudio>

Theme repositories

- **rstudiothemes**: <https://github.com/max-alletsee/rstudio-themes>
- **rsthemes**: <https://www.garrickadenbuie.com/project/rsthemes/>

RStudio and Editor themes are two different things

- RStudio theme applies to the IDE's framework; including Modern (default), Classic, Sky, and Dark.

The Sky theme is similar to the Modern theme, except for the tab and toolbar headers.

The dark theme is a superset to the Modern and Sky themes that is activated whenever the Editor theme uses a dark palette.

- Editor theme applies to the source pane.

A useful tool to customize your editor theme: <https://tmtheme-editor.github.io/#!/editor/theme/Monokai>

Embedded themes can be found here: <https://github.com/rstudio/rstudio/tree/87e129853121106a87e92df416363f39da95f82e/src/cpp/session/resources/themes>

**Useful elements:**

`.ace_marker-layer .ace_selection` Changes the color and style of the highlighting for the currently selected line or block of lines.

`.ace_marker-layer .ace_bracket` Changes the color and style of the highlighting on matching brackets.

**Recommended highlight color:** `rgba(255, 0, 0, 0.47)`

If you really like one of the default themes RStudio provides, but you want to tweak some small things, you can go the theme directory and change the element's appearance.

RStudio's **default editor theme directory** on Mac:

Right click `RStudio.app`, "Show Package Contents" to navigate to the application folder.

`/Applications/RStudio.app/Contents/Resources/resources/themes/ambiance.rstheme`  
(deprecated)

New editor theme directory: `/Applications/RStudio.app/Contents/Resources/app/resources/themes/ambianc`

You may also find the default themes on GitHub repo: <https://github.com/rstudio/rstudio/tree/master/src/cpp/session/resources/themes>

If you want to install or create a completely new theme, use the **Custom theme (user-defined) folder**:

- `~/.config/rstudio/themes/idle_fingers_2.rstheme` on mac
- `viridis-theme`

```
/* yaml tag */
.ace_meta.ace_tag {
  color: #2499DA;
}
/* quoted by $...$ and code chunk options */
.ace_support.ace_function {
  color: #55C667;
}
```

See [HERE](#) for common selectors you can use.

A collection of screenshots of default RStudio themes: <https://www.trifields.jp/list-of-rstudio-editor-themes-2520>

Q: The margin line is too bright.

A: Change the `.ace_print-margin` element.

```
.ace_print-margin {
  width: 1px;
  background: #e8e8e8;
}
```

`#e8e8e8` is the culprit here, and should be darkened. I changed it to `#2F3941`.

Source: <https://github.com/rstudio/rstudio/issues/3420#issuecomment-453154475>

### Install custom themes

- Using `rstudiothemes` pkg

Go to the repository to see which theme you want to use. Then install the theme. Themes can be applied to RStudio via “Tools” - “Global Options” - “Appearance” - “Add Theme”.

```
# install the pseudo-package from this Github repository
devtools::install_github("max-alletsee/rstudio-themes")

library(rstudiothemes) # ... then load the library

# example 1: bulk-install all light themes
install_rstudio_themes(theme = "all_light")

# example 2: install two specific light themes
install_rstudio_themes(theme = c("Ayu Light", "Github {rsthemes}"))

# example 3: install one specific dark theme
install_rstudio_themes(theme = "49th Parallel")
```

- Using `rstudioapi` package’s “addTheme” function

```
# create temporary download directory
theme_49th_parallel <- fs::path_temp("49th_parallel-RStudio",
                                     ext = "rstheme")

# download theme from github
download.file("https://raw.githubusercontent.com/wvictor14/rstudio_themes/master/49th_parallel",
             theme_49th_parallel)

# apply the theme
rstudioapi::addTheme(theme_49th_parallel,
                    apply = TRUE)
```

---

## 2.2 Update R

Q: How to tell which version of R you are running?

A: In the R terminal, type `R.version`.

The key thing to be aware of is that when you update R, if you just download the latest version from the website, you will lose all your packages!

The easiest way to update R and not cause yourself a huge headache is to use the `installr` package. When you use the `updateR()` function, a series of dialogue boxes will appear. These should be fairly self-explanatory but there is a full step-by-step guide available for how to use `installr`, the important bit is to select “Yes” when it asked if you would like to copy your packages from the older version of R.

```
# Install the installr package
install.packages("installr")

# Load installr
library(installr)

# Run the update function
updateR()
```

---

## 2.3 Packages Management

### 2.3.1 Load packages

Q: What is the difference btw `library(package)` and `require(package)`?

A:

- `library(package)` returns an error if the package doesn't exist.
- `require(package)` returns `FALSE` if the package is not found and `TRUE` if the package is loaded. `require` is designed for use inside other functions, such as using the value it returns in some error checking loop, as it outputs a warning and continues if the package is not found.

Q: How to reload a package after updating?

A: Call `detach(package:pkg, unload = TRUE)` or `unloadNamespace` first, then use `library(pkg)` to reload. If you use `library` on a package whose namespace is loaded, it attaches the exports of the already loaded namespace. So detaching and re-attaching a package may not refresh some or all components of the package, and is inadvisable. The most reliable way to completely detach a package is to restart R.

For example, if we want to detach `ggplot2` package, we can use

```
detach(package:ggplot2, unload=TRUE)
```

`requireNamespace` can be used to *test* if a package is installed and loadable because it comes back with either `TRUE` (if found the pkg) or `FALSE` (if failed to find the pkg).

```
> !requireNamespace("ggplot2")
[1] FALSE
> !requireNamespace("ggplot3")
Loading required namespace: ggplot3
Failed with error: 'there is no package called 'ggplot3''
[1] TRUE
```

To see whether need to install some packages:

```
# install the package if it is not available
if (!requireNamespace("devtools")) install.packages("devtools")
# or equivalently
if (!require("devtools")) install.packages("devtools")
```

You can also use `require(devtools)` to check whether the required package is available, but note that it will load the package as a side effect.

Alternatively,

```
# short command
"ggplot2" %in% installed.packages()
# full command
"ggplot2" %in% rownames(installed.packages())
```

`installed.packages()` Finds details of all packages installed in the specified library path `lib.loc`. Returns a matrix of package names, library paths and version numbers.

```
> installed.packages() %>% class()
[1] "matrix" "array"

> installed.packages() %>% str()
chr [1:355, 1:16] "abind" "alphavantager" "anytime" "askpass" "assertthat" "backports"
- attr(*, "dimnames")=List of 2
..$ : chr [1:355] "abind" "alphavantager" "anytime" "askpass" ...
..$ : chr [1:16] "Package" "LibPath" "Version" "Priority" ...
```

The following code can be used to load packages for your project and set up the working environment.

```
# load the pkg, if not found, install then load
require(dplyr) || {install.packages("dplyr"); require(dplyr)}
require(odbc) || {install.packages("odbc"); require(odbc)}
require(DBI) || {install.packages("DBI"); require(DBI)}
```

If using `library()`, will return error if some package is not installed and interrupt the program.

If it is a list of packages you want to check, use `lapply` to loop through all



packages.

```
## First specify the packages of interest
packages = c("MASS", "nlme")

## Now load or install&load all
package.check <- lapply(
  packages,
  FUN = function(x) {
    if (!require(x, character.only = TRUE)) {
      install.packages(x, dependencies = TRUE)
      library(x, character.only = TRUE)
    }
  }
)
```

You can then use `search()` to determine whether all the packages have loaded.

```
search()
[1] ".GlobalEnv"      "package:nlme"      "package:MASS"
[4] "package:stats"    "package:graphics"  "package:grDevices"
[7] "package:datasets" "renv:shims"         "package:utils"
[10] "package:methods" "Autoloads"         "package:base"
```

Q: `dplyr` has many conflicts with `plyr`.

A: Specify pkg using `::`. Or set library priority by

- changing the order in which you load the packages.

```
# load dplyr last so that it has priority
library(plyr)
library(dplyr)
```

- with the `{needs}` package

```
library(needs)
# prioritize the functions in dplyr
prioritize(dplyr)
```

Q: How to unload a package without restarting R?

A: `detach("package:ggplot2", unload=TRUE)` or uncheck the checkbox button in Packages pane.

## 2.3.2 Install packages

Install R packages from source

```
# From local tarball
install.packages(
  # indicate path of the package source file
  "~/Documents/R/UserPackages/shoRtcut2_0.1.0.tar.gz",
  # indicate it is a local file
  repos = NULL)

# From github
install.packages("Rcpp", repos="https://rcppcore.github.io/drat")
```

### Install from GitHub

```
devtools::install_github(repo, ref="HEAD", subdir = NULL)
```

- **repo** repository address in the format `username/repo[/subdir][@ref|#pull]`. Alternatively, you can specify `subdir` and/or `ref` using the respective parameters. If both are specified, the values in `repo` take precedence.
- **ref** Desired git reference. Could be a commit, tag, or branch name, or a call to `github_pull()` or `github_release()`. Defaults to "HEAD", which means the default branch on GitHub and for git remotes.

Ex

```
# install version 3.5.1
install_github("tidyverse/ggplot2", ref="ggplot2 3.5.1")
```

---

Check installed packages

```
# print all installed packages
rownames(installed.packages())
# check if `ggplot2` is installed
"ggplot2" %in% rownames(installed.packages())
```

```
installed.packages(lib.loc=NULL, priority=NULL)
```

- `lib.loc` character vector describing the location of **R** library trees to search through
- `priority` used to select packages; "high" is equivalent to `c("base", "recommended")`

```
# list all bases packages from your `R.Version`
> rownames(installed.packages(priority="base"))
[1] "base"      "compiler" "datasets"  "graphics"  "grDevices" "grid"
[7] "methods"   "parallel"  "splines"   "stats"     "stats4"    "tcltk"
[13] "tools"     "utils"

# what R loads on startup
> c(getOption("defaultPackages"), "base")
[1] "datasets" "utils"     "grDevices" "graphics"  "stats"     "methods"   "base"
```

`getOption("defaultPackages")` is what R loads on startup although the `basepackage` is not counted.

Check package version

```
packageVersion("ggplot2") # check package version
```

Q: How do I know if I have the latest version?

A: You can go to GitHub repo to check release notes. You will find the latest version of packages there.

### 2.3.3 Update packages

- Update an individual package

- Using `install.packages`

```
install.packages("ggplot2") # update one specific package
```

- Using `update.packages`

```
update.packages(oldPkgs = "ggplot2")
```

Note that you need to specify `oldPkgs` explicitly as it is a named argument.

- Update ALL outdated packages

```
## update all installed packages in a stated library location, default to `.libPaths()`
update.packages(lib.loc = .libPaths(), ask = TRUE)
```

`update.packages` updates ALL outdated packages in a stated library location. That library location is given by the first argument (if not supplied it works on all known library locations for the current R session).

It will ask you for every package if you want to update.

To just say **yes** to everything, use `ask = FALSE`.

```
update.packages(ask = FALSE)
```

Unfortunately this won't update packages installed by `devtools::install_github()`

#### Troubleshooting

Q: I ran `update.packages("ggplot2")`, but nothing happened. No output on console, no error, nothing.

A: The first argument specifies the library location you want to search through (and update packages therein). `update.packages("ggplot2")` means you want

to update the packages in library location `ggplot2`, which is most unlikely to exist on your R installation.

---

Q: I tried to update `ggplot2` with `install.packages("ggplot2")`, but nothing happened.

A: If `ggplot2` is already loaded, then you can't install `ggplot2` in the current session now. If you need to, save any objects you can't easily recreate, and quit out of R. Then start a new R session, immediately run `install.packages("ggplot2")`, then once finished, load the package and reload in any previously saved objects.

---

More about `update.packages`:

- `update.packages(lib.loc = NULL, repos = getOption("repos"), ask = TRUE)`: First a list of all packages found in `lib.loc` is created and compared with those available at the repositories. If `ask = TRUE` (the default) packages with a newer version are reported and for each one the user can specify if it should be updated. If so the packages are downloaded from the repositories and installed in the respective library path (or `instlib` if specified).
- You can specify one specific package to update using `update.packages(oldPkgs = "ggplot2")`. It will check updates only for that package and ask you if you want to update.

The easiest way to update an individual package is just to use `install.packages`. It is a one step command, compared to `update.packages`, which first checks and then asks.

- `update.packages` returns `NULL` invisibly.
  - Be aware that some package updates may cause your previous code to stop working. For this reason, we recommend updating all your packages once at the beginning of each academic year (or semester) – don't do it before an assessment or deadline just in case!
- 

### Updating all Packages after R update

R packages are missing after updating. So you have to save the installed packages and re-install them after updating.

- Alternatively, `installr::updateR()` automatically updates R and installs your packages.

Here is how to do it manually.

```
## get packages installed
packs <- as.data.frame(installed.packages(.libPaths()[1]), stringsAsFactors = F)
# Save to local
f_name <- "~/Documents/R/packages.csv"
rownames(packs)
write.csv(packs, f_name, row.names = FALSE)
packs <- read_csv(f_name)
packs
## Re-install packages using install.packages() after updating R
install.packages(packs$Package)
```

R library path /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/library

- use `find.package("ggplot2")` to find the location to where the given package is found.
- alternatively, you can run `.libPaths()`
  - `.libPaths()` without an argument will return a list of all the places R will currently look for loading a package when requested.
  - `.libPaths("a/b/c")` with an argument will add that new directory ("a/b/c") to the ones R was already using. If you use that directory often enough, you may wish to add that call to `.libPaths("a/b/c")` in your `.Rprofile` startup file in your home directory.

---

## Put your R package on GitHub

Reference: [https://jennybc.github.io/2014-05-12-ubc/ubc-r/session2.4\\_github.html](https://jennybc.github.io/2014-05-12-ubc/ubc-r/session2.4_github.html)

- Change to the package directory
- Initialize the repository with `git init`
- Add and commit everything with
  1. `git add .` stage changes;
  2. `git status` optional check staged changes, but yet to submit;
  3. and `git commit` submit staged changes.
- Create a new repository on GitHub
- Connect your local repository to the GitHub one

```
# add repo name "origin" to the remote repo at the URL
git remote add origin https://github.com/username/reponame
```

- Push everything to github

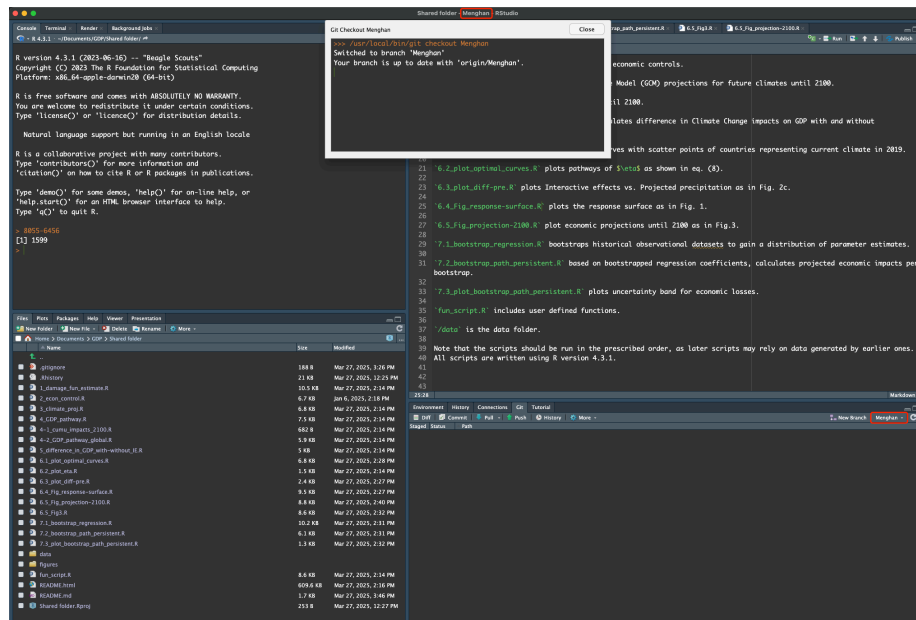
```
# rename the current local branch to "main"
git branch -M main
```

```
# creates a remote branch "origin" and sets it upstream of the "main" branch
git push -u origin main
```

## 2.4 Using Git with RStudio

Before you start coding, make sure that you are on the correct branch. You may check

- from the Git tab on the Environment, History, Connections, ... pane
- you can also see from the status bar on the very top of the window. The words are formatted as “Projection Name – Branch – RStudio”.



Choose a License for your repo

Q: Which open source license is appropriate for my project?

A: See <https://opensource.guide/legal/#which-open-source-license-is-appropriate-for-my-project>.

Q: How to add a license to my repo?

A: Follow the instructions here.

## 2.5 Copilot

Copilot offers autocomplete-style suggestions as you code as “ghost text”.

GitHub Copilot primarily relies on the context in the file you are actively editing. Any comments, code, or other context provided within the active document will be used as a “prompt” that Copilot will then use to provide a suggested completion.

- To expand the scope of the context used by Copilot beyond just the active document, there is a setting to also index and read from other R, Python, or SQL files in the current project. This setting can be toggled on or off in the Tools > Global Options > Copilot > “Index project files with GitHub Copilot” setting.
- At times, normal autocomplete and Copilot may seem to conflict with each other. In these cases, it is best to review the Copilot suggestion and determine if it is appropriate for the current context. If it is, you can accept the suggestion by pressing **Tab**. If it is not, you can ignore the suggestion and continue typing or force the normal autocomplete to show by pressing **Ctrl+Space**.

Issue: **Ctrl+Space** crash with spotlight shortcut.

- Only show Copilot suggestions when evoke manually using **Ctrl + Backslash (\)**.

The Copilot suggestions can be very distracting and clutter your script.

For general advice on how to use Copilot, please see:

- [RStudio Copilot User Guide](#)
- [How to use GitHub Copilot: Prompts, tips, and use cases](#)

---

## 2.6 Save R Workspace

If you want to save all objects in your work space, use `save.image()`. It will create an image of your current variables and functions, and save them to a file called `.RData`. When R next loads, objects stored in this image are by default restored.

This sounds convenient, however, you do NOT want to do this because this corrupt reproducibility of your project.

You want to start from a clean slate every time.

It is suggested change RStudio Global Options to

- *not* “restore `.RData` into workspace at startup”, and
- *never* “save workspace to `.RData` on exit”.

In case you do feel the need to save the workspace, use the following cmd.

```
save.image(file = ".RData", version = NULL, ascii = FALSE, compress
= !ascii, safe = TRUE)

## save current workspace ##
f_name <- "RImage/TCR_2023-05-09.RData"
f_name
save.image(f_name)
# load(f_name)
```

Q: Can I save the loaded packages in the current session/workspace?

A: The workspace is for *objects* like data and functions. Starting R with particular packages loaded is what your `.Rprofile` file is for, and you can have a different one in each directory. But I'd recommend not saving anything between R sessions and instead recreate it all using code. This is much more likely to lead to reproducible results.

### History

When you quit a project, `.Rhistory` is automatically written to the project directory unless you opt out to. It contains a history of all of the commands that you have sent to the R console in this session.

---

### Pop out an editor

Click the **Show in New Window** button in any source editor tab.

To return a document to the main window, click the **Return to Main Window** button on the editor toolbar.

---

### Environment Pane

By default, the Environment pane is located in the top-right and includes the Environment, History, Connections, Build, and Version Control System (VCS) tabs.

#### Version Control System (VCS)

The VCS tab will change based on the version control system you have enabled for that session. For example, using Git will change the tab name to Git and provide some common commands for viewing diffs, committing changes, pull and push ...

**Output pane**  
The Output pane displays various outputs such as plots, HTML content, or on-disk files. It contains the Files, Plots, R Packages, Help, Viewer, and Presentation tabs.

Ref: RStudio Pane Layout



---

### Global Options that make coding easier

- Syntax highlight and matched parentheses.

Under “Tools -> Global Options -> Code -> Display”, under **Syntax section**, check the boxes for **highlight R function calls** and **use rainbow parentheses**. The second is especially useful to mark matching opening and closing brackets.

- Show whitespace characters.

In “Tools -> Global Options -> Code -> Display”, check “Show whitespace characters”. This will let you see spaces and newlines in the editor.

References:

<https://coding-club.rostools.org/posts/tips-and-tricks/>

---

## 2.7 Options

**getOption(x)** Allow the user to set and examine a variety of global *options* which affect the way in which **R computes and displays its results**. Use **getOption** to check default values of global options.

- **x** a character string holding an option name, must be quoted in quotes
- Can only query one option at a time. If multiple options are given, will return the value of the first option.

**options(...)** query and modify global options.

- ... any options can be defined, using **name = value**.

Note that you do NOT need to quote your option name here!

- **options()** with no arguments returns a list with the current values of the options.
- **options("name")** can be used to examine options' current value too; return a *list*, whereas **getOption("name")** returns the value only.
  - Note that you need to quote the option name when you do queries.
  - You can query more than one options at a time.

```
> options("width", "digits")
$width
[1] 90
$digits
```

```
[1] 7

> getOption("width", "digits")
[1] 90
```

`?options` to get the help page of global options. To check which options are available and their definitions.

### Use examples

```
## Two ways checking default option values
> options("width")
$width
[1] 81

> getOption("width")
[1] 81

## Change option values
# use name=value
> options(width=80, digits=15) # set print width, digits to print for numeric values u
# use a named list
> options(list(width=80, digits=15))
```

### Commonly used global options:

Option	Description
width	Controls the maximum number of columns on a line used in printing vectors, matrices and arrays, and when filling by <code>cat</code> . Defaults to 80. Don't change this if you want to print more columns. Use <code>options(tibble.width=400)</code> instead.
pillar.sigfig	Tibbles print numbers with three significant digits by default, switching to scientific notation if the available space is too small. <code>options(pillar.sigfig = 4)</code> to increase the number of digits printed out

## 2.8 R Startup

`Sys.getenv(x)` get the values of the environment variables. Returns a vector of the same length as `x`.

- `x` a character vector

Environment Variables examples:

```
> Sys.getenv(c("HOME", "R_HOME", "R_PAPERSIZE", "R_PRINTCMD"))
      HOME                               R_HOME
"/Users/menghan" "/Library/Frameworks/R.framework/Resources"
      R_PAPERSIZE                       R_PRINTCMD
      "a4"                               "lpr"
```

Rstudio doesn't load Rprofile or Renvirom

I store my Rprofile and Renvirom in non-default places (i.e. ~/.config/R). When opening R in a normal shell, my environment is loaded perfectly fine. When opening Rstudio, it doesn't load my options, settings or paths.

- Have to wrap your option settings in `rstudio.sessionInit`

<https://damien-datasci-blog.netlify.app/post/2020-12-31-pimp-your-r-startup-message/>

- Open .Rprofile

```
usethis::edit_r_profile()
```

- wrap up your options in the following snippet

```
setHook("rstudio.sessionInit", function(newSession) {
  # any code included here will be run at the start of each RStudio session
  options(buildtools.check = function(action) TRUE )
}, action = "append")
```

- Understanding R's startup

<https://rviews.rstudio.com/2017/04/19/r-for-enterprise-understanding-r-s-startup/>

`usethis` is a workflow package: it automates repetitive tasks that arise during project setup and development, both for R packages and non-package projects.

What is .Rprofile?

.Rprofile is a startup file to set options and environment variables. .Rprofile files can be either at the user or project level.

- User-level .Rprofile files live in the base of the user's home directory, and
- project-level .Rprofile files live in the base of the project directory.

Quitting R will erase the default theme setting. If you load `ggplot2` in a future session it will revert to the default gray theme. If you'd like for `ggplot2` to always use a different theme (either yours or one of the built-in ones), you can set a load hook and put it in your .Rprofile file. For example, the following hook sets the default theme to be `theme_minimal()` every time the `ggplot2` package is loaded.

```
setHook(packageEvent("ggplot2", "onLoad"),
        function(...) ggplot2::theme_set(ggplot2::theme_bw()))
```

Of course, you can always override this default theme by adding a theme object to any of your plots that you construct in `ggplot2`.

Rcpp compilation breaks in R 4.1.0

```
devtools::build("my_package")
Error: Could not find tools necessary to compile a package
Call `pkgbuild::check_build_tools(debug = TRUE)` to diagnose the problem.
```

- In RStudio, I am continually prompted to install additional build tools and I can't install the build tool. → Bypass the option `options(buildtools.check = function(action) TRUE)`.
- Turns out R was pointing to an old clang version in my Makevars.

I just deleted it using [in Terminal]

```
sudo rm ~/.R/Makevars
```

### Install SDK command line tool

Download from [developer.apple.com](https://developer.apple.com). Software development kit.

<https://developer.apple.com/download/all/>

### R compiler tools for cpp on MacOS

- <https://thecoatlessprofessor.com/programming/cpp/r-compiler-tools-for-rcpp-on-macos/>
- install OpenMP enabled `clang` from the terminal  
<https://rpubs.com/Kibalnikov/776164>

## Chapter 3

# Knit Rmd

**R Markdown** is a powerful tool for combining analysis and reporting into the same document. R Markdown has grown substantially from a package that supports a few output formats, to an extensive and diverse ecosystem that supports the creation of books, blogs, scientific articles, websites, and even resumes.

Nice documentations

- R markdown: The definitive guide. provides detailed references
- R markdown cookbook concise and covers essential functions, with examples.

Q: What is the difference between Rmd and R script?

A:

- An R script (`.R`) is used for developing and troubleshooting code; a place where you can store reusable code fragments.
- An R Markdown file (`.Rmd`) is used to integrate R commands with explanatory text and output, making it useful for creating reports.

---

**Quick takeaways:**

- Can still use horizontal separator `ctrl + shift + S` for dashed lines and `ctrl + shift + =` for equals
- Headers must have one empty line above and below to separate it from other text

---

YAML metadata

Q: What is YAML?

A: YAML is a human-friendly data serialization language for all programming languages.

Q: What does YAML do?

A: It is placed at the very beginning of the document and is read by each of Pandoc, **rmarkdown**, and **knitr**.

- Provide metadata of the document.
- located at the top of the file.
- adheres to the YAML format and is delimited by lines containing three three dashes (---).

It can set values of the template variables, such as **title**, **author**, and **date** of the document.

- The **output** field is used by rmarkdown to apply the output format function `rmarkdown::html_document()` in the rendering process.

There are two types of output formats in the **rmarkdown** package: documents (e.g., `pdf_document`), and presentations (e.g., `beamer_presentation`).

Supported output format examples: `html_document`, `pdf_document`.

R Markdown documents (`html_documents`) and R Notebook documents (`html_notebook`) are very similar; in fact, an R Notebook document is a special type of R Markdown document. The main difference is using R Markdown document (`html_documents`) you have to knit (render) the entire document each time you want to preview the document, even if you have made a minor change. However, using an R Notebook document (`html_notebook`) you can view a preview of the final document without rendering the entire document.

- Many aspects of the LaTeX template used to create PDF documents can be customized using *top-level* YAML metadata (note that these options do not appear underneath the **output** section, but rather appear at the top level along with **title**, **author**, and so on). For example:

```
---
title: "Crop Analysis Q3 2013"
output: pdf_document
fontsize: 11pt
geometry: margin=1in
---
```

A few available metadata variables are displayed in the following (consult the Pandoc manual for the full list):

Variable	Description
<code>lang</code>	Document language code
<code>fontsize</code>	Font size (e.g., <code>10pt</code> , <code>11pt</code> , or <code>12pt</code> )
<code>documentclass</code>	LaTeX document class (e.g., <code>article</code> )
<code>classoption</code>	Options for documentclass (e.g., <code>oneside</code> )
<code>geometry</code>	Options for geometry class (e.g., <code>margin=1in</code> )
<code>mainfont</code> , <code>sansfont</code> , <code>monofont</code> , <code>mathfont</code>	Document fonts (works only with <code>xelatex</code> and <code>lualatex</code> )
<code>linkcolor</code> , <code>urlcolor</code> , <code>citecolor</code>	Color for internal links (cross references), external links (link to websites), and citation links (bibliography)
<code>linestretch</code>	Options for line spacing (e.g. 1, 1.5, 3).

- In PDFs, you can use code, typesetting commands (e.g., `\vspace{12pt}`), and specific packages from LaTeX.

1. The `header-includes` option loads LaTeX packages.

```

---
output: pdf_document
header-includes:
- \usepackage{fancyhdr}
---

\pagestyle{fancy}
\fancyhead[LE,RO]{Holly Zaharchuk}
\fancyhead[LO,RE]{PSY 508}

# Problem Set 12

```

2. Alternatively, use `extra_dependencies` to list a character vector of LaTeX packages. This is useful if you need to load multiple packages:

```

---
title: "Untitled"
output:
  pdf_document:
    extra_dependencies: ["bbm", "threeparttable"]
---

```

if you need to specify options when loading the package, you can add a second-level to the list and provide the options as a list:

```

---
title: "Untitled"
output:
  pdf_document:
    extra_dependencies:
      caption: ["labelfont={bf}"]
      hyperref: ["unicode=true", "breaklinks=true"]
      lmodern: null
---

```

Here are some examples of LaTeX packages you could consider using within your report:

- \* `pdfpages`: Include full PDF pages from an external PDF document within your document.
  - \* `caption`: Change the appearance of caption subtitles. For example, you can make the figure title italic or bold.
  - \* `fancyhdr`: Change the style of running headers of all pages.
- Some options are passed to Pandoc, such as `toc`, `toc_depth`, and `number_sections`. You should consult the Pandoc documentation when in doubt.

```

---
output:
  pdf_document:
    toc: true
    keep_tex: true
---

```

- \* `keep_tex: true` if you want to keep intermediate TeX. Easy to debug. Defaults to `false`.

We can include variables and R expressions in this header that can be referenced throughout our R Markdown document. For example, the following header defines `start_date` and `end_date` parameters, which will be reflected in a list called `params` later in the R Markdown document.

Thus, if we want to use these values in our R code, we can access them via `params$start_date` and `params$end_date`.

Should I use quotes to surround the values?

- Whenever applicable use the unquoted style since it is the most readable.
- Use quotes when the value can be misinterpreted as a data type or the value contains a `:`.

```

# values need quotes
foo: '{{ bar }}' # need quotes to avoid interpreting as `dict` object
foo: '123'       # need quote to avoid interpreting as `int` object

```



```
foo: 'yes'           # avoid interpreting as `boolean` object
foo: "bar:baz:bam"  # has colon, can be misinterpreted as key

# values need not quotes
foo: bar1baz234
bar: 123baz
```

---

## Document dependency

By default, R Markdown produces standalone HTML files with no external dependencies, using `data:URIs` to incorporate the contents of linked scripts, stylesheets, images, and videos. This means you can share or publish the file just like you share Office documents or PDFs. If you would rather keep dependencies in external files, you can specify `self_contained: false`.

Note that even for self-contained documents, MathJax is still loaded externally (this is necessary because of its big size). If you want to serve MathJax locally, you should specify `mathjax: local` and `self_contained: false`.

One common reason to keep dependencies external is for serving R Markdown documents from a website (external dependencies can be cached separately by browsers, leading to faster page load times). In the case of serving multiple R Markdown documents you may also want to consolidate dependent library files (e.g. Bootstrap, and MathJax, etc.) into a single directory shared by multiple documents. You can use the `lib_dir` option to do this. For example:

```
---
title: "Habits"
output:
  html_document:
    self_contained: false
    lib_dir: libs
---
```

---

## Loading LaTeX packages

We can load additional LaTeX packages using the `extra_dependencies` option within the `pdf_document` YAML settings.

This allows us to provide a list of LaTeX packages to be loaded in the intermediate LaTeX output document, e.g.,

```
---
title: "Using more LaTeX packages"
output:
```

```
pdf_document:
  extra_dependencies: ["bbm", "threeparttable"]
---
```

If you need to specify options when loading the package, you can add a second level to the list and provide the options as a list, e.g.,

```
output:
  pdf_document:
    extra_dependencies:
      caption: ["labelfont={bf}"]
      hyperref: ["unicode=true", "breaklinks=true"]
      lmodern: null
```

For those familiar with LaTeX, this is equivalent to the following LaTeX code:

```
\usepackage[labelfont={bf}]{caption}
\usepackage[unicode=true, breaklinks=true]{hyperref}
\usepackage{lmodern}
```

The advantage of using the `extra_dependencies` argument over the `includes` argument introduced in Section 6.1 is that you do not need to include an external file, so your Rmd document can be self-contained.

## Includes

You can do more advanced customization of output by including additional HTML content or by replacing the core Pandoc template entirely. To include content in the document header or before/after the document body, you use the `includes` option as follows:

```
---
title: "Habits"
output:
  html_document:
    includes:
      in_header: header.html
      before_body: doc_prefix.html
      after_body: doc_suffix.html
---
```

An example `header.html` to load a MathJax extension `textmacros`.

```
<script type="text/x-mathjax-config">
  MathJax.Hub.Config({
    loader: {load: ['[tex]/textmacros']},
    tex: {packages: {'[+]' : ['textmacros']}}
  });
</script>
```

## 3.1 Chunk Options

If you want to set chunk options globally, call `knitr::opts_chunk$set()` in a code chunk (usually the first one in the document), e.g.,

```
```{r, label="setup", include=FALSE}
knitr::opts_chunk$set(
  comment = "#>", echo = FALSE, fig.width = 6
)
```
```

Full list of chunk options: <https://yihui.org/knitr/options/>

Chunk options can customize nearly all components of code chunks, such as the source code, text output, plots, and the language of the chunk.

### Other languages are supported in Rmd

You can list the names of all available engines via:

```
names(knitr::knit_engines$get())
## [1] "awk" "bash" "coffee"
## [4] "gawk" "groovy" "haskell"
## [7] "lein" "mysql" "node"
## [10] "octave" "perl" "php"
## [13] "psql" "Rscript" "ruby"
## [16] "sas" "scala" "sed"
## [19] "sh" "stata" "zsh"
## [22] "asis" "asy" "block"
## [25] "block2" "bslib" "c"
## [28] "cat" "cc" "comment"
## [31] "css" "ditaa" "dot"
## [34] "embed" "eviews" "exec"
## [37] "fortran" "fortran95" "go"
## [40] "highlight" "js" "julia"
## [43] "python" "R" "Rcpp"
## [46] "sass" "scss" "sql"
## [49] "stan" "targets" "tikz"
## [52] "verbatim" "theorem" "lemma"
## [55] "corollary" "proposition" "conjecture"
## [58] "definition" "example" "exercise"
## [61] "hypothesis" "proof" "remark"
## [64] "solution" "marginfigure"
```

The engines from `theorem` to `solution` are only available when you use the **bookdown** package, and the rest are shipped with the **knitr** package.

To use a different language engine, you can change the language name in the chunk header from `r` to the engine name, e.g.,

```
```python
x = 'hello, python world!'
print(x.split(' '))
```
```

For engines that rely on external interpreters such as `python`, `perl`, and `ruby`, the default interpreters are obtained from `Sys.which()`, i.e., using the interpreter found via the environment variable `PATH` of the system. If you want to use an alternative interpreter, you may specify its path in the chunk option `engine.path`.

For example, you may want to use Python 3 instead of the default Python 2, and we assume Python 3 is at `/usr/bin/python3`

```
```{python, engine.path = '/usr/bin/python3'}
import sys
print(sys.version)
```
```

- All outputs support markdown syntax.
- If the output is html, you can write in html syntax.

The **chunk label** for each chunk is assumed to be unique within the document. This is especially important for cache and plot filenames, because these filenames are based on chunk labels. Chunks without labels will be assigned labels like `unnamed-chunk-i`, where `i` is an incremental number.

- Chunk label doesn't need a `tag`, i.e., you only provide the `value`.
- If you prefer the form `tag=value`, you could also use the chunk option `label` explicitly, e.g.,

```
```{r, label='my-chunk'}
# one code chunk example
```
```

You may use `knitr::opts_chunk$set()` to change the default values of chunk options in a document.

### Commonly used chunk options

- Complete list here. Or `?opts_chunk` to get the help page.

| Options                | Definitions  |
|------------------------|--|
| <code>echo=TRUE</code> | Whether to display the <b>source code</b> in the output document. Use this when you want to show the output but not the code itself. |
| <code>eval=TRUE</code> | Whether to evaluate the code chunk.  |

| Options                         | Definitions   |
|---------------------------------|---|
| <code>include=TRUE</code>       | Whether to include the chunk <b>output</b> in the output document. If <b>FALSE</b> , nothing will be written into the output document, but the code is still evaluated and plot files are generated if there are any plots in the chunk, so you can manually insert figures later.  |
| <code>message=TRUE</code>       | Whether to preserve messages emitted by <code>message()</code>  |
| <code>warning=TRUE</code>       | Whether to show warnings in the output produced by <code>warning()</code> .   |
| <code>results='markup'</code>   | Controls how to display the text results. When <code>results='markup'</code> that is to write text output as-is, i.e., write the raw text results directly into the output document without any markups. Useful when printing <b>stargazer</b> tables.  |
| <code>comment='##'</code>       | The prefix to be added before each line of the text output. Set <code>comment = ''</code> remove the default <b>##</b> .  |
| <code>fig.keep='high'</code>    | How plots in chunks should be kept. <b>high</b> : Only keep high-level plots (merge low-level changes into high-level plots). <b>none</b> : Discard all plots. <b>all</b> : Keep all plots (low-level plot changes may produce new plots). <b>first</b> : Only keep the first plot. <b>last</b> : Only keep the last plot. If set to a numeric vector, the values are indices of (low-level) plots to keep. If you want to choose the second to the fourth plots, you could use <code>fig.keep = 2:4</code> (or remove the first plot via <code>fig.keep = -1</code> ). |
| <code>fig.align="center"</code> | Figure alignment.   |
| <code>fig.pos="H"</code>        | A character string for the figure position arrangement to be used in <code>\begin{figure}[]</code> .  |
| <code>fig.cap</code>            | Figure caption.   |

`results='markup'` note plural form for results.

- **markup**: Default. Mark up text output with the appropriate environments depending on the output format. For example, for R Markdown, if the text output is a character string "[1] 1 2 3", the actual output that **knitr** produces will be:

```

[1] 1 2 3

```

In this case, `results='markup'` means to put the text output in fenced code blocks ("").

- **asis**: Write text output as-is, i.e., write the raw text results directly into the output document without any markups.

```
```{r, results='asis'}
cat("I'm raw Markdown content.\n")
```
```

Sometime, you encounter the following error messages when you have R codes within `enumerate` environment.

You can't use macro parameter character # in horizontal mode.

By default, knitr prefixes R output with ##, which can't be present in your TeX file.

Solution:

- specify `results="asis"` in code chunks.
- `hold`: Hold all pieces of text output in a chunk and flush them to the end of the chunk.
- `hide` (or `FALSE`): Hide text output.

## 3.2 Print Verbatim R code chunks

### Including verbatim R code chunks inside R Markdown

One solution for including verbatim R code chunks (see below for more) is to insert hidden inline R code (``r``) immediately before or after your R code chunk.

- The hidden inline R code will be evaluated as an inline expression to an empty string by knitr.

Then wrap the whole block within a markdown code block. The rendered output will display the verbatim R code chunk — including backticks.

R code generating the four backticks block:

```
output_code <-
"```\nmarkdown\n```\n{\n  r\n}\nplot(cars)\n```\n\n"
```

Write this code in your R Markdown document:

```
```\nmarkdown\n`r` ```\n{\n  r\n}\nplot(cars)\n```\n\n"
```

```

~~~
~~~~

or

~~~~`markdown
~~~{r}`r ``
plot(cars)
~~~~
~~~~

```

Knit the document and the code will render like this in your output:

```

~~~{r}
plot(cars)
~~~~

```

This method makes use of **Markdown Syntax** for code.

Q: What is the Markdown Syntax for code?

A:

- Inline code use a pair of backticks, e.g., ``code``. To use  $n$  literal backticks, use at least  $n+1$  backticks outside. Note that use a space to separate your outside backticks from your literal backtick(s). For example, to generate ``code``, you use ``` `code` ``` (i.e., two backticks + space + one backtick + `code` + one backtick + space + two backticks). Note that you need to write sequentially.
- Plain code blocks can be written either
  - After three or more backticks (fenced code blocks), or  
Can also use tildes (~)
  - Indent the blocks by four spaces (indented code blocks)

Special characters do not trigger special formatting, and all spaces and line breaks are preserved. Blank lines in the verbatim text need not begin with four spaces.
- Note that code blocks must be separated from surrounding text by blank lines.

If the code itself contains a row of tildes or backticks, just use a longer row of tildes or backticks at the start and end:

```

~~~~~
~~~~~
code including tildes
~~~~~
~~~~~

```

These begin with a row of three or more tildes (~) and end with a row of tildes that must be at least as long as the starting row.

A shortcut form (without braces) can also be used for specifying the language of the code block:

```
```haskell
qsort [] = []
```
```

This is equivalent to:

```
``` {.haskell}
qsort [] = []
```
```

`haskell` is the language class.

You can add more classes, such as `numberLines` for adding line numbers.

This shortcut form may be combined with attributes:

```
```haskell {.numberLines}
qsort [] = []
```
```

Which is equivalent to:

```
``` {.haskell .numberLines}
qsort [] = []
```
```

and

```
<pre id="mycode" class="haskell numberLines" startFrom="100">
  <code>
    primes = filterPrime [2..] where
      filterPrime (p:xs) =
        p : filterPrime [x | x <- xs, x `mod` p /= 0]
  </code>
</pre>
```

If highlighting is supported for your output format and language, then the code block above will appear highlighted, with numbered lines starting with 100, 101, and go on.

---

### References:

<https://yihui.org/en/2017/11/knitr-verbatim-code-chunk/>

<https://support.posit.co/hc/en-us/articles/360018181633-Including-verbatim-R-code-chunks-inside-R-Markdown>



<https://themockup.blog/posts/2021-08-27-displaying-verbatim-code-chunks-in-xaringan-presentations/>

Pandoc's Markdown: <https://pandoc.org/MANUAL.html#fenced-code-blocks>

### 3.3 Rmd Basics

To name a chunk, add the name after `r`, it's not necessary to add `label='chunk-name'`, but it is possible to do so if you prefer the form `tag=value`.

#### The chunk label

- Must be unique within the document. This is especially important for cache and plot filenames, because these filenames are based on chunk labels. Chunks without labels will be assigned labels like `unnamed-chunk-i`, where `i` is an incremental number.
- Avoid spaces (`␣`), periods (`.`), and underscores (`_`) in chunk labels and paths. If you need separators, you are recommended to use hyphens (`-`) instead.

`knitr::opts_chunk$set()` changes the default values of chunk options in a document.

---

#### Unnumbered sections

Add `{-}` at the end of the section title.

```
# Question 1: Variance and Covariance properties {-}
<!-- equivalently, you can use {.unnumbered} -->
# Question 1: Variance and Covariance properties {.unnumbered}
```

Note that the section won't be numbered but will show in the TOC.

If you want to further exclude it from the TOC:

```
# Question 1: Variance and Covariance properties {.unlisted .unnumbered}
```

Headings with `#` will appear in the file outline, which is a convenient feature. So use this method whenever possible.

One exception is level 2 headings in Bookdown:

- By default **Bookdown** starts a new page for each level 2 heading. If you want to keep the style without starting a new page, use an `html` tag. The heading won't be numbered or included in TOC. However, a downside is that the heading won't show up in the file outline either, making them harder to locate.

```
<h2>YAML metadata</h2>
```

---

### Knitting in the global environment

```
rmarkdown::render("/Users/menghan/Library/CloudStorage/OneDrive-Norduniversitet/EK369E,
```

**Advantages:** fast; load and output results in the global environment; easy to inspect afterwards.

Rmd built-in themes for `html` output: <https://rstudio4edu.github.io/rstudio4edu-book/rmd-themes.html>

.Rmd documents can be edited in either source or visual mode. To switch into visual mode for a given document, use the Source or Visual button at the top-left of the document toolbar (or alternatively the `Cmd+Shift+F4` keyboard shortcut).

- Visual mode allows you to preview the effect after having compiled the markdown file.

But it modifies your code silently, be cautious with visual mode.

- More user-friendly in terms of providing dropdown menus for editing.
- Visual mode supports both traditional **keyboard shortcuts** (e.g. `Cmd + B` for bold) as well as markdown shortcuts (using markdown syntax directly). For example, enclose **`**bold**`** text in asterisks or type `##` and press space to create a second level heading.
- One bug for Visual mode is that inside **bullet points**, `$` is automatically escaped as `\$`. In this case, use `cmd+ /` and choose inline math to insert an eqn.
- When type inline equations, first type `$` then the equation, then `$` at last. Do not type `$$` at one time. Otherwise, they will be escaped as regular text.

---

### Comments in Rmd

- In both `html` and `pdf` outputs, use the following to write true comments you don't want to show in the rendered file.

```
<!-- regular html comment -->
```

### Link to an external javascript

```
<SCRIPT language="JavaScript" SRC="my_jxscript.js"></SCRIPT>
```

## 3.4 Citing References

Use `@citationkey` to cite references in text.

Put citations in parentheses, use `[@citationkey]`.

To cite multiple entries, separate the keys by semicolons, e.g., `[@key-1; @key-2; @key-3]`. To suppress the mention of the author, add a minus sign before `@`, e.g., `[-@citationkey]`.

Syntax	Result
<code>@adams1975</code> concludes that ...	Adams (1975) concludes that ...
<code>@adams1975[p.33]</code> concludes that ...	Adams (1975, p. 33) concludes that ...
... end of sentence <code>[@adams1975]</code> .	... end of sentence (Adams, 1975).
<code>[see @adams1975,p.33]</code> .	... end of sentence (see Adams, 1975, p. 33).
delineate multiple authors with colon: <code>[@adams1975; @aberdeen1958]</code>	delineate multiple authors with colon: (Aberdeen, 1958; Adams, 1975)
Check Lo and MacKinlay	Check Lo and MacKinlay (1988,
<code>[-@Lo-Mackinlay1988; -@Lo1989]</code>	1989) for example.
for example.	

### Add an item to a bibliography without using it

By default, the bibliography will only display items that are directly referenced in the document.

If you want to include items in the bibliography without actually citing them in the body text, you can define a dummy `nocite` metadata field and put the citations there.

```
---
nocite: /
      @item1, @item2
---
```

## 3.5 Cross References

### 3.5.1 Using bookdown

You can number and refer to an equation by adding `\begin{equation}` along with a label, provided with `(\#eq:label)`.

- The position of the label matters.
  - For single-lined equations: First write your equation, then append your label `(\#eq:label)`. Otherwise, your equation won't be rendered.

- For multi-lined equations: append `(\#eq:label)` after `\end{split}`, `\end{aligned}` ...
- Note that `\begin{equation}` must NOT be quoted in `$$...$$` for the equation to be rendered.

Otherwise, will cause “Bad math delimiter” error at the time of tex compilation for pdf output. Might be alright for html output though.

Unexpected consequence: Without the `$$...$$`, RStudio won’t provide previews for equations.

- For temporary preview in RStudio at the composing stage, you can enclose the whole math environment in `$$...$$`. But **remember to delete them when you are done** editing the equation.
- See this post for a more efficient workaround.
- You can then refer to the equation in text using `\@ref(eq:CJ)`. Remember to put the label in parentheses ().

General syntax for other environments: `\@ref(type:label)` where `type` is the environment being referenced, and `label` is the chunk label.

This is an equation rendered using bookdown

```
\begin{equation} (\#eq:CJ)
y=\beta_0 + \beta_1x + e_t
\end{equation}
```

will render as

$$y = \beta_0 + \beta_1 x + e_t \quad (3.1)$$

You may refer to it using `eqn \@ref(eq:CJ)`, e.g., see eqn (3.1).

Multilined equations.

```
\begin{equation}
\begin{aligned}
y_i &= f(x_{1i}, x_{2i}, \ldots, x_{Ki}) + \varepsilon_i \\
&= x_{1i} \beta_1 + x_{2i} \beta_2 + \cdots + x_{Ki} \beta_K + \varepsilon_i
\end{aligned} (\#eq:scalar-form)
\end{equation}
```

will render as

$$\begin{aligned} y_i &= f(x_{1i}, x_{2i}, \dots, x_{Ki}) + \varepsilon_i \\ &= x_{1i}\beta_1 + x_{2i}\beta_2 + \dots + x_{Ki}\beta_K + \varepsilon_i \end{aligned} \quad (3.2)$$

You may refer to it using `eqn \@ref(eq:scalar-form)`, e.g., see eqn (3.2) .

Note that

- For HTML output, **bookdown** can only number the equations with labels. Please make sure equations without labels are not numbered by either using the `equation*` environment or adding `\nonumber` or `\notag` to your equations.

---

### Troubleshooting

Issue: Bad math environment delimiter on conversion to pdf when using `equation` or `align`.

Cause: The error happens because I enclosed `\begin{equation}` environment in `$$`. I did this as the dollar signs enable equation rendering and preview in file.

Fix: remove the double signs.

The following equation causes error. Need to remove the dollar signs.

```
$$
\begin{equation}
y=x+2
\end{equation}
$$
```

---

### More examples:

- Headers

```
# Introduction {#intro}

This is Chapter \@ref(intro)
```

- Figures

```
See Figure \@ref(fig:cars-plot)

```{r cars-plot, fig.cap="A plot caption"}
plot(cars) # a scatterplot
```
```

- Tables

```
See Table \@ref(tab:mtcars)

```{r mtcars}
```

```
knitr::kable(mtcars[1:5, 1:5], caption = "A caption")
```

```

- Theorems

```
See Theorem \@ref(thm:boring)

```{theorem, boring}
Here is my theorem.
```
```

- Equations

```
See equation \@ref(eq:linear)

\begin{equation}
a + bx = c (\#eq:linear)
\end{equation}
```

### 3.5.2 Using the LaTeX Way

The **LaTeX** way allows you to assign your own labels by `\tag`. One drawback is that this does not allow preview of equations.

1. Add the following script at the beginning of your document body:

```
<script type="text/x-mathjax-config">
MathJax.Hub.Config({
  TeX: { equationNumbers: { autoNumber: "AMS" } }
});
</script>
```

It configures MathJax to automatically number equations. Source.

2. In the text, use `label{eq:label}`. If you want to provide a specific number to the equation, you can use `\tag{XX.XX}`.
  - Note that `\begin{equation}` is NOT inside `$$ ... $$`!
3. Cite using `$$\ref{eq:label}$$` (no parenthesis) or `$$\eqref{eq:label}$$` (with parenthesis). The dollar sign `$` here around `\ref` and `\eqref` is not essential. Commands work with or without `$`.

Without using the bookdown package.

```
\begin{equation} \label{eq:test} \tag{my custom label}
Y_i = \beta_0 + \beta_1 x_i + \epsilon_i
\end{equation}
```

Cite Equation `\eqref{eq:test}` like this.

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (\text{my label})$$

Refer to the eq (my label)

---

### Reference:

<https://bookdown.org/yihui/bookdown/markdown-extensions-by-bookdown.html#equations>

## 3.6 Equations

Can use `$. . . $` (`$$ . . . $$` for blocks) or `\( . . . \)` (`\[ . . . \]` for blocks) to enclose equations. Difference:

- `$. . . $` provides rendered equation previews in RStudio.
- `\( . . . \)` does not have previews.

Rstudio equation previews do NOT work well with indented equations. → reduce indentation

Rstudio

**Multi-case** functions using `\begin{cases}`

```
\begin{align*}
I_t =
\begin{cases}
1 & \text{if } r_t > 0 \\
0 & \text{if } r_t \leq 0
\end{cases}
\end{align*}
```

will render as

$$I_t = \begin{cases} 1 & \text{if } r_t > 0 \\ 0 & \text{if } r_t \leq 0 \end{cases}$$

For equation numbering support in **bookdown** you need to assign labels.

---

You may refer to an equation using Eq. `\@ref{eq:eq01}`.

```
\begin{align} (\#eq:eq01)
\frac{p(x)}{1-p(x)} = \exp (\beta_0+\beta_1 x) \backslash,.
\end{align}
```

If you want to provide a specific number to the equation, you can use `\tag{XX.XX}`.

- With LaTeX

LaTeX allows custom labels.

$$\frac{p(x)}{1-p(x)} = \exp(\beta_0 + \beta_1 x). \quad (\text{my label latex})$$

My specific label here, see eq (my label latex) (`\eqref{eq:my-label-latex}`).

- 
- With bookdown

bookdown does NOT support custom tag though.

$$\frac{p(x)}{1-p(x)} = \exp(\beta_0 + \beta_1 x). \quad (3.3)$$

My specific label here, see eq (3.3)

---



## Chapter 4

# Basic R

Save & Load R objects

`save(..., f_name)` and `saveRDS()`

`save()` When loaded the named object is restored to the current environment (in general use this is the global environment — the workspace) **with the same name it had when saved**.

`save` writes an external representation of **R** objects to the specified file. The objects can be read back from the file at a later date by using the function `load` or `attach` (or `data` in some cases).

```
save(..., list = character(),      file = stop("'file' must
be specified"),      ascii = FALSE, version = NULL, envir =
parent.frame(),      compress = isTRUE(!ascii), compression_level,
eval.promises = TRUE, precheck = TRUE)
```

- ... The names of the objects to be saved (as symbols or character strings).
- `list` A character vector containing the names of objects to be saved.
  - The names of the objects specified **either** as symbols (or character strings) in ... or as a character **vector** in `list` are used to look up the objects from environment `envir`.
- `file` the name of the file where the data will be saved.

`saveRDS()` doesn't save the both the object and its name it just saves a representation of the object. As a result, the saved object can be loaded into a named object within R that is different from the name it had when originally serialized.

Serialization is the process of converting a data structure or object state into a format that can be stored (for example, in a file or memory buffer, or transmitted across a network connection link) and “resurrected” later in the same or another computer environment.

`saveRDS` works only for saving a single R object, `save()` can save multiple R objects in one file. A workaround for `saveRDS` is to save all target objects in a single R object (e.g., in a `list`), and then use `saveRDS()` to save it at once.

```
datalist = list(mtcars = mtcars, pressure=pressure)
saveRDS(datalist, "twodatasets.RDS")
rm(list=ls())

datalist = readRDS("twodatasets.RDS")
datalist
```

---

### Load R objects

`load(f_name)` to load `.rda` file.

`readRDS(f_name)` to load `.rds` file.

Naming conventions:

- `rda` and `rds` for selected objects
- `.RData` for all objects in your workspace
- The file extensions are up to you; you can use whatever file extensions you want.

An example

```
> require(mgcv)
Loading required package: mgcv
This is mgcv 1.7-13. For overview type 'help("mgcv-package")'.
> mod <- gam(Ozone ~ s(Wind), data = airquality, method = "REML")
> mod

Family: gaussian
Link function: identity

Formula:
Ozone ~ s(Wind)

Estimated degrees of freedom:
3.529 total = 4.529002

REML score: 529.4881
> save(mod, file = "mymodel.rda")
> ls()
[1] "mod"
> load(file = "mymodel.rda")
> ls()
```

```
[1] "mod"

> ls()
[1] "mod"
> saveRDS(mod, "mymodel.rds")
> mod2 <- readRDS("mymodel.rds")
> ls()
[1] "mod" "mod2"
> identical(mod, mod2, ignore.environment = TRUE)
[1] TRUE
```

---

### Save figures in a list

```
p_list <- list(p_ano=p_ano, p_tr=p_tr)
# p_list[[name]] <- p_obj
p_list[[1]]

f_name <- paste0(fig_dir, sprintf("trend_analysis/image_list_%s.rds", con_name))
# saveRDS(p_list, f_name)

# plot in a panel grid
p_allCON <- plot_grid(plotlist=p_list, align="vh", labels=sprintf("(%s)", letters[1:length(p_list)]))
p_allCON
```

---

## 4.1 Data Input & Output

### 4.1.1 Read Data

#### Read Fortran

```
read.fortran(file, format, ..., as.is = TRUE, colClasses = NA)
```

- format Character vector or list of vectors.

#### Read dta

haven::read\_dta() read Stata data file.

```
data <- read_dta("climate_health_2406y1.dta")
# retrieve variable labels/definitions
var_dict <- tibble("name" = colnames(data),
                  "label" = sapply(data, function(x) attr(x, "label")) %>%
                    as.character())
```

```

    )
var_dict

var_label(data$gor) # get variable label
val_labels(data$gor) # get value labels

```

### Read fixed width text files

```
read.fwf(file, widths)
```

- **widths** integer vector, giving the widths of the fixed-width fields (of one line), or list of integer vectors giving widths for multiline records.

`read.table(f_name, header=FALSE, row.names, col.names, sep="", na.strings = "NA")` a very versatile function. Can be used to read `.csv` or `.txt` files.

- **f\_name** path to data.
- **header=FALSE** defaults to `FALSE`, assumes there is no header row in the file unless specified otherwise.
  - If there is a header in the first row, should specify **header=TRUE**.
- **row.names** a vector of row names. This can be
  - a vector giving the actual row names, or
  - a single number giving the column of the table which contains the row names, or
  - character string giving the name of the table column containing the row names.
- **col.names** a vector of optional names for the variables. The default is to use "V" followed by the column number.
- **sep** use white space as delimiter.
  - if it is a `csv` file, use **sep=','** to specify comma as delimiter
- **na.strings = "NA"** a character vector of strings which are to be interpreted as NA values.
  - A useful setting: `na.strings = c("", "NA", "NULL")`

```
read.csv(f_name, header = TRUE, sep = ",", na.strings = "..",
dec=".")
```

- **header = TRUE** whether the file contains the names of the variables as its first line.
- **sep** the field separator string. Values within each row of `x` are separated by this string.
- **na** the string to use for missing values in the data.
- **dec** the string to use for decimal points in numeric or complex columns: must be a single character.
- **fileEncoding** UTF-8

When reading data from github, you need to pass in the raw version of the data in `read.csv()`,

R cannot read the display version.

You can get the URL for the raw version by clicking on the Raw button displayed above the data.

```
# read.table can be used to read txt and csv. Need to specify sep=', ' when reading csv.
data <- read.table("https://raw.githubusercontent.com/my1396/course_dataset/refs/heads/main/boned
data

data <- read.table("https://raw.githubusercontent.com/my1396/course_dataset/refs/heads/main/boned

# can use read_csv or read.csv
data <- read_csv("https://raw.githubusercontent.com/my1396/course_dataset/refs/heads/main/boned
data
```

`read_delim(f_name, delim=";")` allows you to specify the delimiter as ;.

```
readr::read_csv(f_name, na = c("..", NA, ""), locale = locale(encoding =
= "UTF-8"), col_types = cols(Date = col_date(format = "%m/%d/%y"))
) read comma separated values.
```

- `col_types` specify column types. Could be created by `list()` or `cols()`.  
`read_csv` will automatically guess, if you don't explicitly specify column types. You can override column types by providing the argument `col_types`. You don't need to provide all column types, just the ones you want to override.
- By default, reading a file without a column specification will print a message showing what `readr` guessed they were. To remove this message,
  - set `show_col_types = FALSE` for one time setting, or
  - set `options(readr.show_col_types = FALSE)` for the current sessions' global options setting. If want to change permanently every-time when R starts, put `options(readr.show_col_types = FALSE)` in `.Rprofile` as global options.

`read_tsv()` read tab separated values.

`read_csv2(f_name, na = c("..", NA, ""))` use semicolon ; to separate values; and use comma , for the decimal point. This is common in some European countries.

- `locale` The locale controls defaults that vary from place to place. The default locale is US-centric (like R), but you can use `locale()` to create your own locale that controls things like the default time zone, encoding, decimal mark, big mark, and day/month names.
- `locale(date_names = "en", date_format = "%AD", time_format = "%AT", decimal_mark = ".", grouping_mark = ",", tz = "UTC", encoding = "UTF-8", asciify = FALSE)`
  - `decimal_mark` indicate the decimal place, can only be , or .

- **encoding** This only affects how the file is read - readr always converts the output to UTF-8.

### 4.1.2 Write Data

Save data in `uft8` encoding with special language characters

`write_excel_csv()` include a UTF-8 Byte order mark which indicates to Excel the csv is UTF-8 encoded.

`write.csv(x, f_name, row.names=TRUE, fileEncoding = "UTF-8")`

- `x` a matrix or data frame. If not one of the types, it is attempted to coerce `x` to a data frame.

- `write_csv(x)` `x` can only be data frame or tibble. Doesn't support matrix.

```
mat %>% as_tibble(rownames = "rowname") %>% write_csv("mat.csv")
mat %>% write_csv("mat.csv")
```

- `row.names` whether to write row names of `x`. Defaults to `TRUE`.

#### 4.1.2.1 flextable

`flextable` package create tables for reporting and publications.

The main function is `flextable` which takes a `data.frame` as argument and returns a `flextable`. If you are using RStudio or another R GUI, the table will be displayed in the **Viewer** panel or in your default browser.

The package provides a set of functions to easily create some tables from others objects.

The `as_flextable()` function is used to transform specific objects into `flextable` objects. For example, you can transform a crosstab produced with the 'tables' package into a `flextable` which can then be formatted, annotated or augmented with footnotes.

## Chapter 5

# Machine Learning

Parametric models such as generalized linear regression and logistic regression has advantages and disadvantages.

**Strength:**

- The effects of individual predictors on the outcome are easily understood
- Statistical inference, such as hypothesis testing or interval estimation, is straightforward
- Methods and procedures for selecting, comparing, and summarizing these models are well-established and extensively studied

**Disadvantages** in the following scenarios:

- Complex, non-linear relationships between predictors and the outcome
- High degrees of interaction between predictors
- Nominal outcome variables with several categories

In these situations, non-parametric or algorithmic modeling approaches have the potential to better capture the underlying trends in the data.

Here we introduce three models: classification and regression trees (CART), random forests, k-nearest neighbors.

- Classification and regression trees (CART) are “trained” by recursively partitioning the  $d$ -dimensional space (defined by the explanatory variables) until an acceptable level of homogeneity or “purity” is achieved within each partition.
- A major issue with tree-based models is that they tend to be high variance (leading to a high propensity towards over-fitting). Random forests are a non-parametric, tree-based modeling algorithm that is built upon the idea that averaging a set of independent elements yields an outcome with lower variability than any of the individual elements in the set.

This general concept should seem familiar. Thinking back to your introductory statistics course, you should remember that the sample mean,  $\bar{x}$ , of a dataset has substantially less variability ( $\frac{\sigma}{\sqrt{n}}$ ) than the individual data-points themselves ( $\sigma$ ).

Q: What is Bias-Variance Trade-Off in Machine Learning?

A:

- Bias refers to error caused by a model for solving complex problems that is over simplified, makes significant assumptions, and misses important relationships in your data.
- Variance error is variability of a target function's form with respect to different training sets. Models with small variance error will not change much if you replace couple of samples in training set. Models with high variance might be affected even with small changes in training set. High variance models fit the data too well, and learns the noise in addition to the inherent patterns in the data.

## 5.1 Random Forest

### Averaging of independent trees

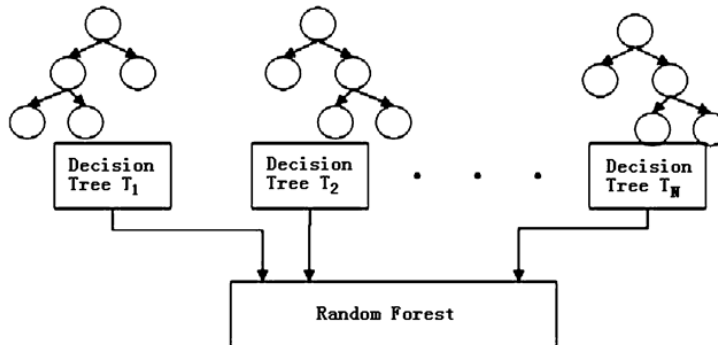
The goal of bagging is to produce  $B$  separate training datasets that are independent of each other (typically  $B$  is in the hundreds). The model of interest (in this case classification and regression trees) is trained separately on each of these datasets, resulting in  $B$  different estimated “models”. These are then averaged to produce a single, low-variance estimate.

Bagging is a general approach, but its most well-known application is in the random forest algorithm:

1. Construct  $B$  bootstrap samples by sampling cases from the original dataset with replacement (this results in  $B$  unique datasets that are similar to the original)
2. Fit a classification and regression tree to each sample, but randomly choose a subset of  $m$  variables that can be used in the construction of that tree (this results in  $B$  unique trees that are fit to similar datasets using different sets of predictors)
3. For a given data-point, each of the  $B$  trees in the forest contributes a prediction or “vote”, with the majority (or average) of these votes forming the random forest's final prediction,  $\hat{y}_i$



```
knitr::include_graphics("images/rf.png")
```



A downside of both the CART and random forest algorithms (as well as many other algorithmic modeling approaches) is an inability to clearly quantify the roles played by individual variables in making predictions. However, the importance of individual variables in a random forest can still be expressed using a measure known as variable importance.

The random forest algorithm requires the following tuning parameters be specified in order to run:

- **ntree** - the number of bagged samples,  $B$ , onto which trees will be grown
- **mtry** - the number of variables that are randomly chosen to be candidates at each split
- Some sort of stopping criteria for individual trees, this can be:
  - **nodesize**, which sets the minimum size of terminal nodes
    - \* larger **nodesize** leads to shallower trees
    - \* smaller node size allows for deeper, more complex trees
  - **maxnodes**, which sets the maximum number of terminal nodes an individual tree can have.

### Applications of Random Forest

Some of the applications of Random Forest Algorithm are listed below:

- **Banking:** It predicts a loan applicant's solvency. This helps lending institutions make a good decision on whether to give the customer loan or not. They are also being used to detect fraudsters.
- **Health Care:** Health professionals use random forest systems to diagnose patients. Patients are diagnosed by assessing their previous medical history. Past medical records are reviewed to establish the proper dosage for the patients.
- **Stock Market:** Financial analysts use it to identify potential markets for stocks. It also enables them to remember the behaviour of stocks.
- **E-Commerce:** Through this system, e-commerce vendors can predict the preference of customers based on past consumption behaviour.

### When to Avoid Using Random Forests?

Random Forests Algorithms are not ideal in the following situations:

- **Extrapolation:** Random Forest regression is not ideal in the extrapolation of data. Unlike linear regression, which uses existing observations to estimate values beyond the observation range.
- **Sparse Data:** Random Forest does not produce good results when the data is sparse. In this case, the subject of features and bootstrapped sample will have an invariant space. This will lead to unproductive spills, which will affect the outcome.

### FAQ

Q: Is RF a linear or non-linear model?

A: RF can capture complex, non-linear relationships.

Q: Is RF sensitive to Imbalanced Data?

A: Yes. It may perform poorly if the dataset is highly imbalanced like one class is significantly more frequent than another.

Q: What is the loss function?

A: Entropy/gini or any other loss function you want.

Q: Difference btw RF and a linear model?

A: A major difference is that a decision tree does not have “parameters”, whereas the linear models need to create a functional form and find the optimal parameters.

### Implementation in R

`ranger` package offers a computation efficient function for RF.

```
RF_ranger <- ranger(formula = formula,
                    data = data_before[idx,],
                    probability = TRUE,
                    importance = "permutation",
                    scale.permutation.importance = TRUE,
                    )
# print(RF_ranger)

rf.pred.test <- predict(RF_ranger, data=data_before[-idx,])$predictions
```

Parameters controlling the general process of RF:

- `probability=FALSE`: Whether to forecast a probability forest.

The hyperparameters `mtry`, `min.node.size` and `sample.fraction` determine the degree of randomness, and should be tuned.

- **mtry=500**: Number of variables to possibly split at in each node in one tree. In plain language, it indicates how many predictor variables should be used in each tree.
  - Default is the (rounded down) square root of the number variables. Alternatively, a single argument function returning an integer, given the number of independent variables.
  - Range btw 1 to the number of predictors.
  - If all predictors are used, then this corresponds in fact to bagging.
- **min.node.size**: The number of observations a terminal node should at least have.
  - Default 1 for classification, 5 for regression, 3 for survival, and 10 for probability. For classification, this can be a vector of class-specific values.
  - Range between 1 and 10
- **sample.fraction**: Fraction of observations to be used in each tree. Default is 1 for sampling with replacement and 0.632 for sampling without replacement. For classification, this can be a vector of class-specific values.
  - Smaller fractions lead to greater diversity, and thus less correlated trees which often is desirable.
  - Range between 0.2 and 0.9

Parameters controlling what and how intermediate results are saved:

- **keep.inbag = FALSE**: Whether to save how often observations are in-bag in each tree.  
Set to **TRUE** if you want to check sample composition in each tree.
- **importance = 'none'|'impurity'|'impurity\_corrected'|'permutation':** Variable importance mode.
- **scale.permutation.importance = FALSE**: Whether to scale permutation importance by standard error as in (Breiman 2001). Only applicable if **'permutation'** variable importance mode selected.
- **write.forest = TRUE**: Whether to save **ranger.forest** object, required for prediction. Set to **FALSE** to reduce memory usage if no prediction intended.
  - Set to **FALSE** when you do parameter tuning.

Q: How to tune hyperparameters?

A: Check out **mlr3** package. Here is an example.

---

## Imbalance Classification

You can balance your random forests using case weights. Here's a simple example:

```

library(ranger)

# Make a dataste
set.seed(43)
nrow <- 1000
ncol <- 10
X <- matrix(rnorm(nrow * ncol), ncol=ncol)
CF <- rnorm(ncol)
Y <- (X %*% CF + rnorm(nrow))[,1]
Y <- as.integer(Y > quantile(Y, 0.90))
table(Y)

# Compute weights to balance the RF
w <- 1/table(Y)
w <- w/sum(w)
weights <- rep(0, nrow)
weights[Y == 0] <- w['0']
weights[Y == 1] <- w['1']
table(weights, Y)

# Fit the RF
data <- data.frame(Y=factor(ifelse(Y==0, 'no', 'yes')), X)
model <- ranger(Y~., data, case.weights=weights)
print(model)

```

Code Source: <https://stats.stackexchange.com/a/287849>

Fixed proportion sampling: <https://github.com/imbs-hl/ranger/issues/167>

---

## References:

[https://remiller1450.github.io/m257s21/Lab10\\_Other\\_Models.html](https://remiller1450.github.io/m257s21/Lab10_Other_Models.html)

## 5.2 Neural Network

Neural networks are made up objects called “layers” and “neurons” and these things connect to each other in a specific way. Each layer has some number of neurons. For example, the first layer might have 10 neurons, the second might have 15, and so on. The number of layers and the number of neurons in each layer is a “hyperparameter”, the user picks how many of each. Let’s take a look at a single neuron.

$$v_3^{(1)} = g(w_3^{(1)}x + b_3^{(1)})$$

- The LHS,  $v_3^{(1)}$ , will be the output. The superscript (1) refers to the layer number; the subscript 3 refers to the neuron.

An output here means just a single number. If we have say 15 neurons (subscript) for this first layer (superscript), then we will have 15 numbers come out of this first layer:  $v^{(1)} = \{v_1^{(1)}, v_1^{(1)}, \dots, v_{15}^{(1)}\}$  where the bolded  $v$  means a vector.

- $x$  is our input vector.
- $w$  is the weight/coefficient vector.
- $b$  is a bias or the intercept term, shifting the value of  $w \cdot x$  up or down.
- $g$  refers to a “non-linear” function, often called as the “activation function”.