# Data technology

---

**Kia Ora!**

- 🎓 I earned my PhD (Stats) @ Monash University, Australia.
- ❤️ My research interests lie in exploratory data analysis, data visualisation, software design, …
- 🧑‍💻 I turn ☕ into > 10 `#rstats` 📦.
- Outside of work, I play 🎾 and make ☕.

---

**Contact**

- ✉️ earo.wang@auckland.ac.nz
- 📌 Office 303.323
- 🕐 Thursday 2-3pm

---

# Data + Technology

➜ **https://stats220.earo.me**

---

## What I mean by "data"

🥫 | 🍅

✖ Stale, uninteresting, convenient
✖ Highly processed and archived
✖ Example: `student tests`, `titanic`, `wages`

✔ Fresh, interesting, challenging
✔ Locally collected and impactful
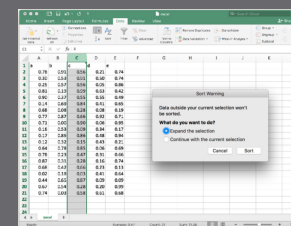✔ Example: Modelling the travel time of transit vehicles in real-time

---

## How I learn new technology

> 🗣️ **Get hands dirty!!**
> 📖 Documentation! Documentation! Documentation!
> 🔍 (Not surprisingly) Learn to google: what that error message means (I google a lot 🫠)

---

# You can't do data science in a GUI

*reference: You can't do data science in a GUI*

## Why programme for data science?

> Programming languages are **languages**.

```r
library(dplyr)
starwars %>%
  group_by(species) %>%
  summarise(
    n = n(),
    mass = mean(mass, na.rm = TRUE)
  ) %>%
  filter(n > 1, mass > 50)
```

> It's just **text**!
>> reproducible, readable, sharable
>> expressive

---

# Why R

> A general-purpose programming language
> Originated by statisticians, a language for statistical analysis
> 292995 + packages on CRAN (Comprehensive R Archive Network, the official repository), Github, etc.
> The tidyverse, a domain specific language in R for data scientists

---

## What R can do?
**- for fun**

### 📦 {cowsay} for generating ASCII picture

```r
library(cowsay)
say("Kia Ora!")
```

```
#>
#>  --------------
#> Kia Ora!
#>  --------------
#>      \
#>       \
#>        \
#>            |\___/|
#>         ==) ^Y^ (==
#>          \  ^  /
#>           )=*=(
#>          /     \
#>          |     |
#>         /| | | |\
#>         \| | |_|/\
#>      jgs //_// ___/
#>           \_)
#>
```

---

## What R can do?
- for fun
**- for data**

### The data science workflow

## What R can do?
- for fun
- for data
**- for communication**

### R Markdown



❯ {rmarkdown} for assignments/reports/papers in `.html` and `.pdf`

❯ {blogdown} for blogs

❯ {bookdown} for books

❯ {xaringan} for slides (220 slides!)

---

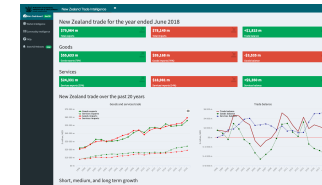**R Markdown documents are fully reproducible: weaving narrative text and code together.**

## What R can do?
- for fun
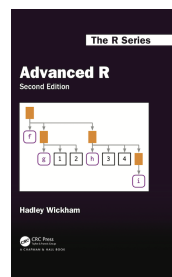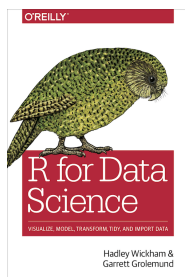- for data
**- for communication**

### R shiny dashboard

❯ Shiny is an R package that makes it easy to build interactive web apps straight from R.



👆 click the image above will take you to the web app, and try to interact with the app.

## Textbook 📚

## At first, you may be like this...    But you can do it!

**Assessments**

- 11 weekly labs 10% (best 10 out of 11)
- 3 assignments 30% (each 10%)
- 1 mid-term test 10% (TBD, possibly week 8)
- 1 final exam 50%

# Project-oriented workflow

> If R were an airplane, RStudio would be the airport, providing many, many supporting services that make it easier for you, the pilot, to take off and go to awesome places. Sure, you can fly an airplane without an airport, but having those runways and supporting infrastructure is a game-changer.
> -- *Julie Lowndes*

**RStudio interface**



This is the scripts pane - the place where you write and save your R code

The environment is where you see the things you create in R. Things like data, objects, and models will appear here as a result of code being run

Files are located here, just like in macOS "Finder" or Window's "My Computer"

The console is where your code gets run. You can write code directly in here, but you might forget the code you wrote, so it's a good idea to write your code above in the scripts pane
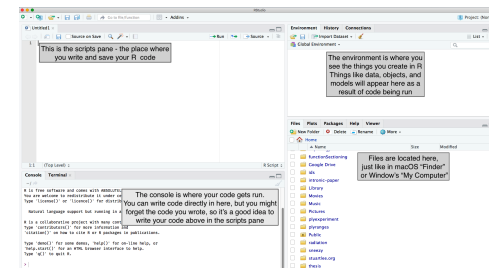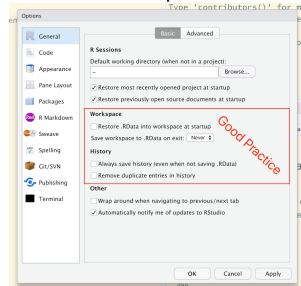
*image credit: Stuart Lee*

## Setting up RStudio (do this once)

Go to **Tools** > **Global Options**:
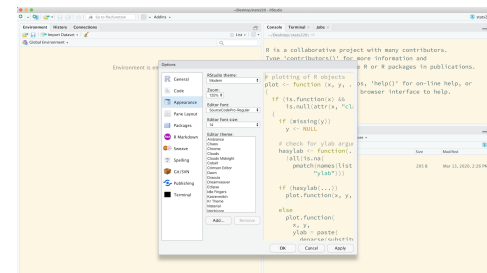


Uncheck `Workspace` and `History`, which helps to keep R working environment fresh and clean every time you switch between projects.

## Your turn

01:00

Change the RStudio appearance up to your taste

## What is a project?

- Each university course is a project, and get your work organised.
- A self-contained project is a folder that contains all relevant files, for example my `stats220/` 🗂 includes:
  - `stats220.Rproj`
  - `data/`
    - `*.csv, *.xlsx`
  - `lectures/`
    - `01-intro.Rmd, 02-import-export.Rmd`
  - `labs/`
    - `lab01.R, lab02.R`
- All working files are relative to the **project root** (i.e. `stats220/`).
- The project should just work on a different computer.

## 🛑 STOP DOING THIS!

Jenny Bryan will set your computer on fire 🔥
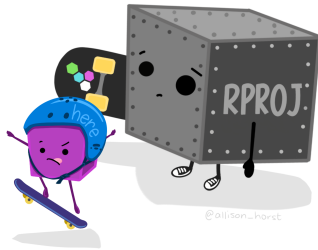
1. if the first line of your R script is

```
setwd("C:\Users\jenny\path\that\only\I\have")
```

2. if the first line of your R script is

```
rm(list = ls())
```

## Create an RStudio project `.Rproj`

1. Click the **Project** icon on the top right corner

2. **New Directory**/**Existing Directory** > **New Project** > **Create Project**

3. Open the project

# 101: syntax and semantics

## Get started
**- assignment**

```
akl_lon <- 174.76
akl_lat <- -36.85
```

⬆ read as "assign the value of `174.76` to an object called `akl_lon`".

An assignment consists of:

> left-hand side: variable names or symbols (`akl_lon`)
> assignment operator: `<-` (RStudio shortcut: `Alt + -`)
> right-hand side: values (`174.76`)

**Get started**

```
akl_lon
```

```
#> [1] 174.76
```

```
akl_lat
```

```
#> [1] -36.85
```

> Names are case sensitive.

```
akl_Lon
```

```
#> Error in eval(expr, envir, enclos): object 'akl_Lon' not found
```

---

**Get started**

**Perform calculations and comparisons**

> Infix operators:
>> +, -, *, /, ^, %% (modulo), %/% (integer division)
>> ==, !=, >, <, >=, <=, %in%

```
akl_lon_region <- akl_lon + c(-1, 1)
akl_lat_region <- akl_lat + c(-.5, .5)

akl_lon_region
```

```
#> [1] 173.76 175.76
```

```
akl_lat_region
```

```
#> [1] -37.35 -36.35
```

---

**Coding style**

> '
>
> *Good coding style is like correct punctuation: you can manage without it, butitsuremakesthingseasiertoread.*
> -- *The tidyverse style guide*

**R style guide**

✔ snake_case          ✘ camelCase (Javascript)
                      ✘ PascalCase (Python)

---

**R** **101: data structures**
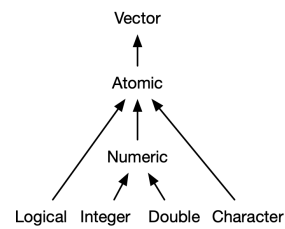
## Atomic vectors



image credit: *Hadley Wickham's **Advanced R***

**Scalars:** *length of 1*

› Logicals: TRUE or FALSE
› Doubles: 174.76, 1.7476e2, Inf, -Inf, NaN (Not a Number)
› Integers: 174L
› Strings: "hello", 'world'

**Vectors:** *values must all be the same type*

```
lgl_vec <- c(TRUE, FALSE)
int_vec <- c(174L, -36L)
dbl_vec <- c(174.76, -36.85)
chr_vec <- c("long", "lat")
```

## Special values

### Missing values

```
NA # Not Applicable
```
```
#> [1] NA
```
```
c(174.76, NA, -36.85)
```
```
#> [1] 174.76     NA -36.85
```
```
length(NA)
```
```
#> [1] 1
```

### The NULL object

```
NULL
```
```
#> NULL
```
```
c(174.76, NULL, -36.85)
```
```
#> [1] 174.76 -36.85
```
```
length(NULL)
```
```
#> [1] 0
```

## Atomic vectors

## Subsetting vectors with []

```
x <- c(akl_lon_region, akl_lat_region)
x
```
```
#> [1] 173.76 175.76 -37.35 -36.35
```

### Positive indices

```
x[c(1, 3)]
```
```
#> [1] 173.76 -37.35
```

### Negative indices

```
x[-c(3, 1)]
```
```
#> [1] 175.76 -36.35
```

## Subsetting vectors with `[]`

**Logical indices**

```
x[c(TRUE, FALSE, TRUE, FALSE)]
```

```
#> [1] 173.76 -37.35
```

```
x[lgl_vec] # recycling
```

```
#> [1] 173.76 -37.35
```

```
x[x > 0]
```

```
#> [1] 173.76 175.76
```

**Special subsetting**

```
x[0]
```

```
#> numeric(0)
```

```
x[]
```

```
#> [1] 173.76 175.76 -37.35 -36.35
```

## Modifying vectors with `[]` on the LHS

```
y <- x
y
```

```
#> [1] 173.76 175.76 -37.35 -36.35
```

```
y[1:3] <- y[1:3] %/% 2
y
```

```
#> [1]  86.00  87.00 -19.00 -36.35
```

> RHS `[]` subsets vector y
> LHS `[]` modifies vector y

## 101: functions

## Function

A function call consists of the function name followed by one or more argument within parentheses.

```
mean(x = x)
```

```
#> [1] 68.955
```

> function name: `mean()`, a built-in R function to compute mean of a vector
> argument: the first argument (LHS x) to specify the data (RHS x)

## Function help page

Check the function's help page with `?mean`

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

> Read **Usage** section
>> What arguments have default values?
> Read **Arguments** section
>> What does `trim` do?
> Run **Example** code

## Function arguments

Match by **positions**

```
mean(x, 0.1, TRUE)
```

```
#> [1] 68.955
```

Match by **names**

```
mean(x, na.rm = TRUE, trim = 0.1)
```

```
#> [1] 68.955
```

## Use functions from packages

```
# install.packages("dplyr")
library(dplyr)
cummean(x)
```

```
#> [1] 173.7600 174.7600 104.0567  68.9550
```

```
first(x)
```

```
#> [1] 173.76
```

```
last(x)
```

```
#> [1] -36.35
```



install.packages("light")  library("light")
Install once
Use many times
Images sourced from https://www.wikihow.com/Change-a-Light-Bulb

## Write your own functions

```
# function_name <- function(arguments) {
#    function_body
# }

my_mean <- function(x, na.rm = FALSE) {
  summation <- sum(x, na.rm = na.rm)
  summation / length(x)
}

my_mean(x)
```

```
#> [1] 68.955
```

**Follow the #rstats community**

**Reading**



> Workflow: basics
> Workflow: scripts
> Workflow: project



> Names and values
> Vectors
> Subsetting