# Organising Your Code

## Data Science Workflows

## Dr Zak Varty

# Outline

Directories >> Files >> Contents

- Focus specifically on code

  - High level approaches to coding

  - Naming conventions

  - Style guides

  - Useful packages

MATH-70076
Effective
Data
Science

# 1. Functional and Object Oriented Programming

# Functional Programming

A functional programming style has two major properties:

- Object immutability,

- Complex programs written using function composition.

Mathematicians often find this way of operating quite intuitive:

$$y = g(x) = f_3 \circ f_2 \circ f_1(x).$$

# The Pipe Operator

Function composition gets messy at ~3 functions

```
1  log(exp(cos(sin(pi))))
```

The `{magrittr}` pipe operator facilitates this compositional style: `%>%`

```
1  library(magrittr)
2  pi %>%
3    sin() %>%
4    cos() %>%
5    exp() %>%
6    log()
```

Recent addition of base R pipe `|>` (R >= 4.1), mostly behaves the same.

```
1  iris |>
2    head(n = 3)
```

MATH-70076
**E**ffective
**D**ata
**S**cience

# When not to pipe !%>%
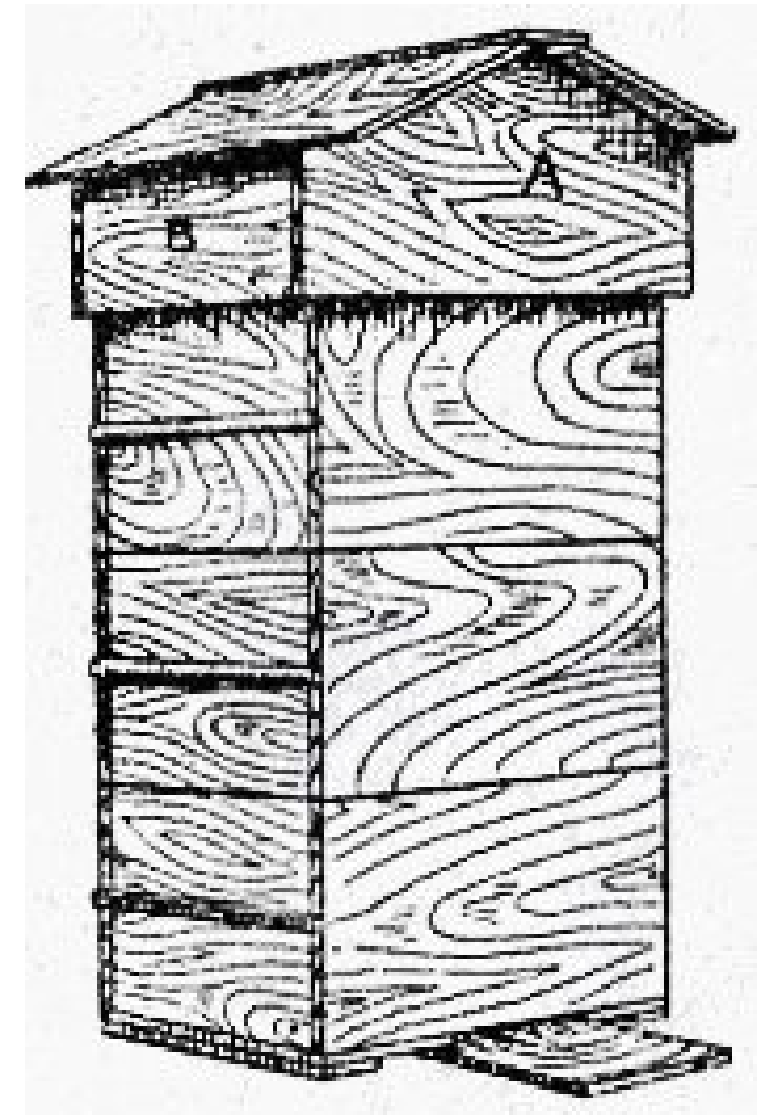
Avoid using the pipe when:

- Manipulating more than one object at a time. (Pipes are for a sequence of steps applied to a single primary object.)

- There are meaningful intermediate objects that could be given informative names.

# Object Oriented Programming

The alternative to functional programming is object oriented programming.

- Solve problems by using lots of simple objects

- R has 3 OOP systems: S3, S4 and R6.

- Objects belong to a class, have methods and fields.

- Agent based simulation of beehive.

# 2. Structuring R scripts

# Top of the Script

- Start script with a comment of 1-2 sentences explaining what it does.

- `setwd()` and `rm(ls())` are the devil's work.

  - "Session" > "Restart R" or Keyboard shortcut: crtl/cmd + shift + 0

- Polite to gather all `library()` and `source()` calls.

- Rude to mess with other people's set up using `install.packages()`.

- Portable scripts use paths relative to the root directory of the project.

# Portable file paths with {here}

```
 1   # Bad - breaks if project moved
 2   source("zaks-mbp/Desktop/exciting-new-project/src/helper_functions/rolling_mean.R")
 3
 4   # Better - breaks if Windows
 5   source("../../src/helper_functions/rolling_mean.R")
 6
 7   # Best - but use here:here() to check root directory correctly identified
 8   source(here::here("src","helper_functions","rolling_mean.R"))
 9
10   # For more info on the here package:
11   vignette("here")
```

For even more on {here}, try: r-wtf chapter, r4ds chapter or project oriented workflow blog post.

MATH-70076
Effective
Data
Science

# 3. Style Gudie Summary

MATH-70076
**E**ffective
**D**ata
**S**cience

# The Body of the Code - Comments

Well commented and organised code easier to read and understand.

```r
 1  # This is an example script showing good use of comments and sectioning
 2
 3  library(here)
 4  source(here("src","helper_functions","rolling_mean.R"))
 5
 6  #========================================================================    <- 80 characters max for readability
 7  # Major Section on Comments ----
 8  #========================================================================
 9
10  #------------------------------------------------------------------------
11  ##  Minor Section on inline comments ----
12  #------------------------------------------------------------------------
13  x <- 1:10 # this is an inline comment
14
15  #------------------------------------------------------------------------
16  ##  Minor Section on full line comments ----
17  #------------------------------------------------------------------------
18  rolling_mean(x)
19  # This is an full line comment
```

MATH-70076
Effective
Data
Science

# Naming Things Revisited: Objects = Nouns

Object names should use only lowercase letters, numbers, and _.

Use underscores (_) to separate words within a name. (snake_case)

Use singular over plural names.

```
1  # Good
2  day_one
3  day_1
4
5  # Bad
6  first_day_of_the_month
7  DayOne
8  dayone
9  djm1
```

MATH-70076
Effective
Data
Science

# Naming Things Revisited: Functions = Verbs

Function names should use only lowercase letters, numbers, and **_**.

Use underscores (**_**) to separate words within a name. (`snake_case`)

Suggest imperative mood, as in a recipe.

Break long functions over multiple lines. 4 vs 2 spaces.

```
1  # Good
2  add_row()
3  permute()
4
5  # Bad
6  row_adder()
7  permutation()
8
9  long_function_name <- function(
10     a = "a long argument",
11     b = "another argument",
12     c = "another long argument") {
13   # As usual code is indented by two spaces.
14  }
```

MATH-70076
**E**ffective
**D**ata
**S**cience

# Casing Consistently

Many options for separating words within names:

- `CamelCase`

- `pascalCase`

- `snakecase`

- `underscore_separated`❤️

- `hyphen-separated`

- `point.separated` 💀

# Style guide in brief

1. Use comments to structure your code

2. Objects = Nouns

3. Functions = Verbs

4. Use snake case and friendly grammar

# Further tips for friendly coding

- Write your code to be easily understood by humans.

- Use informative names, typing is cheap.

```
1  # Bad
2  for(i in dmt){
3    print(i)
4  }
5
6  # Good
7  for(temperature in daily_max_temperature){
8    print(temperature)
9  }
```

- Divide your work into logical stages, human memory is expensive.

MATH-70076
Effective
Data
Science

# Tidyverse Style Guide

For further details on writing clean, readable code see the Tidyverse Style Guide.

# 4. Reduce, Reuse, Recylce

# Writing reusable code (DRY coding)

- **If you do something twice, write a function.**

  - when you write a function, document it.

  - when you write a function, test it.

- **If your function is used in two scripts, it gets it's own script.**

  - name that script after the function & save it in the `src/` directory of your project

- **If your function is used across projects, add it to a package.**

# Remembering how to use your own code

**When you write a function, document it.** What should the documentation contain?

- Inputs

- Outputs

- Example use cases

- Author (if not obvious or working in a team)

# Roxygen for code documentation

The R package **{Roxygen}** can help with this.

1. `install.packages("Roxygen")`

2. With cursor inside function: Code > Insert Roxygen Skeleton

3. Keyboard shortcut: **cmd + option + shift + r** or **crtl + option + shift + r**

4. Fill out relevant fields

MATH-70076
**E**ffective
**D**ata
**S**cience

# Roxygen Example

```r
 1  #' Title
 2  #'
 3  #' @param x
 4  #' @param remove_NA
 5  #'
 6  #' @return
 7  #' @export
 8  #'
 9  #' @examples
10  geometric_mean <- function(x, remove_NA = FALSE){
11    # Function body goes here
12  }
```

# Roxygen Example

```r
 1  #' Calculate the geometric mean of a numeric vector
 2  #'
 3  #' @param x numeric vector
 4  #' @param remove_NA logical scalar, indicating whether NA values should be stripped be
 5  #'
 6  #' @return the geometric mean of the values in `x`, a numeric scalar value.
 7  #'
 8  #' @examples
 9  #' geometric_mean(x = 1:10)
10  #' geometric_mean(x = c(1:10, NA), remove_NA = TRUE)
11  #'
12  geometric_mean <- function(x, remove_NA = FALSE){
13      # Function body goes here
14  }
```

For more on Roxygen, see the package documentation or the chapter of R packages on function documentation.

MATH-70076
**E**ffective
**D**ata
**S**cience

# Checking Your Code

**If you write a function, test it.** Testing code has two main purposes:

- To warn or prevent user misuse (e.g. strange inputs),

- To catch edge cases.

Testing is fun but painful - you are trying as hard as you can to break your beautiful new creation!

MATH-70076
Effective
Data
Science

# An Informal Testing Workflow

1. Write a function

2. Experiment with the function in the console, try to break it

3. Fix the break and repeat.

**Problems:** Time consuming and not reproducible.

# Formalising Our Testing Workflow

- Test in a script, named after the function and stored in the `tests/` directory.

- The R package `{testthat}` provides useful functions for unit testing your code. Check out the {testthat} function reference for more examples.

- We'll explore this in more detail during the live session.

```
1  testthat::expect_equal(
2    object = geometric_mean(x = c(1, NA), remove_NA = FALSE),
3    expected = NA)
4
5  # Error: geometric_mean(x = c(1, NA), remove_NA = FALSE) not equal to NA.
6  # Types not compatible: double is not logical
```

MATH-70076
**Effective**
**Data**
**Science**

# Summary

- Functional and Object Oriented Programming

- Structuring your scripts

- Styling your code

- Reduce, reuse, recycle

- Documenting and testing