



# C91AR | ADVANCED STATISTICS USING R

## Lecture 3: Tidy Code Principles

Dr Pete McKenna

2025-01-28

# LET'S CREATE AN R MARKDOWN FOR TODAY

- Call the file “Lecture3\_tidy-code.rmd”

# PACKAGES FOR TODAY'S SESSION

- Add this package installation to your Markdown file
- Remember to annotate inside your code chunks using `#`
- Also, you can add text summaries outside the code chunks to add clarity and interpretation.

```
1 # Let's stick to pacman to keep loading and installation simple
2 pacman::p_load(tidyverse,
3                   nycflights13,
4                   summarytools)
```

# READING FOR TODAY'S SESSION

- Chapter 2: Workflow Basics
  - covers basics of writing R code (e.g., commenting and functions)
- Chapter 4 Workflow: code style
  - Covers the basics of code formatting
- Why did I ask you to read this?
  - These chapters provide valuable advice about how to set up your project and format your code.
  - The principles covered here will be the building blocks for your R journey and will make life a lot easier for you in the long-term.

# ENCOUNTERING MESSY CODE



- What can you tell me about this bedroom?
- How does this image make you feel?

# WHAT ARE WE GOING TO COVER TODAY?

- This is a hybrid session, involving both lecturing and coding
- We are going to use the `nycflights13` data to learn about and apply tidy code principles
- In doing so, we will get our first glimpse of *tidyverse* code
- Outline plans for Tutorial 2

# WHAT IS THE *TIDYVERSE*?

- The *tidyverse* is a collection of R packages designed for data science.
- These packages share an underlying design philosophy, grammar, and data structures (e.g., datasets), making it easier to learn and use them together.
- The tidyverse simplifies many common data science tasks, making your code more readable and easier to maintain.

# KEY TIDYVERSE VERBS AND SYNTAX

- `mutate()`: adds new variables or modifies existing variables.
- `select()`: picks (and amends) variables based on their names.
- `filter()`: subsets cases based on specified conditional criteria.
- `count()`: counts the number of occurrences of unique values.
- `summarise()`: reduces multiple values down to a single summary.
- `arrange()`: changes the ordering of the vectors.
- `| >`: the **pipe operator** means ‘then do this’
- [More information about Tidyverse verbs here](#)

# WHY ARE YOU TELLING ME THIS NOW?

The examples of tidy code formatting in this lesson use these verbs. We will examine these functions in more detail in future, but it's useful for you to have some grounding before discussion 'tidy code'.

# flights METADATA

```
1 ?flights
```

# A QUICK TOOL TO EXPLORE THE DATA

```
1 # Use summary tools to get an overview of the data  
2 dfSummary(flights)
```

What can you infer from this summary?

# SOME EXAMPLES

# NAMING VARIABLES

```
1 # Strive for:  
2 short_flights <-  
3   flights |> filter(air_time < 60)  
4  
5 # Avoid:  
6 SHORTFLIGHTS <- flights |> filter(air_time < 60)
```

Given what we have learned so far, what do you think this code does?

# SPACES

```
1 # In this code the author has aligned the "=" sign, giving better readability
2 flights |>
3   mutate(
4     speed      = distance / air_time,
5     dep_hour   = dep_time %/% 100, # %/% = integer division that discards
6     dep_minute = dep_time %% 100 # %% = modulus operator returns the remainder
7   )
```

# PIPES

```
1 # Strive for
2 flights |>
3   filter(!is.na(arr_delay), !is.na(tailnum)) |>
4   count(dest)
5
6 # Avoid
7 flights |>filter(!is.na(arr_delay), !is.na(tailnum)) |>count(dest)
```

If I told you that the operator **!** means **ignore**, what do you think this code does?

# LONG, TIDY COMMAND CHAINS

```
1 flights |>
2   group_by(dest) |>
3   summarize(
4     distance = mean(distance),
5     speed = mean(distance / air_time, na.rm = TRUE)
6   ) |>
7   ggplot(aes(x = distance, y = speed)) +
8   geom_smooth(
9     method = "loess",
10    span = 0.5,
11    se = FALSE,
12    color = "white",
13    linewidth = 4
14  ) +
15  geom_point()
```

# CLEANUP

```
1 # Clear data
2 rm(list = ls())  # Removes all objects from environment
3
4 # Clear packages
5 p_unload(all)  # Remove all contributed packages
6
7 # Clear plots
8 graphics.off()  # Clears plots, closes all graphics devices
9
10 # Clear console
11 cat("\014")  # Mimics ctrl+L
```

# RECAP

- You read about tidy code principles
- I introduced you to the *tidyverse*
- We did a little a little exploration and summarising of `nycflights13` dataset
- Showed you how to clean up after yourself.

# NEXT SESSION

- We will be applying these skills to an experiment dataset
- Learn how to read in datasets
- Perform more data cleaning using the *tidyverse*

