

Module 1: Big Picture Ideas

Rebecca C. Steorts (slide and course adaptation from Maria Tackett)

Announcements

- ▶ Resources for extra R review
 - ▶ Learn R: An interactive introduction to data analysis R (focus on Chapters 4 - 6)
- ▶ Readings for next week will be posted later this week

Questions from last class?

Topics

- ▶ Data analysis life cycle
- ▶ Reproducible data analysis
- ▶ Analyzing multivariable relationships

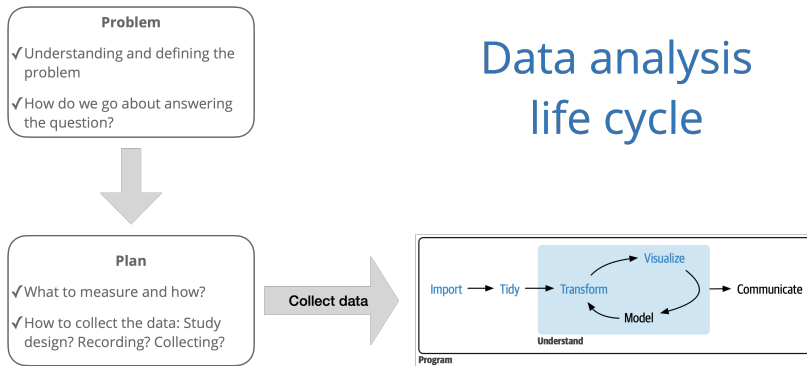


Figure 1: Source: *R for Data Science* with additions from *The Art of Statistics: How to Learn from Data*.

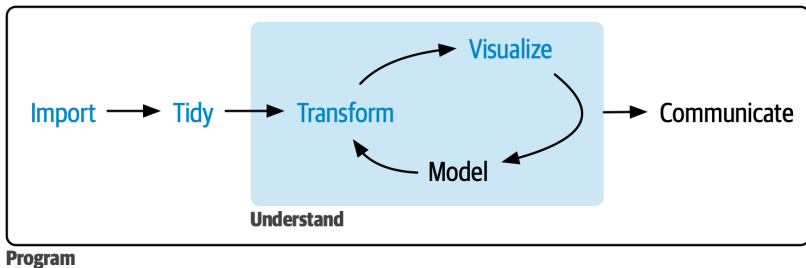


Figure 2: Source: *R for Data Science*

Reproducibility

Reproducibility checklist

What does it mean for an analysis to be reproducible?

Reproducibility checklist

What does it mean for an analysis to be reproducible?

Near term goals:

- ▶ Can the tables and figures be exactly reproduced from the code and data?
- ▶ Does the code actually do what you think it does?
- ▶ In addition to what was done, is it clear *why* it was done?

Reproducibility checklist

What does it mean for an analysis to be reproducible?

Near term goals:

- ▶ Can the tables and figures be exactly reproduced from the code and data?
- ▶ Does the code actually do what you think it does?
- ▶ In addition to what was done, is it clear *why* it was done?

Long term goals:

- ▶ Can the code be used for other data?
- ▶ Can you extend the code to do other things?

Why is reproducibility important?

- ▶ Results produced are more reliable and trustworthy [ostblom2022]
- ▶ Facilitates more effective collaboration [ostblom2022]
- ▶ Contributing to science, which builds and organizes knowledge in terms of testable hypotheses [alexander2023]
- ▶ Possible to identify and correct errors or biases in the analysis process [alexander2023]

When things go wrong

Reproducibility error	Consequence	Source(s)
Limitations in Excel data formats	Loss of 16,000 COVID case records in the UK	(Kelion 2020)
Automatic formatting in Excel	Important genes disregarded in scientific studies	(Ziemann, Eren, and El-Osta 2016)
Deletion of a cell caused rows to shift	Mix-up of which patient group received the treatment	(Wallensteen et al. 2018)
Using binary instead of explanatory	Mix-up of the intervention with the	(Aboumatar and Wise 2019)

Toolkit

- ▶ **Scriptability** → R
- ▶ **Literate programming** (code, narrative, output in one place)
→ Rmarkdown
- ▶ **Version control** → Git / GitHub

R and RStudio

- ▶ R is a statistical programming language
- ▶ RStudio is a convenient interface for R (an integrated development environment, IDE)

R: Engine



RStudio: Dashboard



Figure 3: Source: Statistical Inference via Data Science

RStudio IDE

The screenshot shows the RStudio IDE interface with the following components and annotations:

- Source Editor:** The main area for editing R scripts. It contains a file named `ae-02-bikeshare.qmd`. The text in the editor is:

Data

Our dataset contains daily rentals from the Capital Bikeshare in Washington, DC in 2011 and 2012. It was obtained from the `dcbikeshare` data set in the `dsbox` R package.

We will focus on the following variables in the analysis:

 - `count`: total bike rentals
 - `temp_orig`: Temperature in degrees Celsius
 - `season`: 1 – winter, 2 – spring, 3 – summer, 4 – fall

Click [here](#) for the full list of variables and definitions.

```
1 {r load-data}
2 #| message: false
3 bikeshare <- read_csv("data/dcbikeshare.csv")
```
- Environment:** The top-right pane showing the current environment. It lists the loaded package `ae-02-bikeshare` and the data frame `bikeshare`.
- History:** The top-right pane showing the history of R commands executed.
- Git:** The top-right pane showing the Git status and commit history.
- Files:** The bottom-right pane showing the file explorer. It lists the files in the current project, including `.gitignore`, `.RData`, `.Rhistory`, `ae-02-bikeshare.pdf`, `ae-02-bikeshare.qmd`, `ae-02-bikeshare.Rproj`, `data`, and `README.md`.
- Console:** The bottom-left pane showing the R console output. It displays the R version `R 4.2.0` and the current working directory `~/my-files/teaching/sta210-fa22/ae-repos/ae-02-bikeshare/`.

Colored annotations highlight the following areas:

- Source editor:** The main area for editing R scripts.
- environment - history - git:** The top-right panes.
- console:** The bottom-left pane.
- files - plots - packages - help - viewer:** The bottom-right pane.

Rmarkdown

- ▶ Fully reproducible reports – the analysis is run from the beginning each time you render
- ▶ Code goes in chunks and narrative goes outside of chunks
- ▶ Visual editor to make document editing experience similar to a word processor (Google docs, Word, Pages, etc.)

The image shows the Quarto editor interface. On the left, the source document 'ae-02-bikeshare.qmd' is open. The top toolbar has a 'Render' button circled in red. The document content includes a YAML header, a code chunk for loading packages, and a narrative section about bike rentals in Washington, DC. On the right, the 'Console' pane shows the rendered output of the document, which includes the title, subtitle, date, and the first code chunk's output.

Source Document (ae-02-bikeshare.qmd):

```
1 ---
2 title: "AE 02: Bike rentals in Washington, DC"
3 subtitle: "Exploring and modeling relationships"
4 date: "Sep 05, 2022"
5 format: pdf
6 editor: visual
7 ---
```

Code Chunk:

```
1 {r load-packages}
2 #| message: false
3 library(tidyverse)
4 library(tidymodels)
```

Narrative:

Data

Our dataset contains daily rentals from the Capital Bikeshare in Washington, DC in 2011 and 2012. It was obtained from the `dcbbikeshare` data set in the `dsbox` R package.

We will focus on the following variables in the analysis:

- count: total bike rentals
- temp_orig: Temperature in degrees Celsius

Rendered Output:

AE 02: Bike rentals in Washington, DC
Exploring and modeling relationships

Sep 05, 2022

```
library(tidyverse)
library(tidymodels)
```

Data

Our dataset contains daily rentals from the Capital Bikeshare in Washington, DC in 2011 and 2012. It was obtained from the `dcbbikeshare` data set in the `dsbox` R package.

We will focus on the following variables in the analysis:

- count: total bike rentals
- temp_orig: Temperature in degrees Celsius
- season: 1 - winter, 2 - spring, 3 - summer, 4 - fall

Click [here](#) for the full list of variables and definitions.

```
bikeshare <- read_csv("data/dcbbikeshare.csv")
```

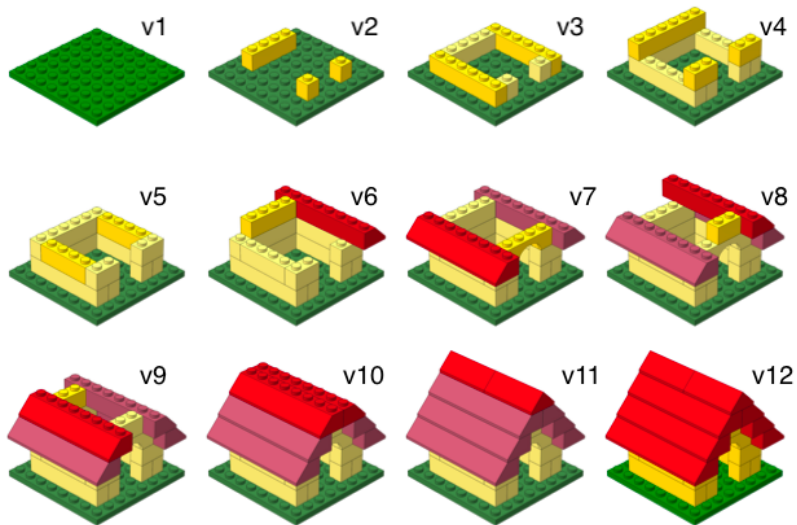
Daily counts and temperature

How will we use Quarto?

- ▶ Every assignment is written in Rmarkdown document
- ▶ You'll have a template in Rmarkdown to start with

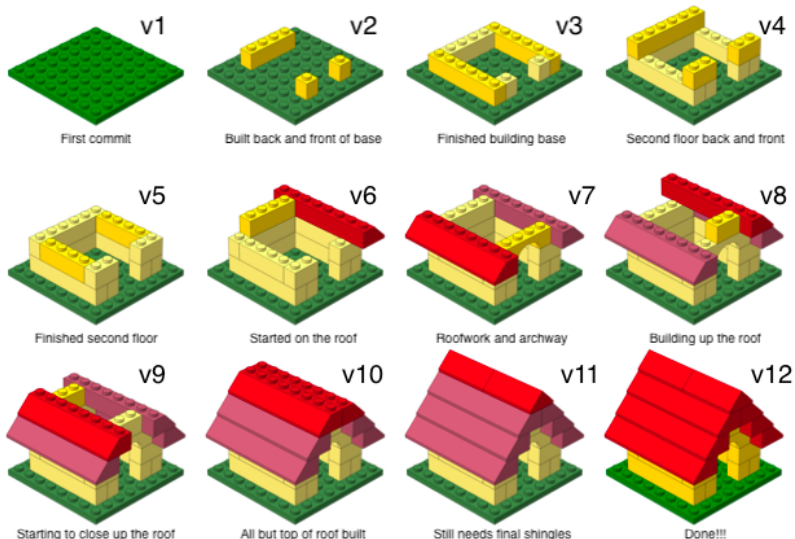
Version control with git and GitHub

What is versioning?



What is versioning?

with human readable messages



Why do we need version control?

Provides a clear record of how the analysis methods evolved. This makes analysis auditable and thus more trustworthy and reliable.
[@ostblom2022]



- ▶ **git** is a version control system – like “Track Changes” features from Microsoft Word.
- ▶ **GitHub** is the home for your git-based projects on the internet (like DropBox but much better).
- ▶ There are a lot of git commands and very few people know them all. 99% of the time you will use git to add, commit, push, and pull.

Multivariable relationships

Carbohydrates in Starbucks food

- ▶ Starbucks often displays the total calories in their food items but not the other nutritional information.
- ▶ Carbohydrates are a body's main fuel source. The Dietary Guidelines for America recommend that carbohydrates make up 45% to 65% of total daily calories.¹
- ▶ Our goal is to understand the relationship between the amount of carbohydrates and calories in Starbucks food items. We'd also like to assess if the relationship differs based on the type of food item (bakery, salad, sandwich, etc.)

¹Source: Mayo Clinic

Starbucks data

- ▶ **Observations:** 77 Starbucks food items
- ▶ **Variables:**
 - ▶ carb: Total carbohydrates (in grams)
 - ▶ calories: Total calories
 - ▶ bakery: 1: bakery food item, 0: other food type

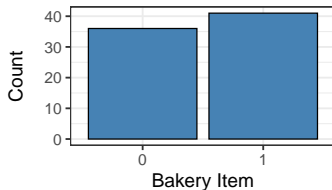
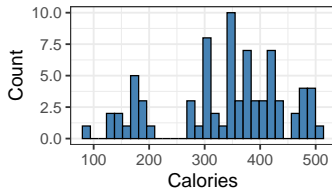
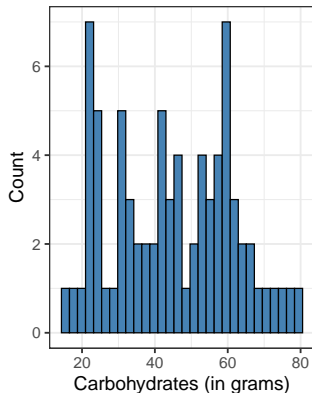
Terminology

- ▶ carb is the **response variable**
 - ▶ variable whose variation we want to understand / variable we wish to predict
 - ▶ also known as *outcome* or *dependent* variable

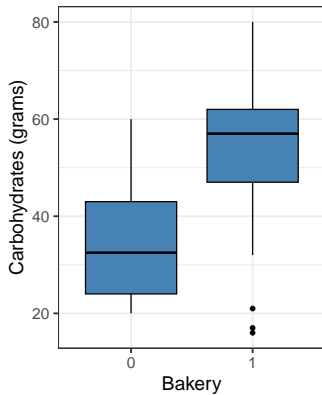
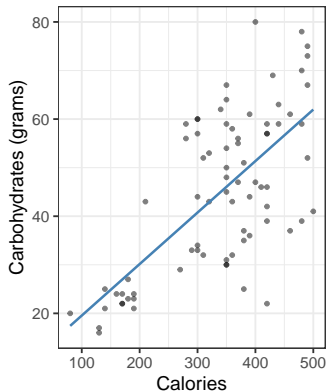
Terminology

- ▶ carb is the **response variable**
 - ▶ variable whose variation we want to understand / variable we wish to predict
 - ▶ also known as *outcome* or *dependent* variable
- ▶ calories, bakery are the **predictor variables**
 - ▶ variables used to account for variation in the response
 - ▶ also known as *explanatory*, *independent*, or *input* variables

Univariate exploratory data analysis



Bivariate exploratory data analysis



Function between response and predictors

$$\text{carb} = f(\text{calories}, \text{bakery}) + \epsilon$$

- ▶ **Goal:** Determine f
- ▶ How do we determine f ?
 - ▶ Make an assumption about the functional form f (parametric model)
 - ▶ Use the data to fit a model based on that form

Determine f

- 1) Choose the functional form of f , i.e., **choose the appropriate model given the response variable**
 - ▶ Suppose f takes the form of a linear model

$$y = f(\mathbf{X}) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon$$

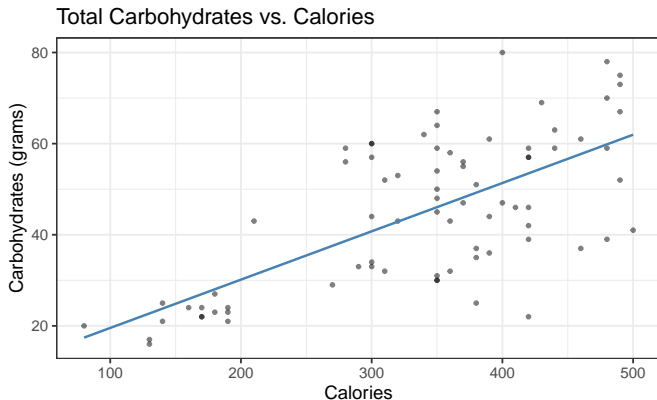
Determine f

- 1) Choose the functional form of f , i.e., **choose the appropriate model given the response variable**
 - ▶ Suppose f takes the form of a linear model

$$y = f(\mathbf{X}) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon$$

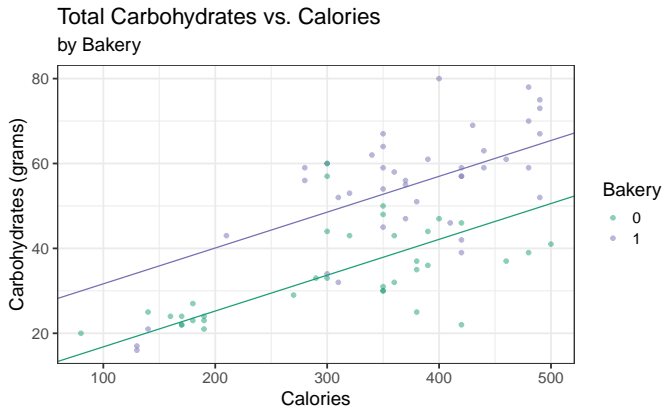
2. Use the data to fit (or train) the model, i.e, **estimate the model parameters**, $\beta_0, \beta_1, \dots, \beta_p$

Carb vs. Calories



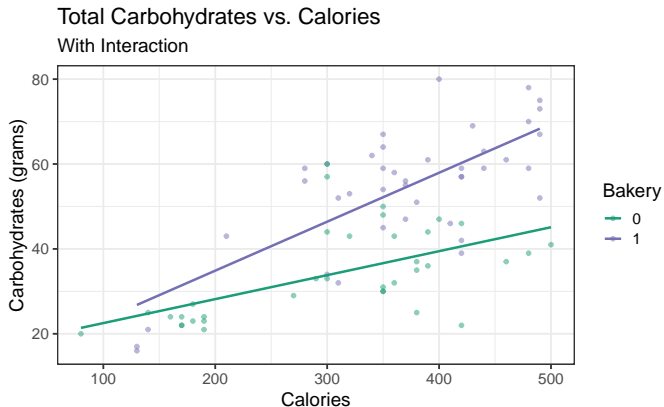
$$\text{carb} = \beta_0 + \beta_1 \text{ calories} + \epsilon$$

Carb vs. Calories + Bakery



$$\text{carb} = \beta_0 + \beta_1 \text{ calories} + \beta_2 \text{ bakery} + \epsilon$$

Carb vs. Calories + Bakery (with interaction)



$$\text{carb} = \beta_0 + \beta_1 \text{ calories} + \beta_2 \text{ bakery} + \beta_3 \text{ calories} \times \text{bakery} + \epsilon$$

Statistical model vs. regression equation

Statistical model (also known as data-generating model)

$$\text{carb} = \beta_0 + \beta_1 \text{ calories} + \beta_2 \text{ bakery} + \beta_3 \text{ calories} \times \text{bakery} + \epsilon$$

Models the process for generating values of the response in the population (function + error)

Regression equation

Estimate of the function using the sample data

$$\hat{\text{carb}} = \hat{\beta}_0 + \hat{\beta}_1 \text{ calories} + \hat{\beta}_2 \text{ bakery} + \hat{\beta}_3 \text{ calories} \times \text{bakery}$$

Why fit a model?

- ▶ **Prediction:** Expected value of the response variable for given values of the predictor variables
- ▶ **Inference:** Conclusion about the relationship between the response and predictor variables
- ▶ What is an example of a **prediction** question that can be answered using the model of carb vs. calories and bakery?
- ▶ What is an example of an **inference** question that can be answered using the model of carb vs. calories and bakery?

Recap

▶ **Reproducibility**

- ▶ It is best practice conduct all data analysis in a reproducible way
- ▶ We will implement a reproducible workflow using R, Quarto, and git/GitHub

▶ **Multivariable relationships**

- ▶ We can use exploratory data analysis to describe the relationship between two variables
- ▶ We make an assumption about the relationship between variables when doing linear regression
- ▶ The two main objectives for fitting a linear regression model are (1) prediction and (2) inference

References