

# Introduction Data Management and

Kylie A. Bem

Northeastern University  
Khoury College of Computer Science



Northeastern University

# Basics

- DS 5110 Introduction to Data M
- Sec 3:Tue 11:45 am - 1:25 pm,Th
- Sec 4:Thu 6:00 pm - 9:20 pm
- Dates: Jan 8, 2023 - Apr 27, 2023
- Course website:
  - ◆ <https://kuwisdelu.github.io/ds5110-spring24/>
  - ◆ Link to all relevant resources available on C

# Intros: Kylie Arie

- PhD Statistics, Purdue University
- Research includes:
  - ◆ Bioinformatics (mass spectrometry imaging)
  - ◆ Statistical computing languages (R development)
- Software (R packages):
  - ◆ “Cardinal” - mass spectrometry imaging tools
  - ◆ “matter” - out-of-memory computing
- 2015 John M. Chambers Statistical Software award for development of “Cardinal” package
- Published author of short science fiction and fantasy

## Intros: TA

*"Keep an open and curious mindset"*



Jayatha

*"Data science is about more than co*



Sakshi

*"Engage in class and*



Shubham

*"Embrace challenges with curiosity and perseverance"*



## Goals for today

- Course expectations and responsibilities
- Tools and resources
- Introduction to data science
- Introduction to R

# COURSE EXPECTATIONS AND POLICIES

## Course materials

- Syllabus:
  - ◆ <https://kuwisdelu.github.io/ds5110/>
- Piazza:
  - ◆ <https://piazza.com/northeastern/spr2020/ds5110>
- Canvas:
  - ◆ <https://northeastern.instructure.com/courses/10327>

## Zoom meetings

- Classes will be hosted synchronously
  - ◆ Meeting links for each class can be found in the course calendar
  - ◆ Meetings will be recorded and made available for review
- Synchronous attendance is encouraged
  - ◆ Projects must be presented during the scheduled meeting time
  - ◆ Quizzes will have a limited time frame to complete

## Announcements, questions

- Outside office hours, all course-related correspondence should use **Piazza**
  - ◆ <https://piazza.com/northeastern/spring2024>
  - ◆ Course material and office hours schedules
  - ◆ Use “Q&A” to ask asynchronous questions
  - ◆ Send a private note for course-related questions
- Do not email for course-related questions
- How to ask a good question:
  - ◆ <https://stackoverflow.com/help/how-to-ask>

## Office hours

- Virtual office hours will be held:
  - ◆ Download at <https://teams.northeastern.edu>
  - ◆ Live online during our scheduled office hours
  - ◆ Chat or video call
- Or schedule a video call by appointment
- In-person appointments available

## Assignments and

- Assignments and grading will use Canvas:
  - ◆ <https://northeastern.instructure.com>
  - ◆ Semester schedule (subject to change)
  - ◆ Assignments must be submitted via Canvas
  - ◆ Grades will be distributed via Canvas
- Do not email assignments — the

## Grade breakdown

- Homework: 40%
- Quizzes: 30%
- Project: 30%

## Homework

- 5-6 individual homework assignments
- Practice on real-world datasets
- Submit on Canvas
  - ◆ R Markdown source code
  - ◆ PDF report (generated from R)
- Due every ~1-2 weeks

## Quizzes

- 2-3 cumulative quizzes
- Complete online during class
- Test your conceptual understanding
- Will replace a class meeting

# Project

- Teams of 2-5 students
- Complete a term project using the following steps:
- Choose a team and submit a proposal
- Present in-class during final week of term
- Submit a final written report
- Detailed guidelines will be discussed in class

## Peer review

- Provide feedback for your peers
- Receive feedback on your work
- Due after projects are submitted
- Rubrics will be shared on Canvas and discussed in class

## Grading policy

- Late assignments will not receive credit
  - ◆ Assignments late by more than 1 day will not be graded.
  - ◆ Submit early to avoid technical issues or loss of work.
- Extensions may be given on assignments
  - ◆ Request at least 48 hours in advance of the due date.
  - ◆ Require a reasonable justification.
- Petitions for re-grades must be submitted in writing
  - ◆ Must be submitted within 1 week and 1 day of the grade being posted.
  - ◆ The new grade may be lower than the original grade.

## General policies

- Academic integrity policy
  - ◆ <http://www.northeastern.edu/osccr/academic-integrity-policy>
  - ◆ Please remember you are here to learn, not just work.
  - ◆ Use work with permission; cite it properly; don't plagiarize.
- Title IX policy
  - ◆ <https://www.northeastern.edu/ouec/>
  - ◆ Northeastern is committed to preventing discrimination.
  - ◆ You may talk to me, but I am a mandatory reporter.
  - ◆ Confidential resources are available if you need them.

## General policies

- Please be kind and respectful to
  - ◆ Use other people's requested names and pronouns
  - ◆ Be considerate of other people's identities and backgrounds
  - ◆ Keep the classroom a safe and healthy environment
- Getting help
  - ◆ Please reach out as early as possible if you are experiencing challenges
  - ◆ I am less able to make accommodations close to the start of the semester
  - ◆ Northeastern's WeCare office is a resource
    - <https://studentlife.northeastern.edu/we-care/>

# TOOLS AND RE

## Software

- Install R from CRAN
  - ◆ <https://cran.r-project.org>
- Install RStudio
  - ◆ <https://rstudio.com/products/rstudio>
- Install LaTeX
  - ◆ <https://www.latex-project.org/guides>



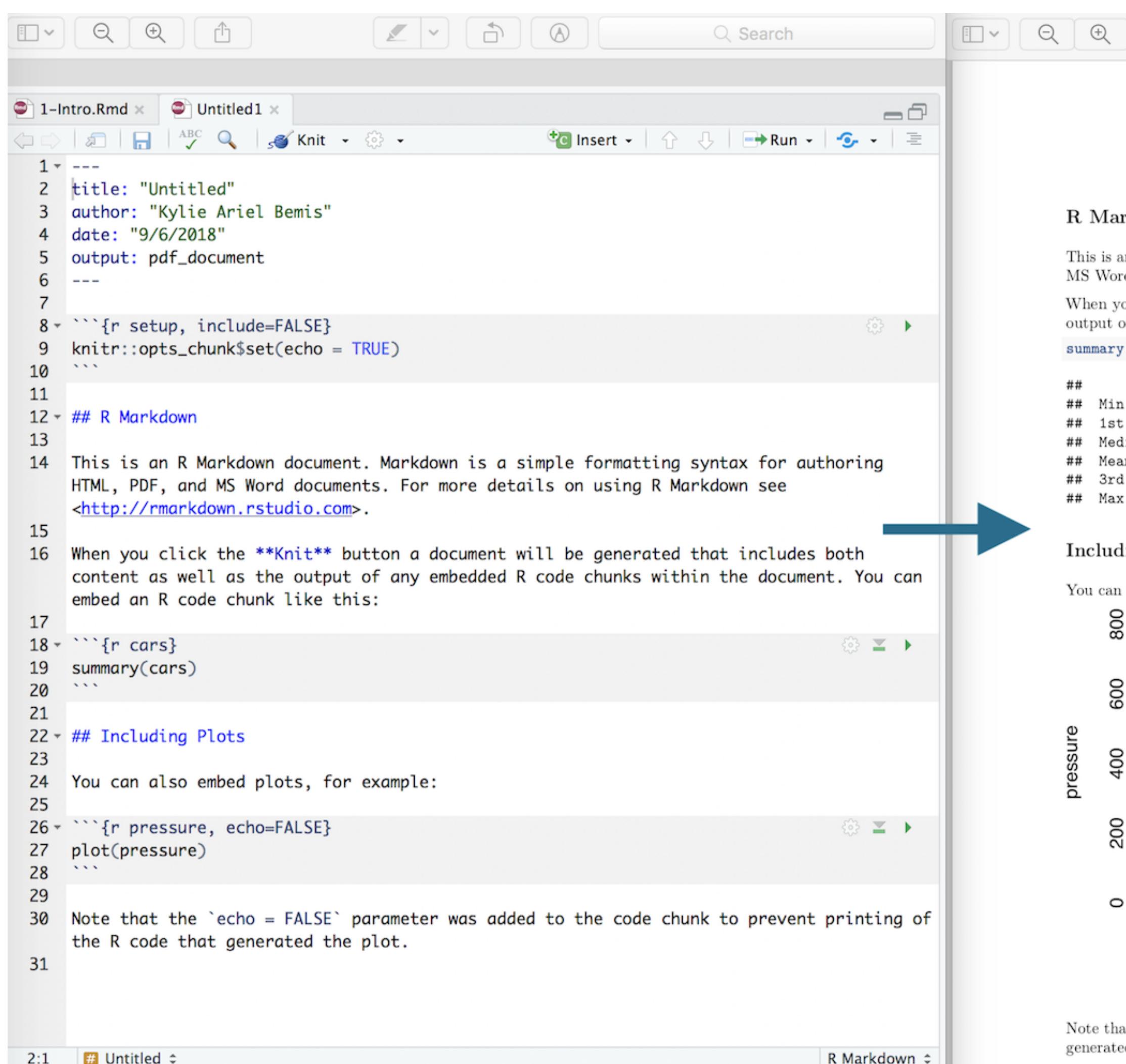
# R packages

- Tidyverse
  - ◆ ggplot2, dplyr, etc.
  - ◆ `install.packages("tidyverse")`
- Datasets
  - ◆ `install.packages("nycflights13")`
  - ◆ `install.packages("gapminder")`
- R Markdown
  - ◆ `install.packages("rmarkdown")`

## R Markdown

- Authoring framework for data analysis
- Combines R programming with writing
- Code for analysis embedded in text
- Documents are fully reproducible
- Required for homework submissions
  - ◆ Rmd — R + Markdown source code
  - ◆ PDF — polished, presentable report

# R Markdown



The screenshot shows the RStudio interface with the 'Untitled1' tab active. The code editor contains the following R Markdown document:

```
1 ---  
2 title: "Untitled"  
3 author: "Kylie Ariel Bemis"  
4 date: "9/6/2018"  
5 output: pdf_document  
6 ---  
7  
8 ```{r setup, include=FALSE}  
9 knitr:::opts_chunk$set(echo = TRUE)  
10 ```  
11  
12 ## R Markdown  
13  
14 This is an R Markdown document. Markdown is a simple formatting syntax for authoring  
HTML, PDF, and MS Word documents. For more details on using R Markdown see  
http://rmarkdown.rstudio.com.  
15  
16 When you click the **Knit** button a document will be generated that includes both  
content as well as the output of any embedded R code chunks within the document. You can  
embed an R code chunk like this:  
17  
18 ```{r cars}  
19 summary(cars)  
20 ```  
21  
22 ## Including Plots  
23  
24 You can also embed plots, for example:  
25  
26 ```{r pressure, echo=FALSE}  
27 plot(pressure)  
28 ```  
29  
30 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of  
the R code that generated the plot.  
31
```

A large blue arrow points from the text 'You can also embed plots, for example:' towards the right side of the screen, where a small plot of 'pressure' data is visible.

## Reading and resources

- *R for Data Science* by Wickham
  - ◆ <http://r4ds.had.co.nz/>
- *Advanced R* by Wickham
  - ◆ <http://adv-r.had.co.nz/>
- *Introduction to Statistical Learning* by Friedman, Hastie, Tibshirani
  - ◆ <https://hastie.su.domains/ISLR2/ISLRv2.pdf>
- *Text Mining with R* by Silge and Robinson
  - ◆ <https://www.tidytextmining.com/>

## How to ask for help

- Write titles that summarize your problem
  - ◆ Bad: “Doubt on homework”
  - ◆ Good: “Why does my function return None?”
- Explain the problem
  - ◆ What is the problem you’re trying to solve?
  - ◆ What have you tried so far to fix it?
- <https://stackoverflow.com>

## How to ask for help

- Provide a reproducible example
  - ◆ Provide the minimum amount of code necessary
  - ◆ Make sure your example reproduces the error
  - ◆ Ideally, the example should run on my machine
  - ◆ Provide the actual code, and not just a screenshot
- <https://stackoverflow.com/questions/ask#reproducible-example>

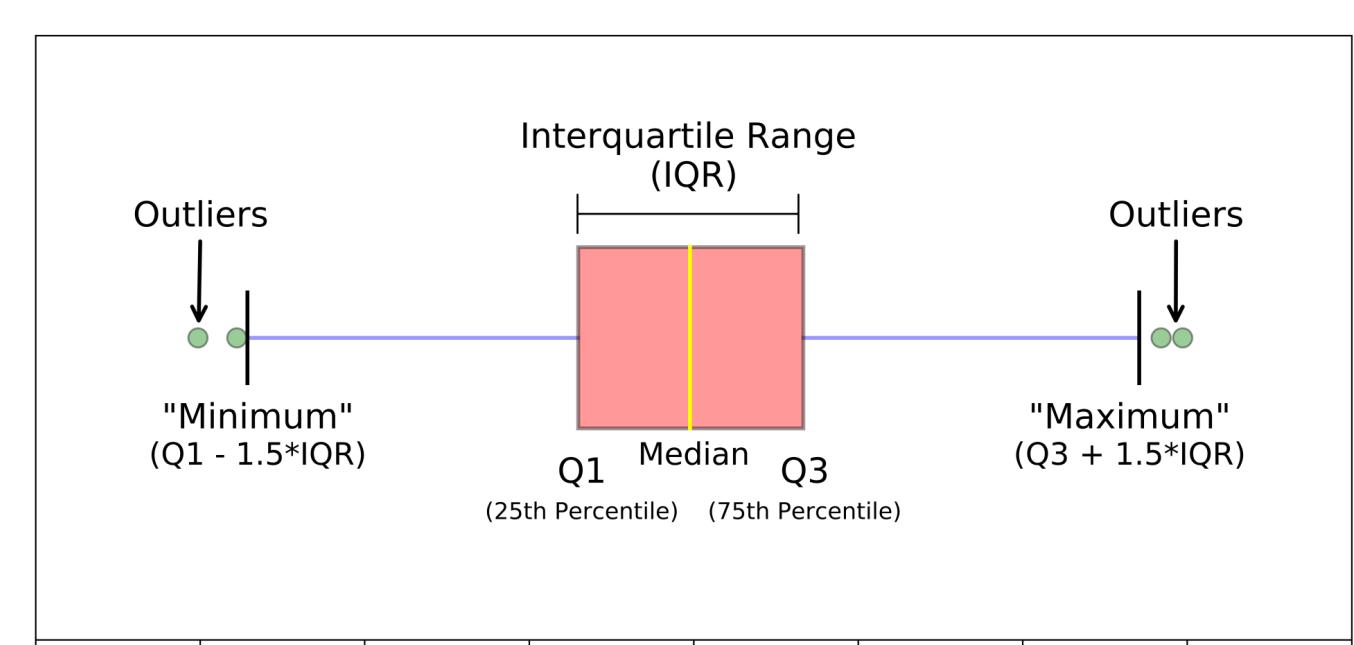
# INTRODUC TO DATA SC

## Learning goals in this module

- What is data science?
- What does a data scientist do?
- What you will learn in this module

# Where did data science come from?

- 1935 — R.A. Fisher's *Design of Experiments*
  - ◆ “Correlation does not imply causation”
- 1974 — Peter Naur coins term “data processing”
  - ◆ Also proposed “datalogy” to replace “computer science”
- 1977 — John Tukey's *Exploratory Data Analysis*
  - ◆ Introduced the box-and-whisker plot and popularized the term “outlier”



- 1977 — The International Association for Statistical Computing is established
- 1989 — The first Knowledge Discovery in Databases (KDD) workshop is held
- 2001 — William S. Cleveland publishes “Data Science: An Action Plan for Expanding the Technical Areas of the Field of Statistics”
  - ◆ First proposes “data science” as a separate discipline distinct from statistics or computer science

- 2002 — *Data Science Journal* launch
- 2003 — *Journal of Data Science* launch
- 2005 - 2014
  - ◆ Many people write articles that struggle to precisely define “data science” but everyone seems to agree it’s important
  - ◆ “Sexy” is a frequent descriptor of the new discipline
- 2015 - present
  - ◆ Northeastern launches Masters in Data Science

## Defining data s

- Drew Conway writes in “The Data Science Venn Diagram” (2010):
  - ◆ “...one needs to learn a lot as they aspire to become a fully competent data scientist. Unfortunately, simply enumerating texts and tutorials does not untangle knots. Therefore... I present the Data Science Venn diagram... hacking skills, math and stats knowledge, and substantive expertise.”

## What is data science?

- A **data scientist** blends several fields of expertise to draw insights
  - ◆ Statistics and machine learning
  - ◆ Computer science and programming
  - ◆ Domain knowledge
- A **data scientist** leverages data to make inferences that drive decisions in a particular field of application

## The data science process

- Understand the problem
  - ◆ What is the purpose of this analysis?
  - ◆ May require domain knowledge
- What data and tools are needed
  - ◆ May need to pull from a database
  - ◆ May need to identify appropriate sources
  - ◆ Most data needs to be cleaned
  - ◆ What additional resources do you need?

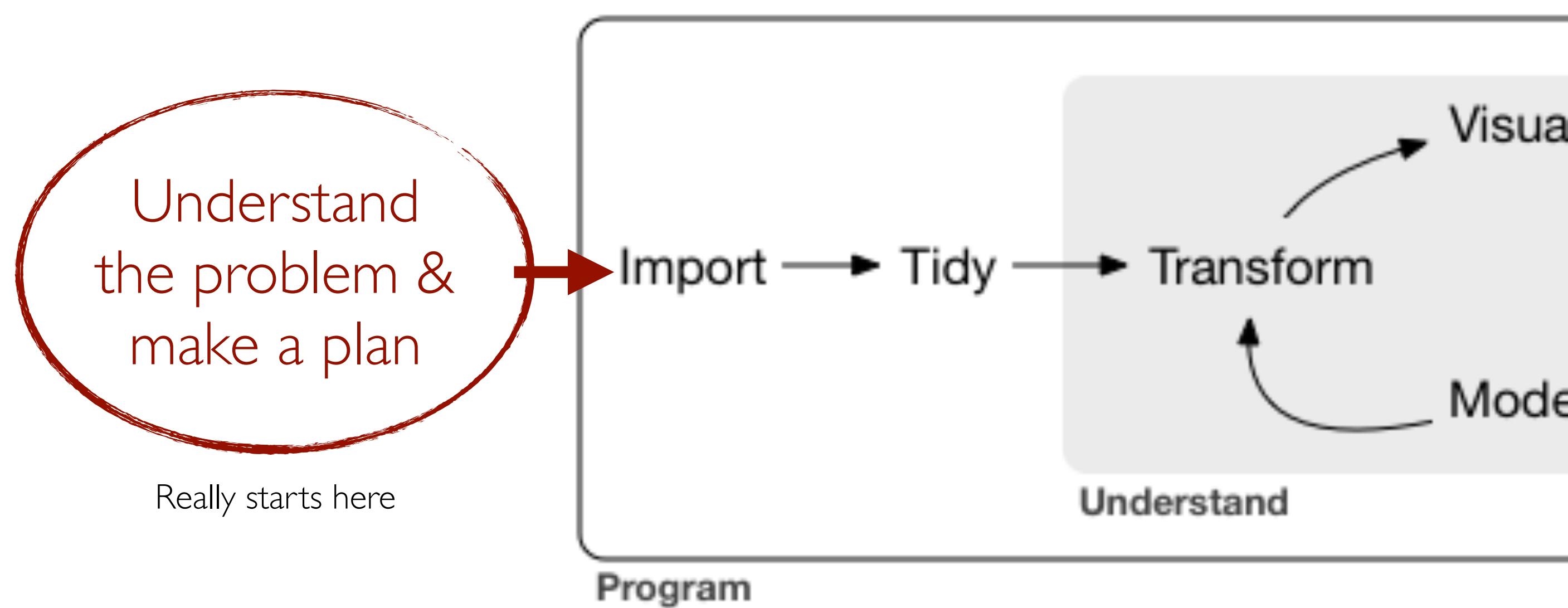
## The data science process

- Perform the analysis
  - ◆ What kind of analysis is appropriate?
    - Regression or classification (supervised)
    - Clustering or dimension reduction (unsupervised)
    - Hypothesis testing? Statistical inference?
  - ◆ Visualize the cleaned data and results
- Communicate the results
  - ◆ What insights or decision-making recommendations can be derived?
  - ◆ Communicate the results in an accessible way
  - ◆ Do you need to deliver a product?

# Data science workflow

Data science is an iterative process that

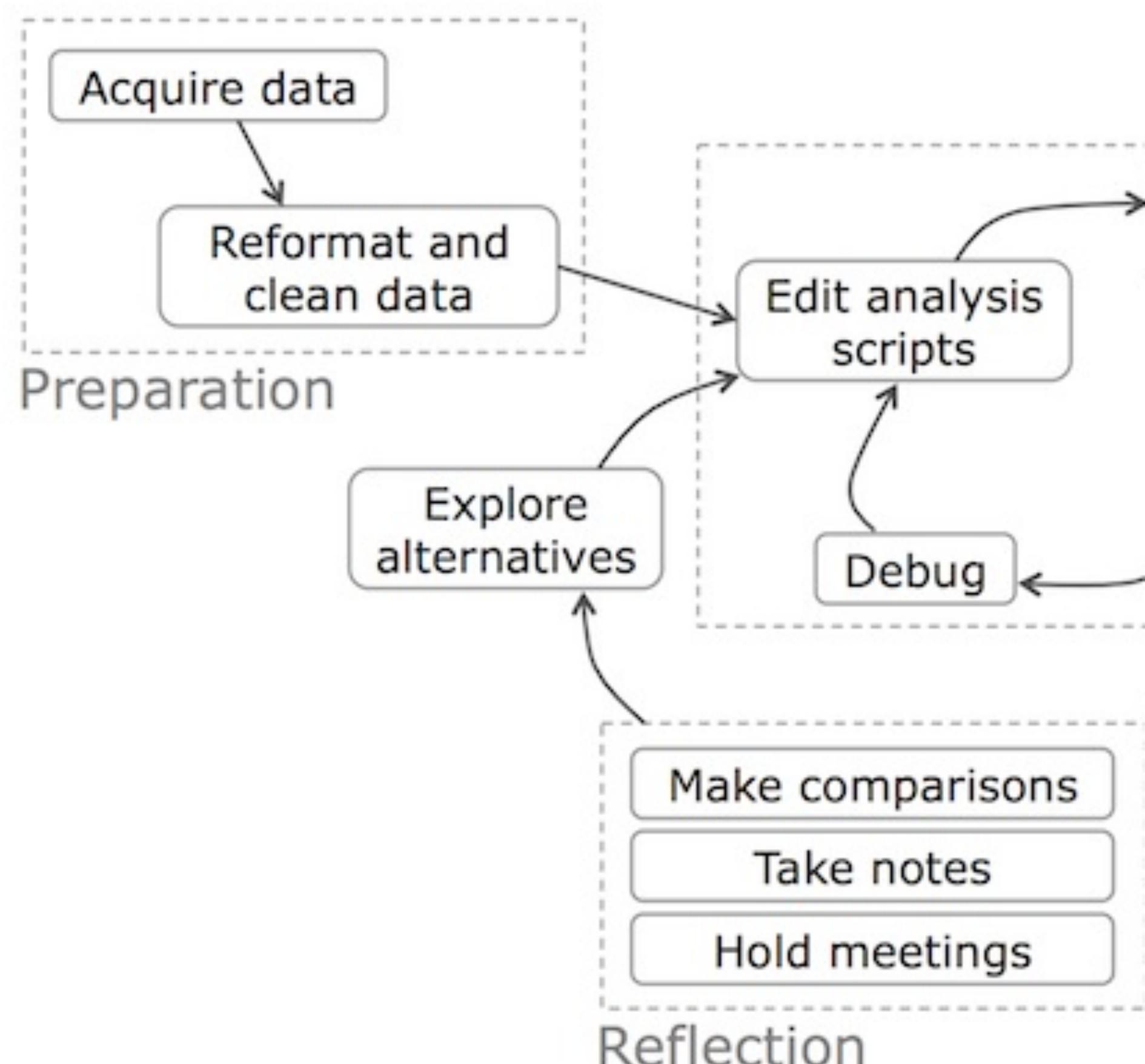
Apply hacking skills with domain knowledge



An example of the data science workflow is shown in the diagram above, taken from *R for Data Science* by Wickham.

# Data science workflow

Another perspective



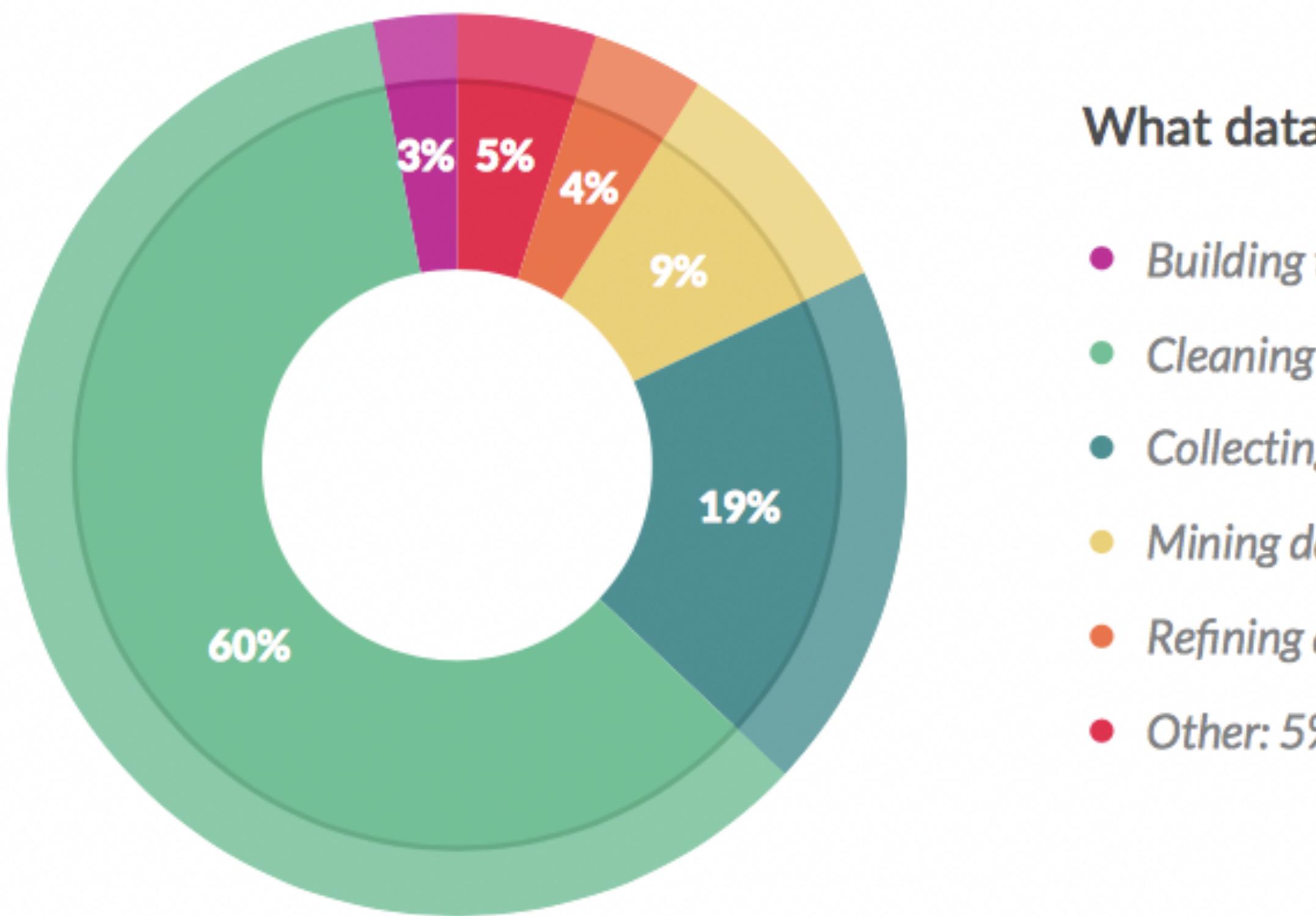
<https://cacm.acm.org/blogs/blog-cacm/169199-data-science-workflow>

## What does a data scientist do?

- visualize data
- pull from a database
- scrape from a website
- clean and wrangle data
- run a regression
- build and train a model
- design an experiment

You should be able to do ~80% of this

## What does a data scientist do?



<https://visit.figure-eight.com/rs/416-ZBE-142/images/Crowdsource%20-%20Data%20Science%20-%20Industry%20Task%20Allocation%20-%20Donut%20Chart%20-%20Large%20-%201000x600px.pdf>

## What you will (and won't)

- This course will discuss the fundamental concepts of data management, processing, and exploration:
  - ◆ How to clean and organize data for analysis
  - ◆ How to visualize and perform basic data analysis
  - ◆ How to interpret results and communicate findings
- This course will not delve into:
  - ◆ theory, database design, or cloud computing
  - ◆ But we will use machine learning tools as part of our analysis

# Course outline

- Data visualization
  - ◆ Looking at data and understanding it
- Data wrangling
  - ◆ Cleaning and organizing data for analysis
- Relational data and databases
  - ◆ How data is stored, accessed, and manipulated
- Building models for data analysis
  - ◆ Basics of statistical modeling and machine learning
- Special topics (TBD)

# INTRODUCTION