

# Strategic and Competitive Intelligence

Data Science Applications in R



**Filippo Chiarello**  
[filippochiarello.90@gmail.com](mailto:filippochiarello.90@gmail.com)

# About me

---



- FILIPPO CHIARELLO
- PHD STUDENT IN MANAGEMENT ENGINEERING
- RESEARCH ON DATA SCIENCE AND TEXT MINING
- MUSICIAN



# Goal of this Module

---

**Data → Knowledge**

# Some Magic...

---



# ... lot of problems

---

- Start solving problems as soon as possible (today)
- Learn the 20% of tools that solves the 80% problems
- Don't solve the same problem multiple times
- Don't solve problem that someone else had solved

# Lessons

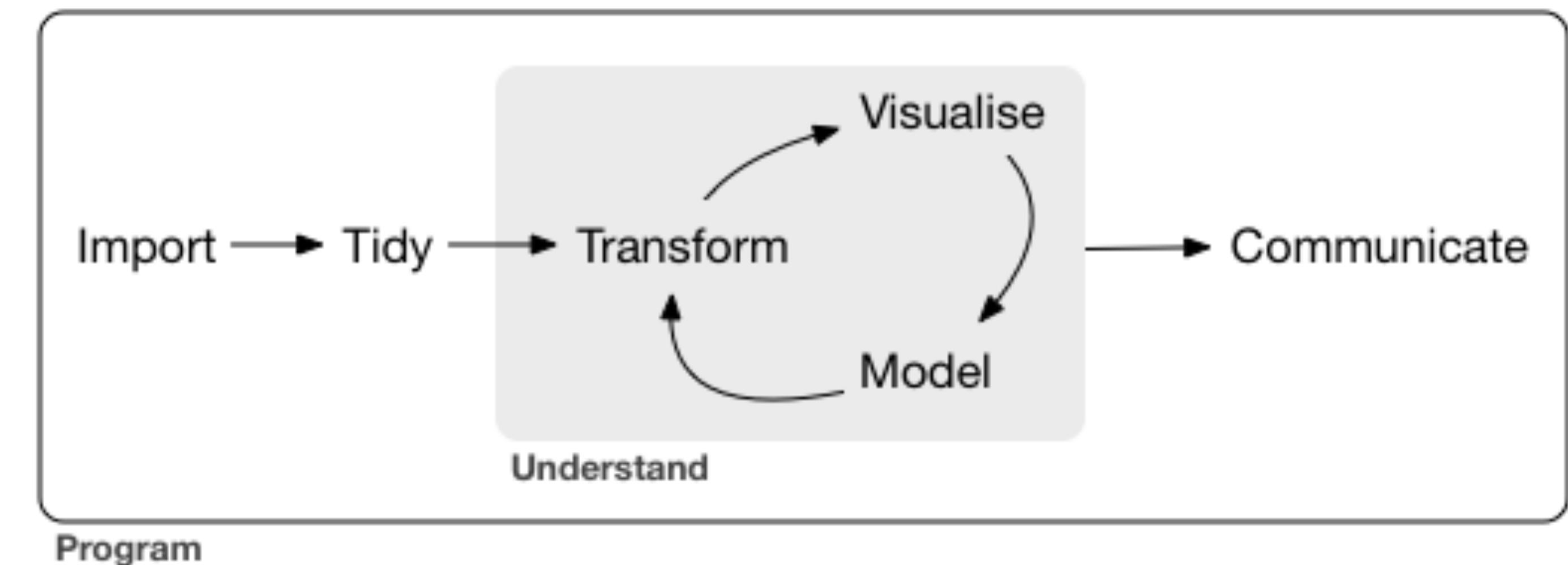
---

1. Introduction to R
2. Data Visualisation with ggplot2
3. Data Transformation with dplyr
4. Theory of Models
5. Text Mining
6. Social Networks Analysis
7. Patents Analysis
8. Papers Analysis
9. Results Communication and Reporting with Rmarkdown

# Mantras

---

# Make Models



# Mantras

---

Ask Why

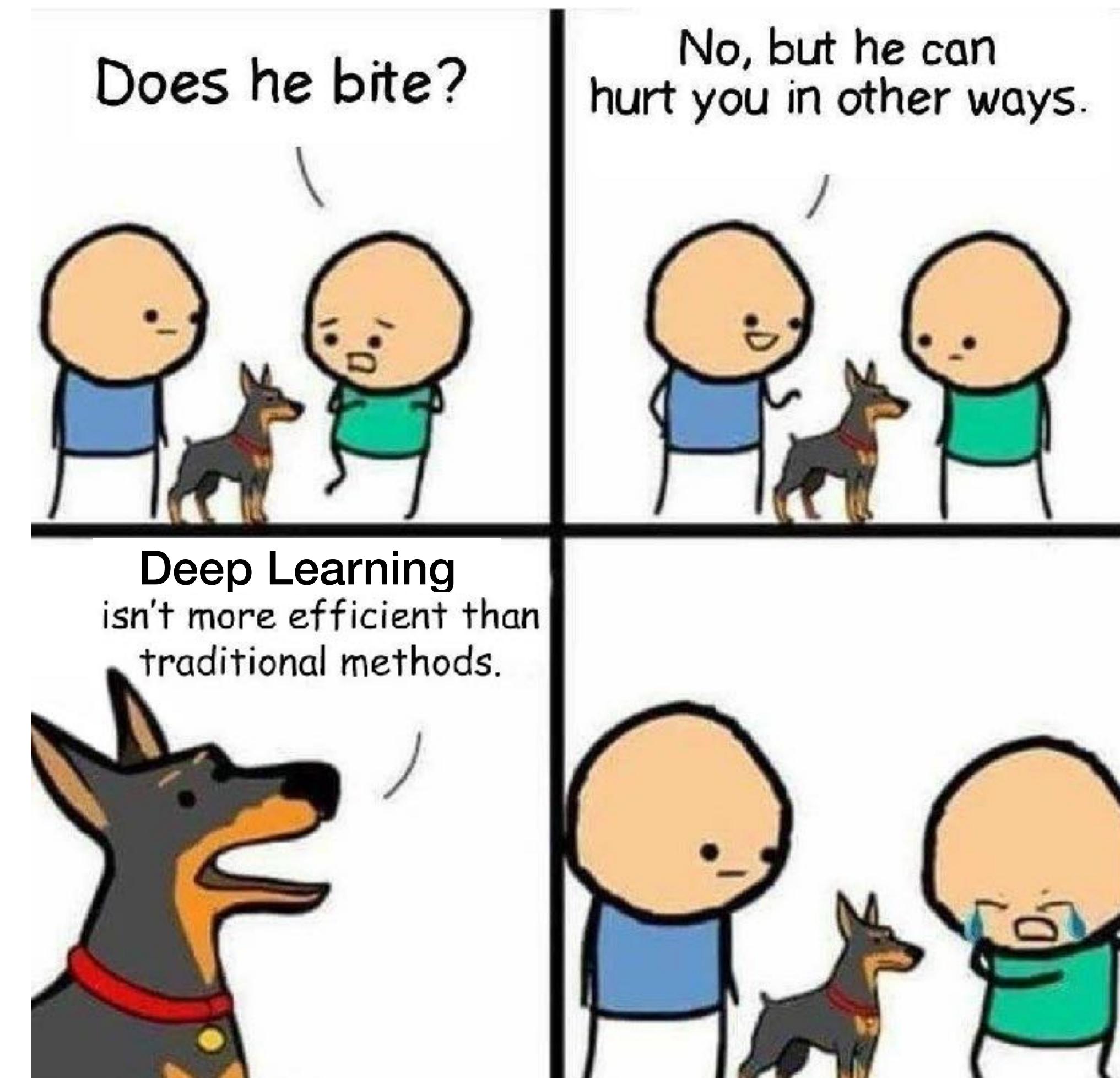


# Mantras

---

**Be Sceptic**

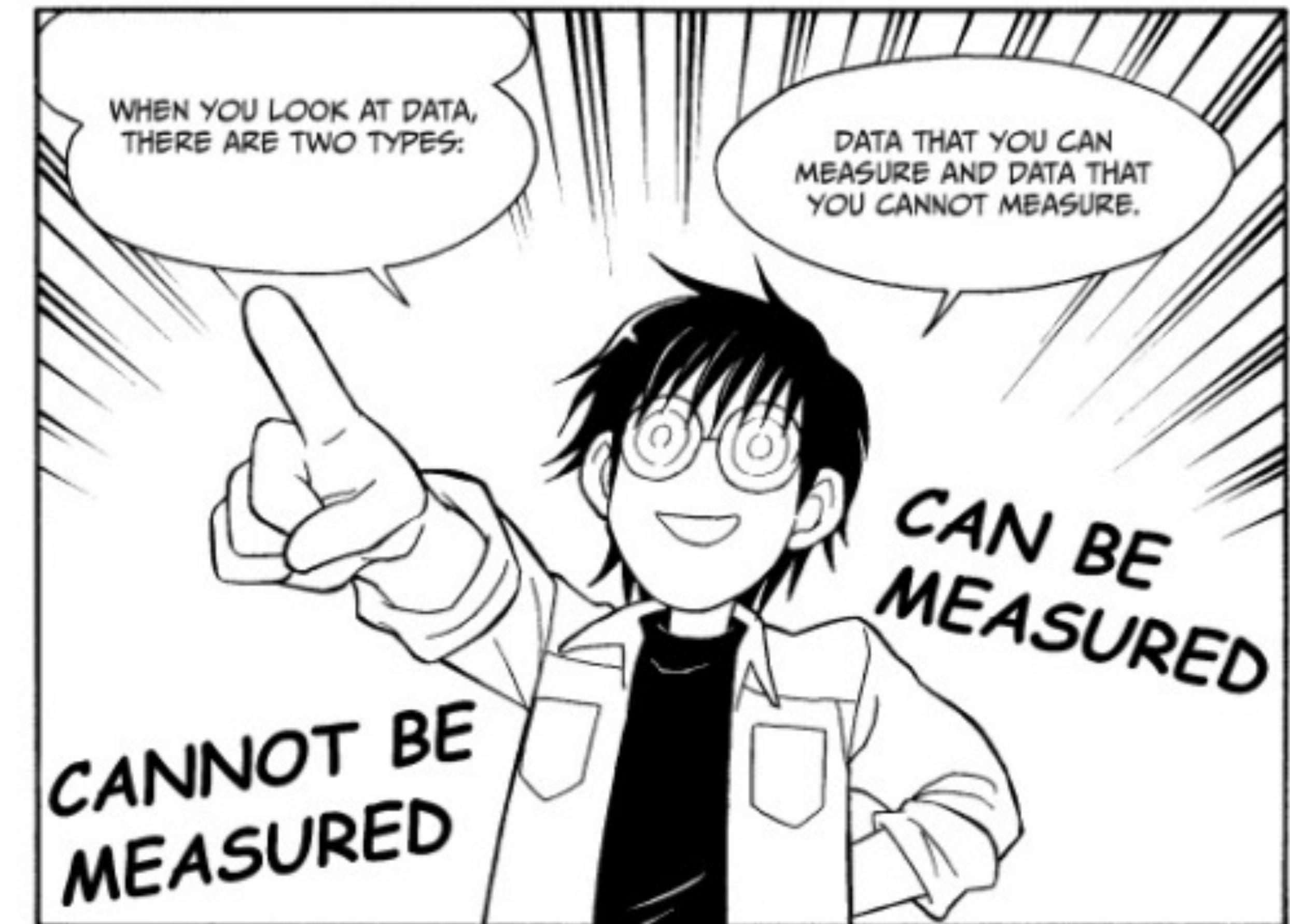
**Be Agnostic**



# Mantras

---

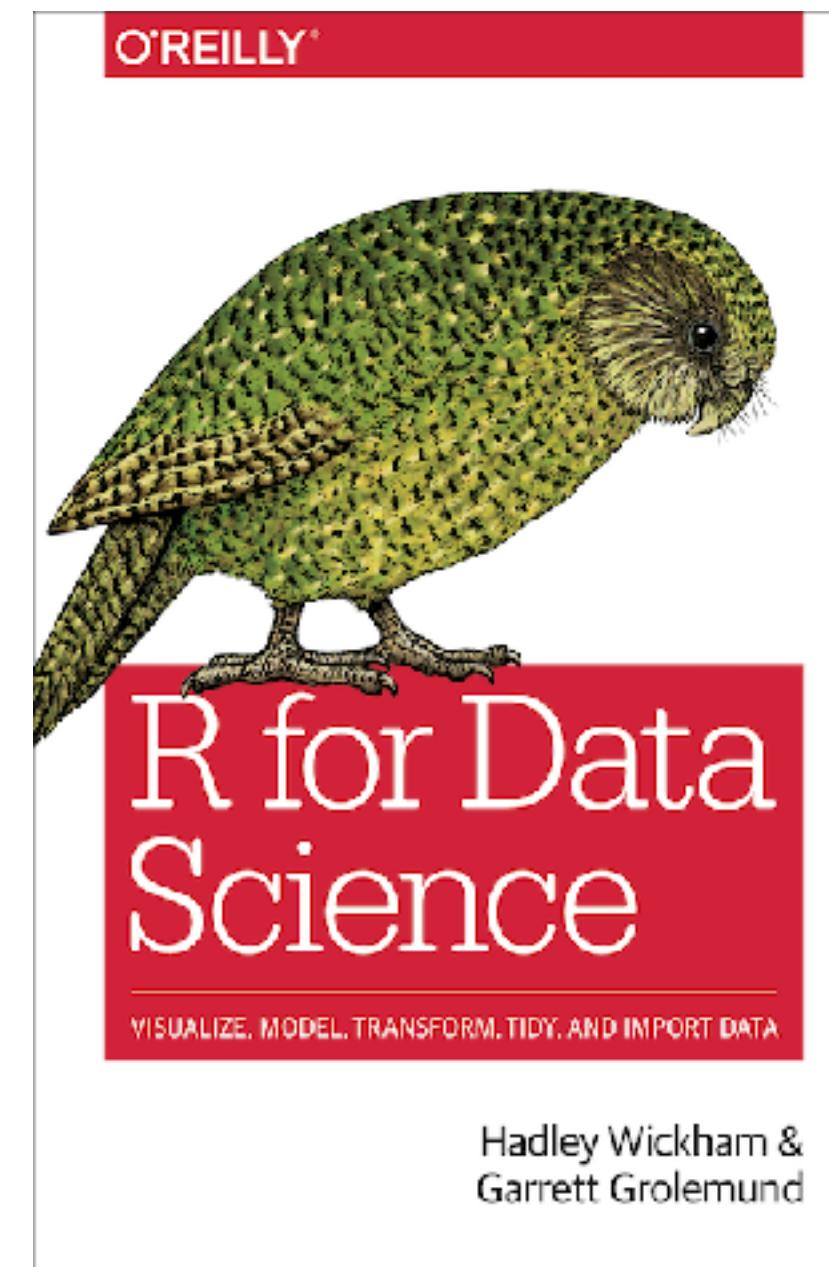
Measure it



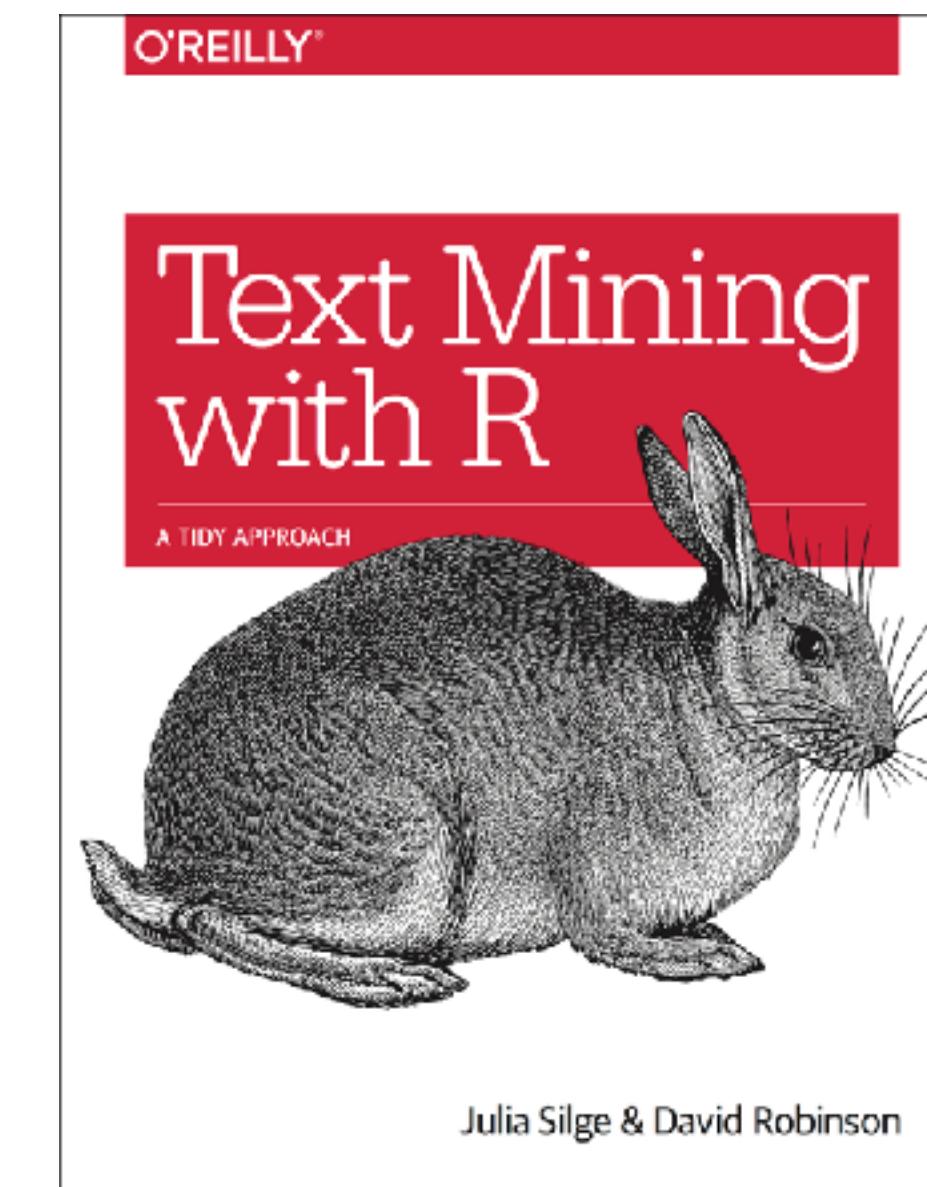
*From "Manga Guide to Statistics", Shin Takahashi, 2008*

# Material

---



<http://r4ds.had.co.nz/>



<https://www.tidytextmining.com/>

# Material

[Home](#)   [Posts](#)   [Courses Material](#)   [Publications](#)   [Contact](#)



## Filippo Chiarello

PhD Student

University of Pisa, DESTEC



## About Me

I study how to use Data Science tools for Innovation Management.

In this blog you can find:

- **Posts:** Projects i'm working on, books that i love, cool blog posts from other researchers, mind-blowing talks and so on.
- **Courses Materials:** News and materials for the my courses.
- **Publications:** My latest research activities.

## Interests

- Data Science
- R
- Engineering Design
- Natural Language Processing
- Artificial Intelligence

## Education

- **MEng in Engineering Management, 2015**  
University of Pisa
- **BSc in Engineering Management, 2013**  
University of Pisa

# Material

The screenshot shows a Twitter profile page for a user named Filippo Chiarello. The profile picture is a black and white portrait of a man with dark hair and a beard, wearing a suit and tie. The header background is blue. Below the profile picture, the user's name "Filippo Chiarello" and handle "@f\_chiare" are displayed. The bio reads: "I study innovation using datascience. #rstat lover. PhD student at @Unipisa." It also shows location "Pisa, Toscana" and a graduation date "Iscrizione a maggio 2015". There are links to "15 foto e video" and a "STATEMENT OF ACCOMPLISHMENTS" document.

Key statistics on the profile card: Tweet 572, Following 850, Follower 186, Mi piace 1.948, Liste 1, Momenti 0. A "Modifica profilo" button is visible.

The main feed area has tabs for "Tweet", "Tweet e risposte", and "Contenuti". The "Tweet" tab is active, showing a single tweet from Albert Y. Kim (@rudeboybert) about markdown URLs. The "Contenuti" tab displays a screenshot of a RStudio interface showing code for knitting documents.

The sidebar on the right contains sections for "Le tue interazioni Tweet" (with 1.876 impressions), "Visualizza i tuoi Tweet più popolari", "Chi segue" (with links to profiles for Matteo @MatteoCip87 and DataEconometria @Data...), and a link to "Aggiorna - Visualizza tutto".

# Material

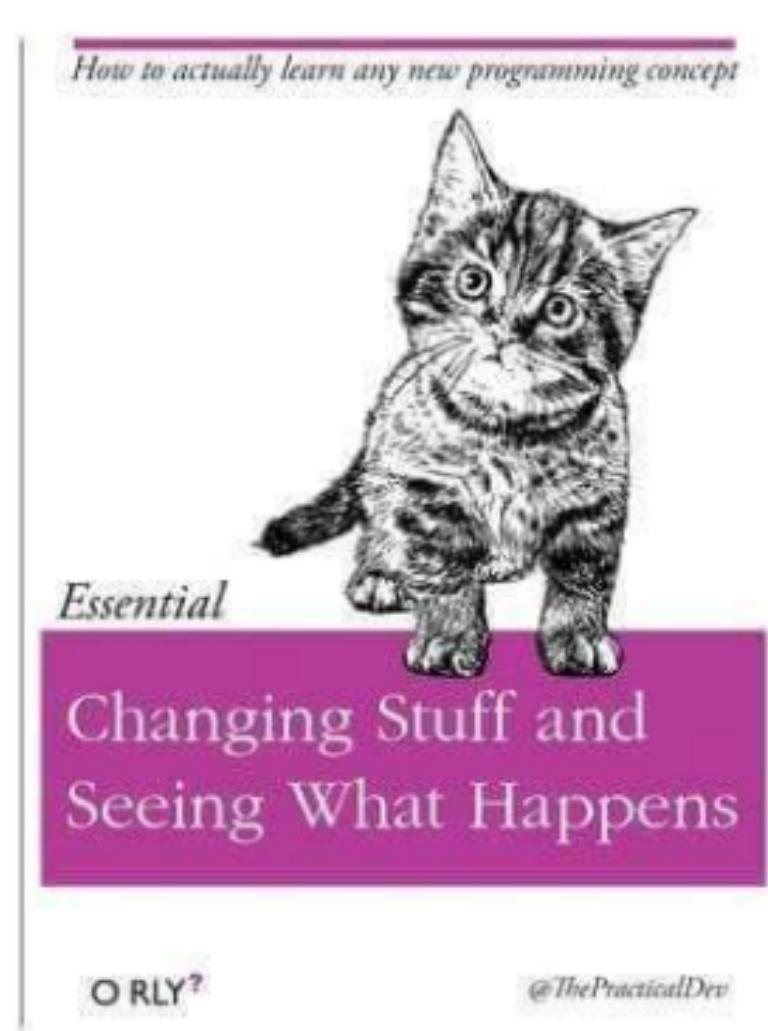
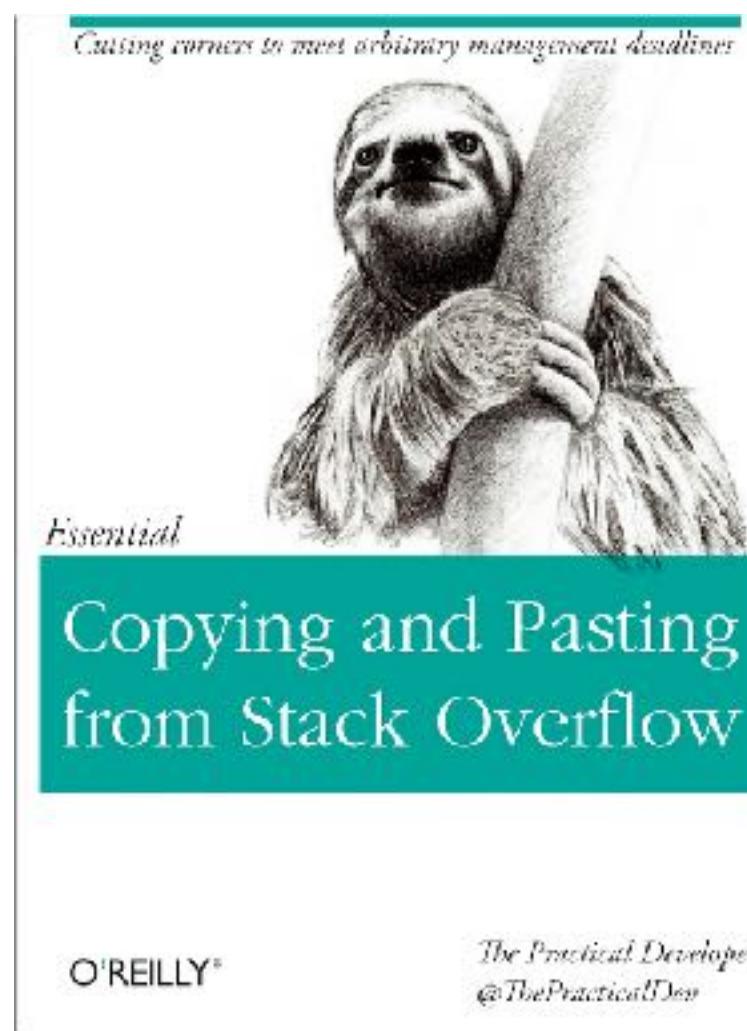
The screenshot shows a Slack interface for a channel named '#general'. The channel has 81 members and 0 messages. The messages are as follows:

- Saturday, September 15th**
  - Antonella Martini** 12:24 PM: WELCOME TO THE NEW COURSE ON STRATEGIC & COMPETITIVE INTELLIGENCE # SCI
- Tuesday, September 18th**
  - Francesco Cucinotta** 7:02 AM: joined #general along with 3 others.
  - Antonella Martini** 10:44 AM: Here you are the link for SCI on e-learn web site -> <https://elearn.ing.unipi.it/course/index.php?categoryid=341>  
Log in with your UniPi personal credential AND use the pass code (DATA\_SCIENCE)  
On Thursday you'll find R4DS file to download for the Friday lesson with **@Filippo Chiarello**
  - Federica Trevisan** 2:40 PM: joined #general along with 3 others.
- Yesterday**
  - Filippo Chiarello** 12:18 PM: Good morning students, in my web site ([filippochiarell.com](http://filippochiarell.com)) you will find the instructions to be ready for tomorrow.  
Please spread the word! 🎉
  - 1 reply** Today at 12:55 PM
- Today**
  - Antonella Martini** 12:47 PM: Dear students, would you please introduce yourself briefly (just ads done in the kick off)?

The sidebar on the left shows the channel list, including '# general' (selected), '# random', 'Direct Messages', and 'Apps'.

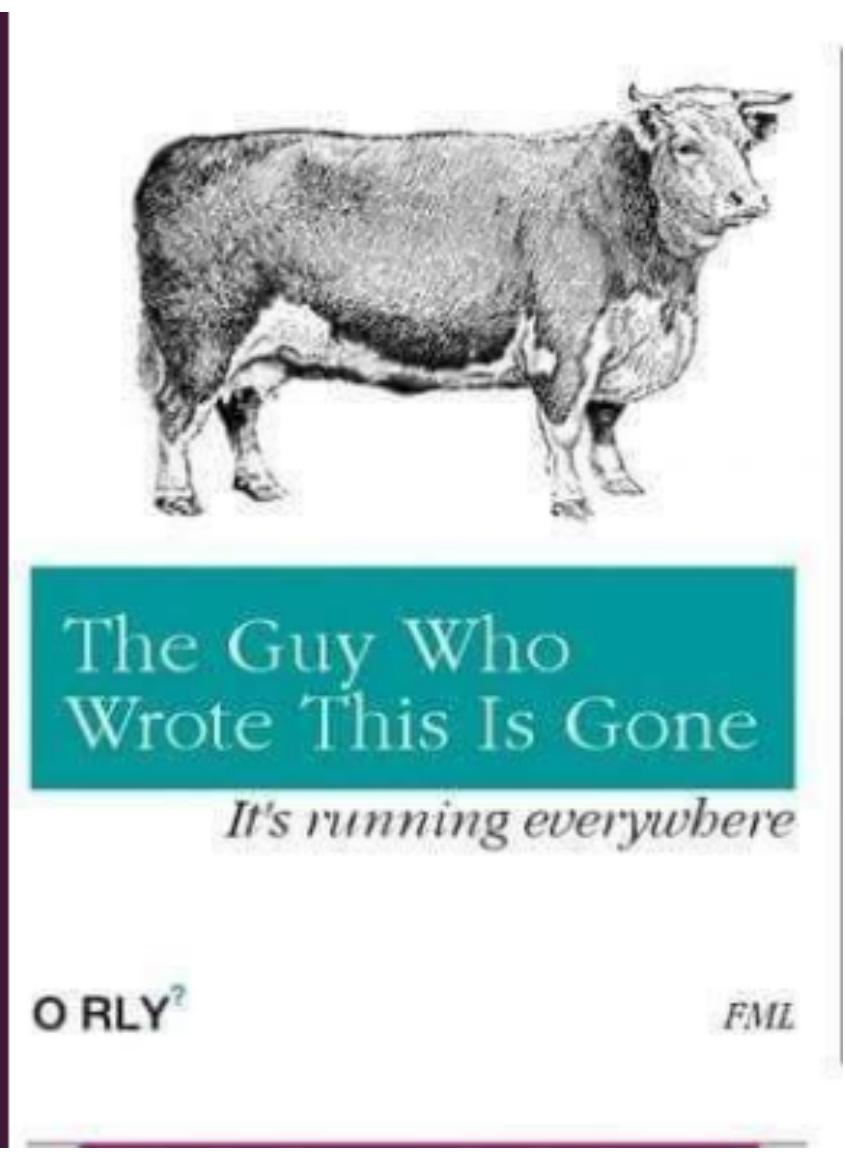
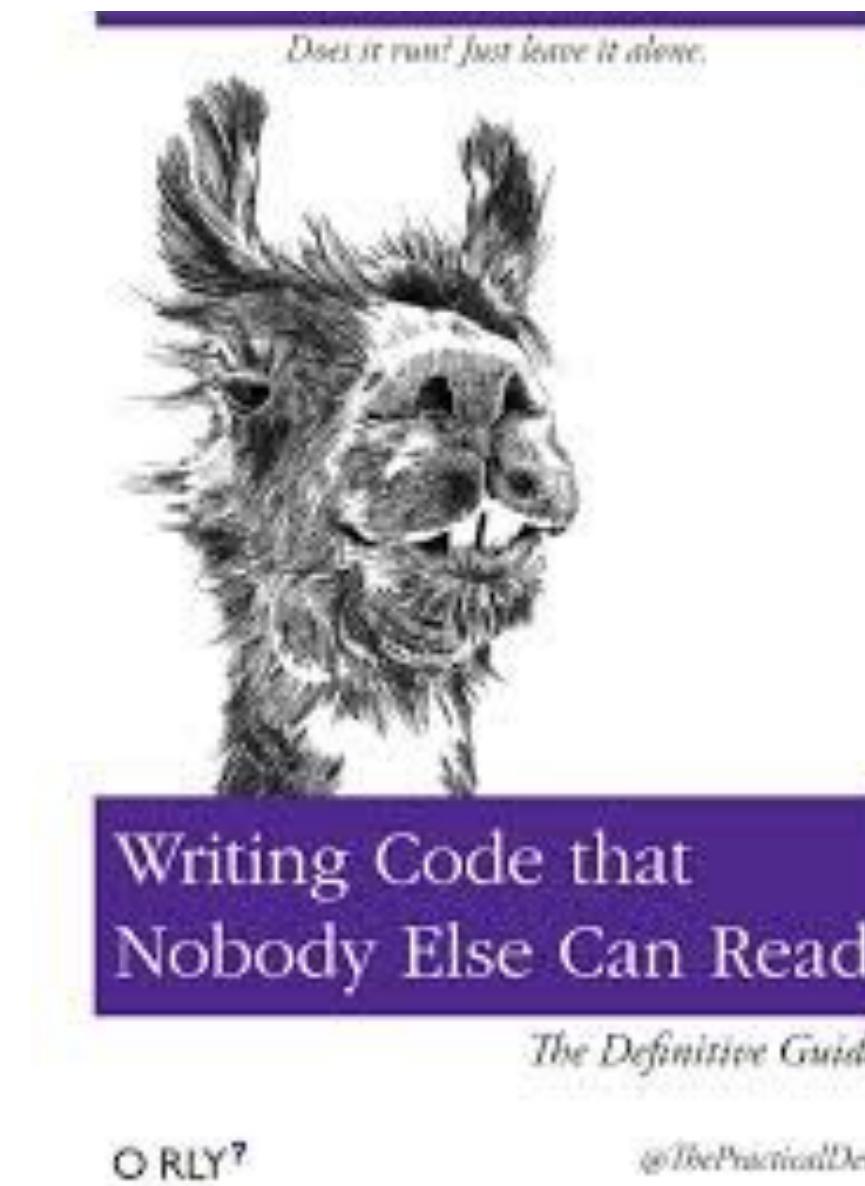
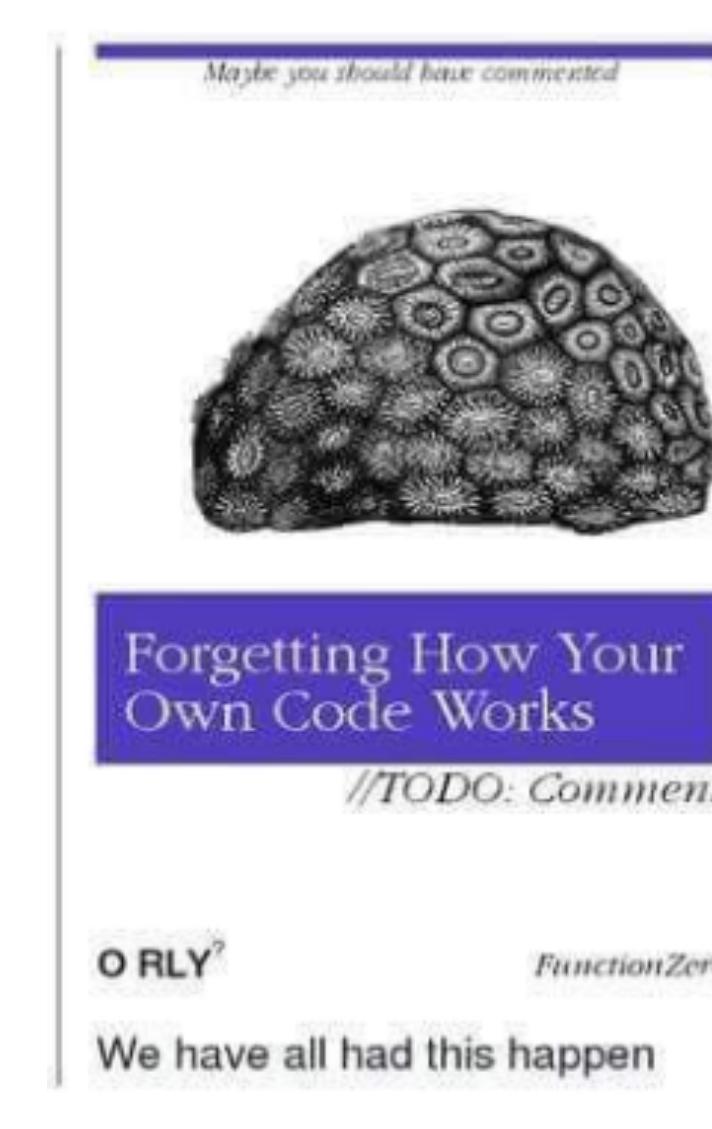
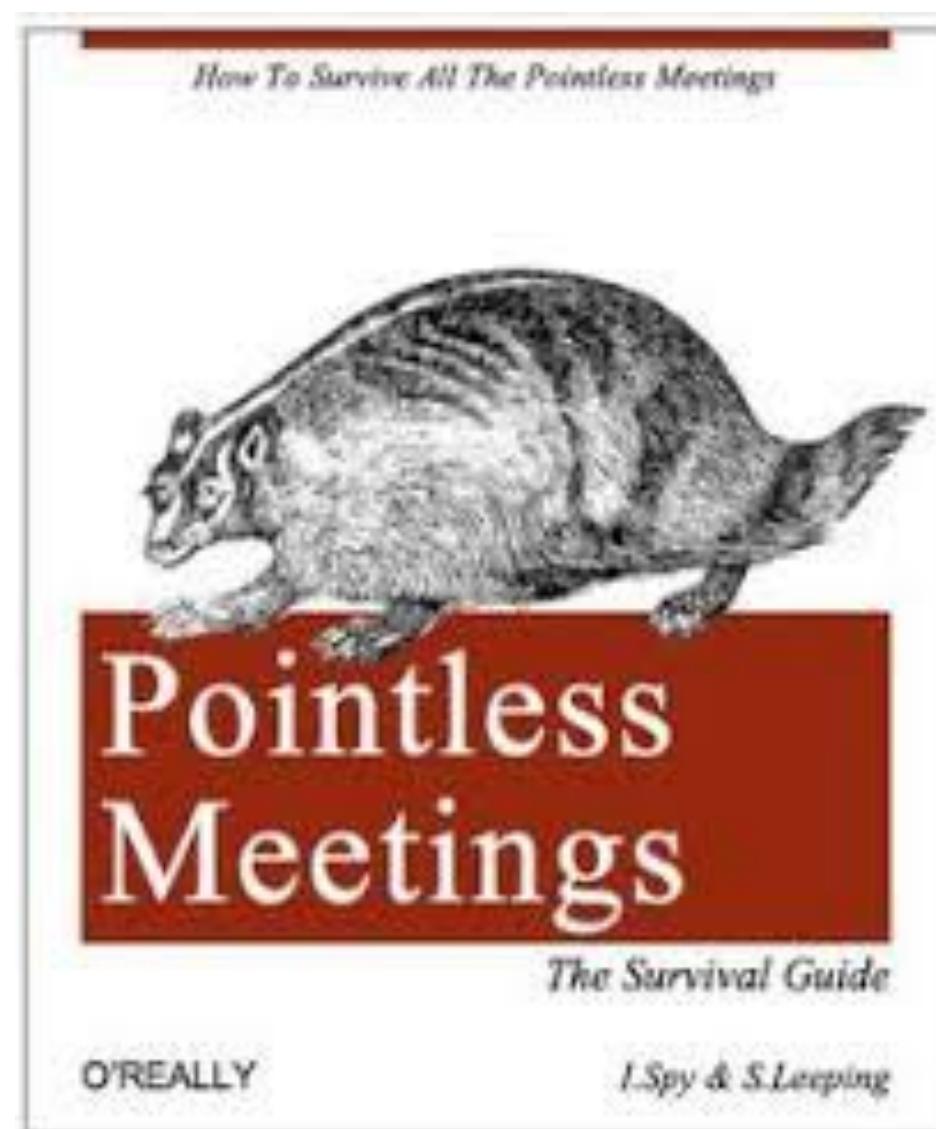
# Material

---



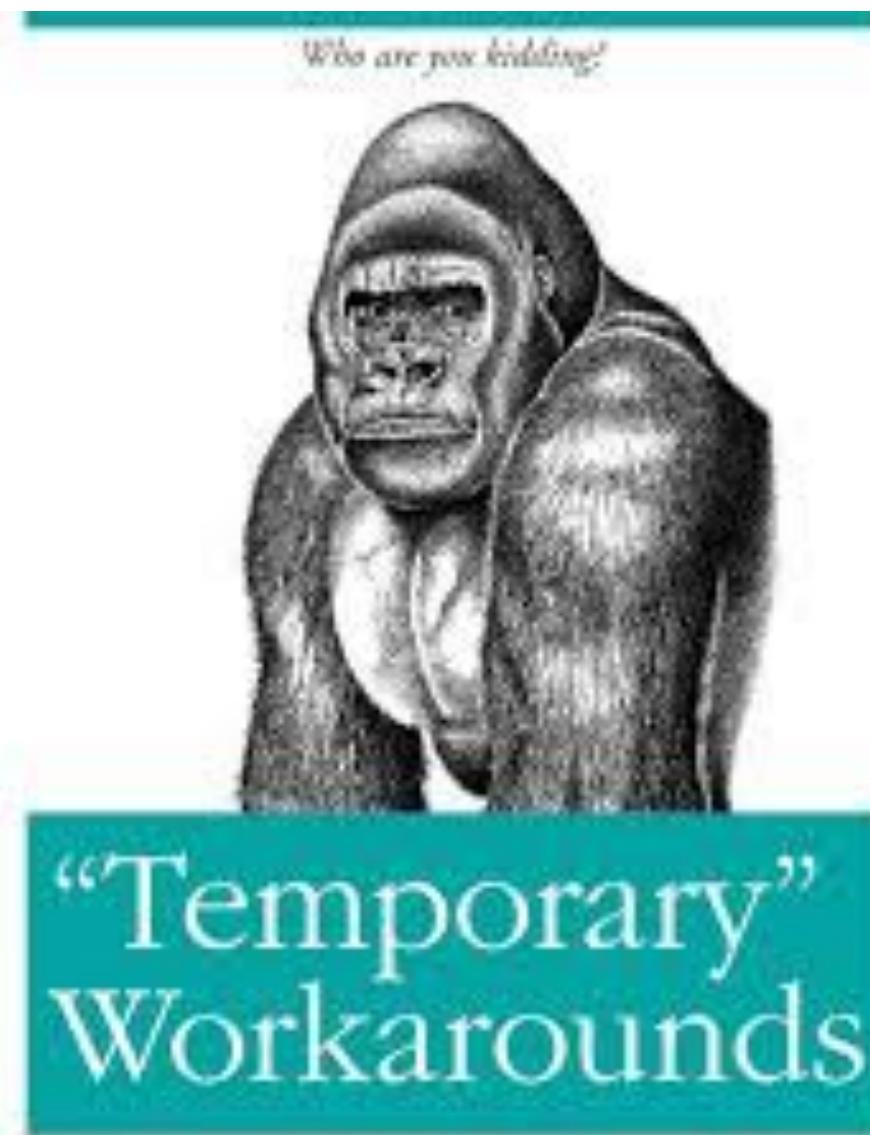
# Material

---



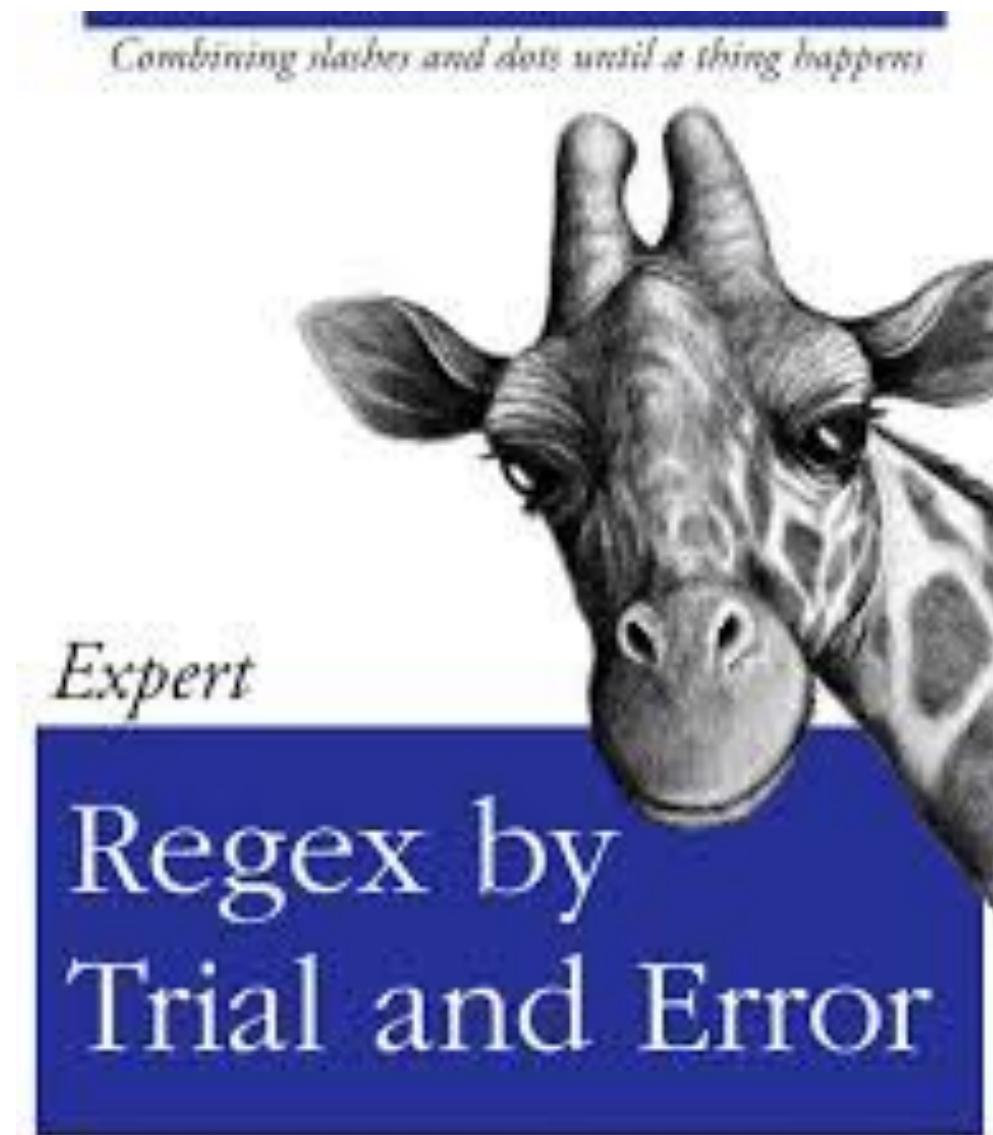
# Material

---



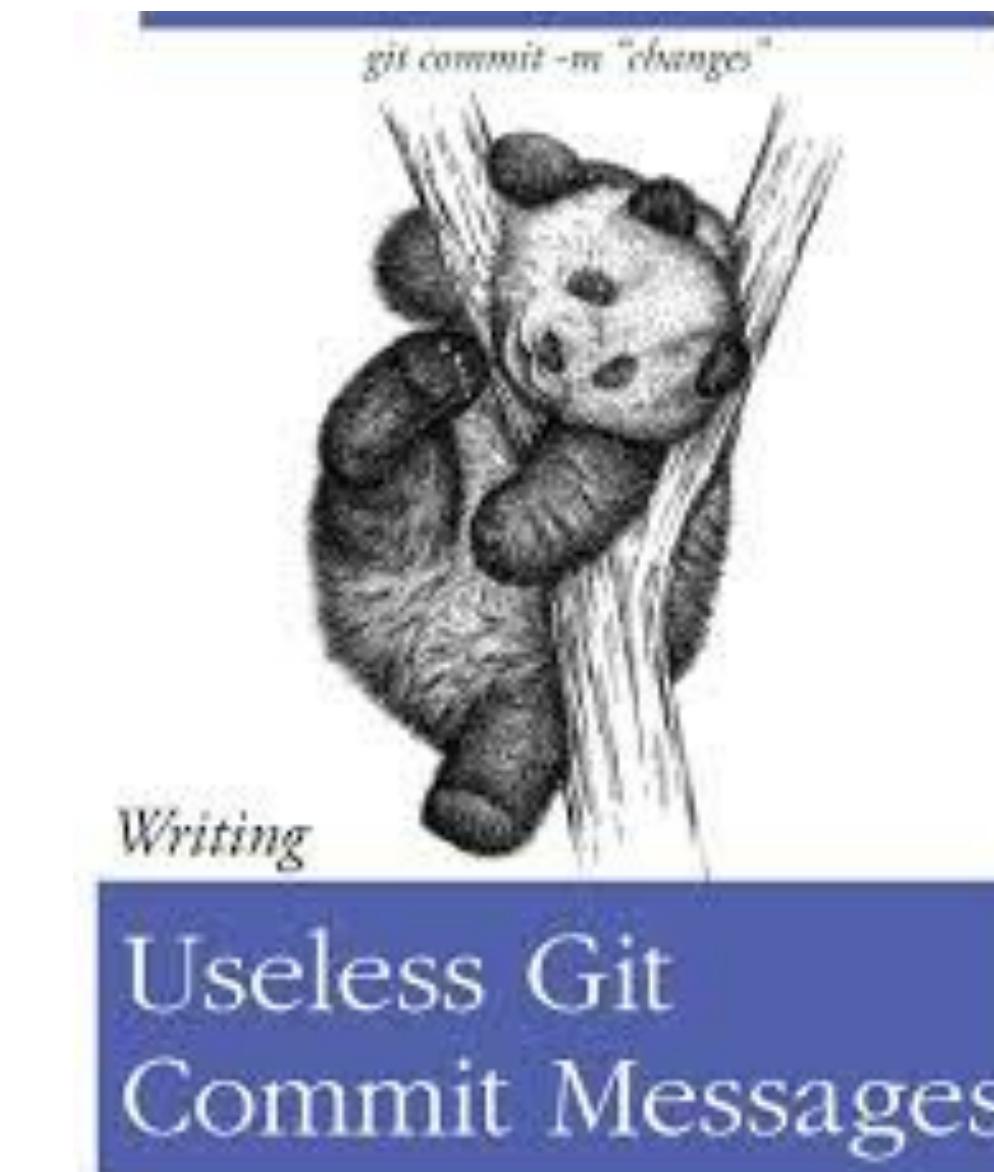
O RLY?

@ThePracticalDev



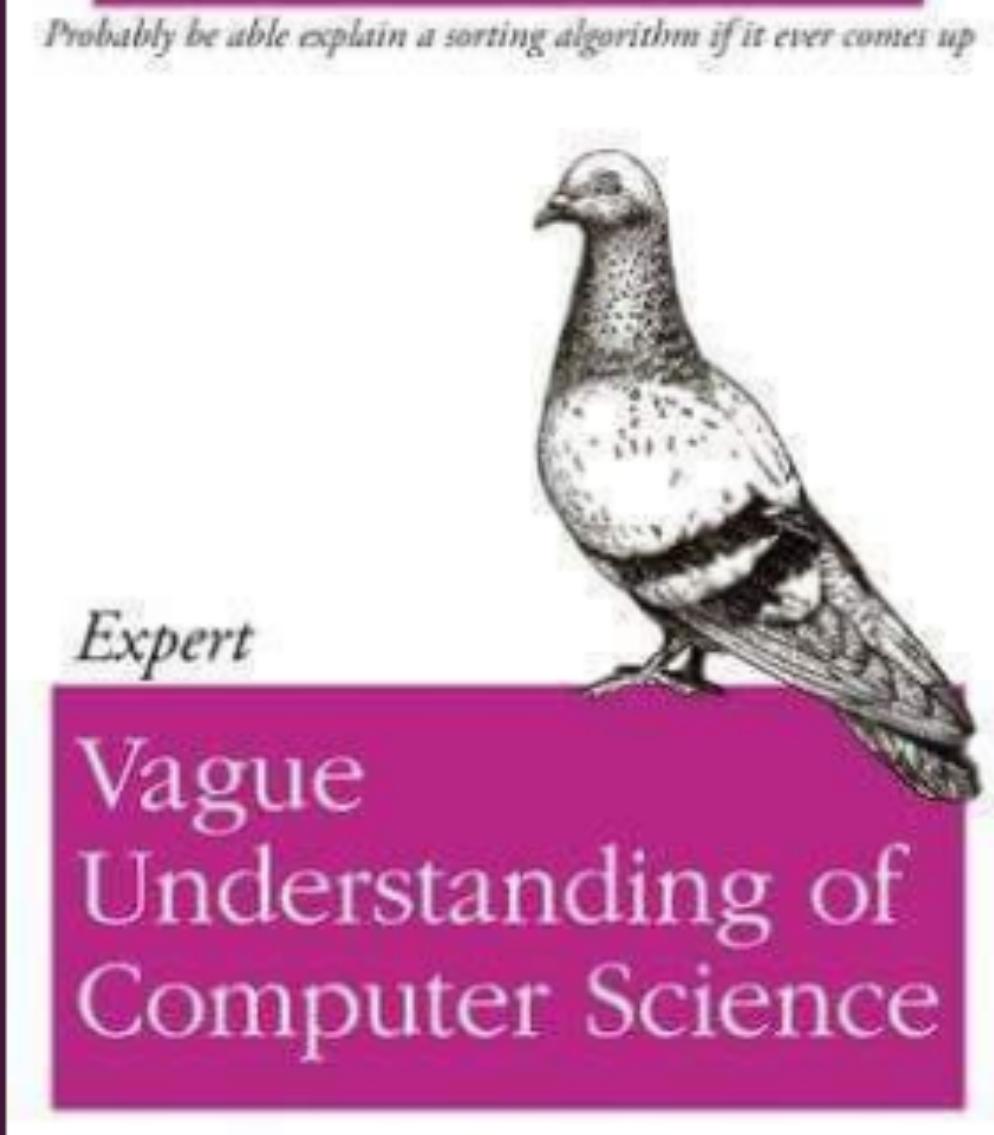
O RLY?

@ThePracticalDev



O RLY?

@ThePracticalDev



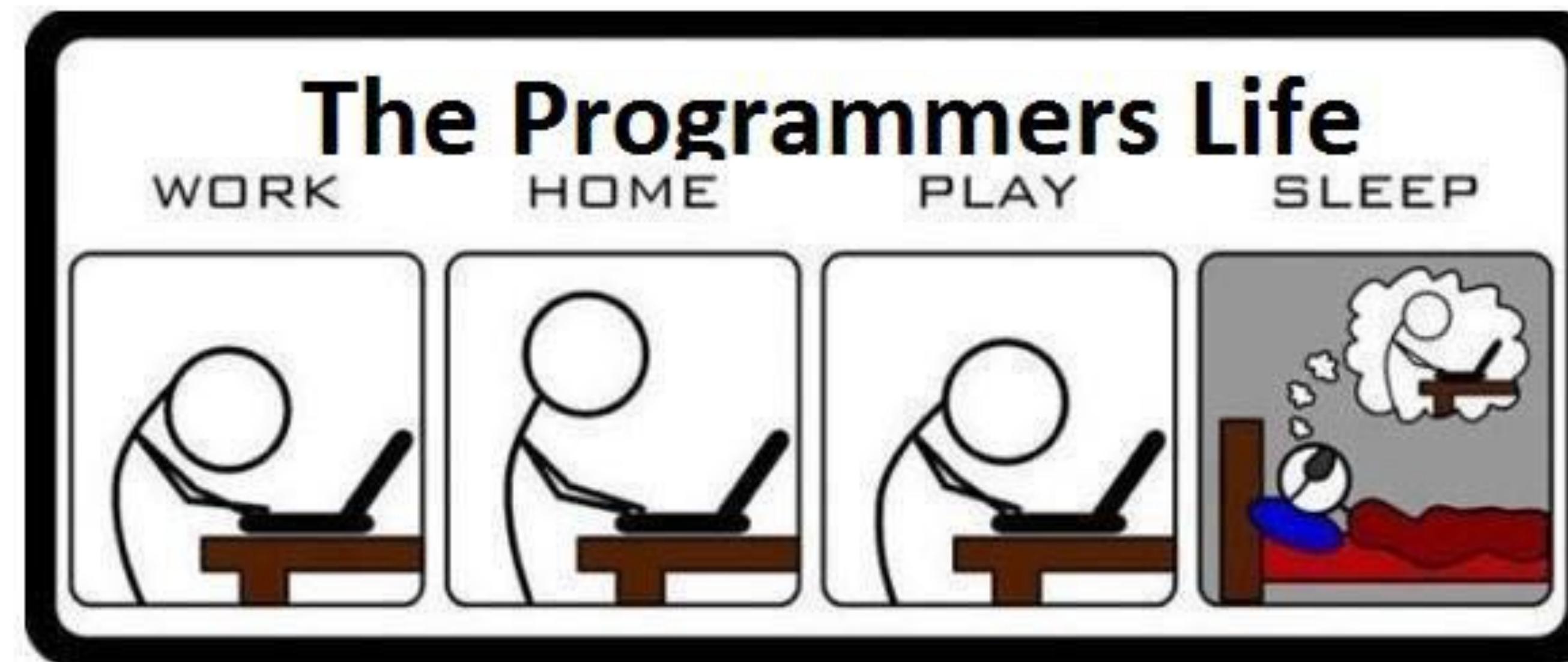
O RLY?

@ThePracticalDev

In other words.....

---

We are not interested in your programming skills....



In other words.....

---

**... but in your critical thinking skills!**



*"If you don't reveal some insights soon, I'm going  
to be forced to slice, dice, and drill!"*

# **Demystifying Data Science**

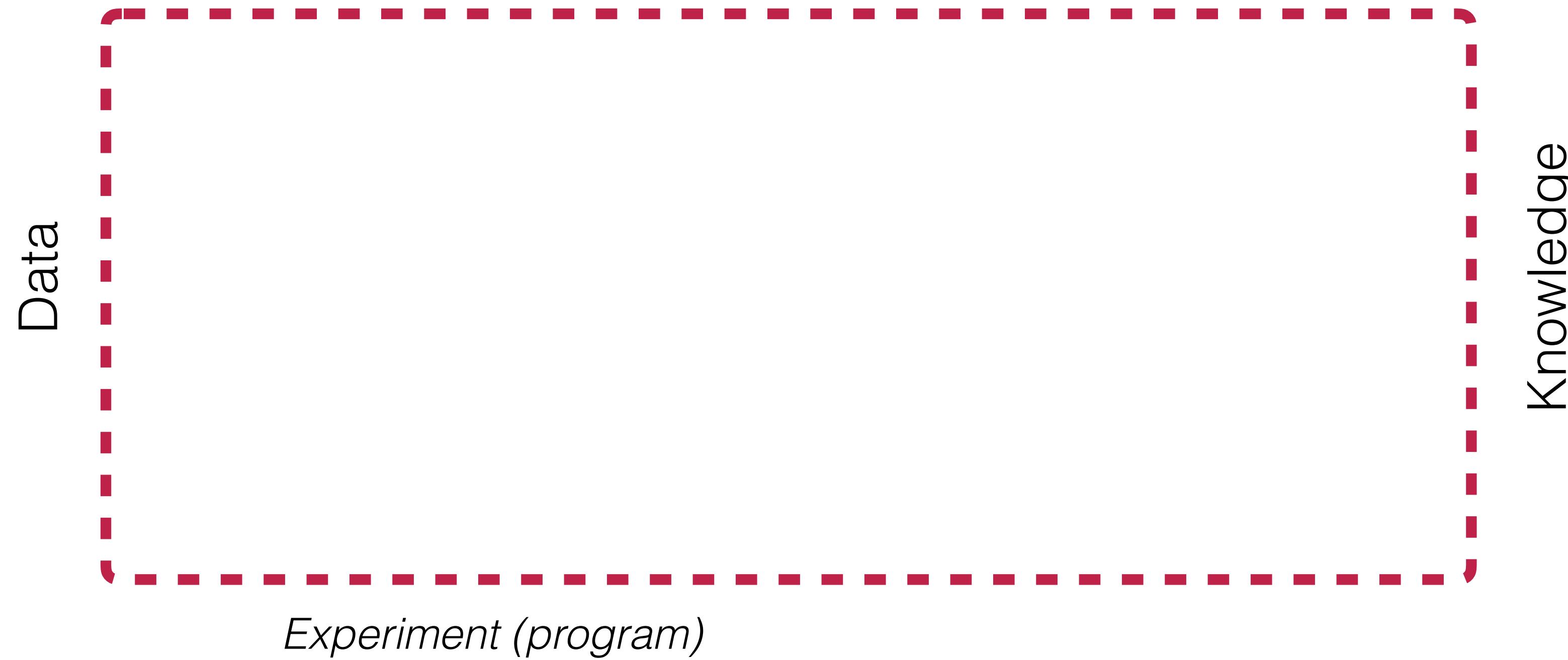
# DS != programming

---



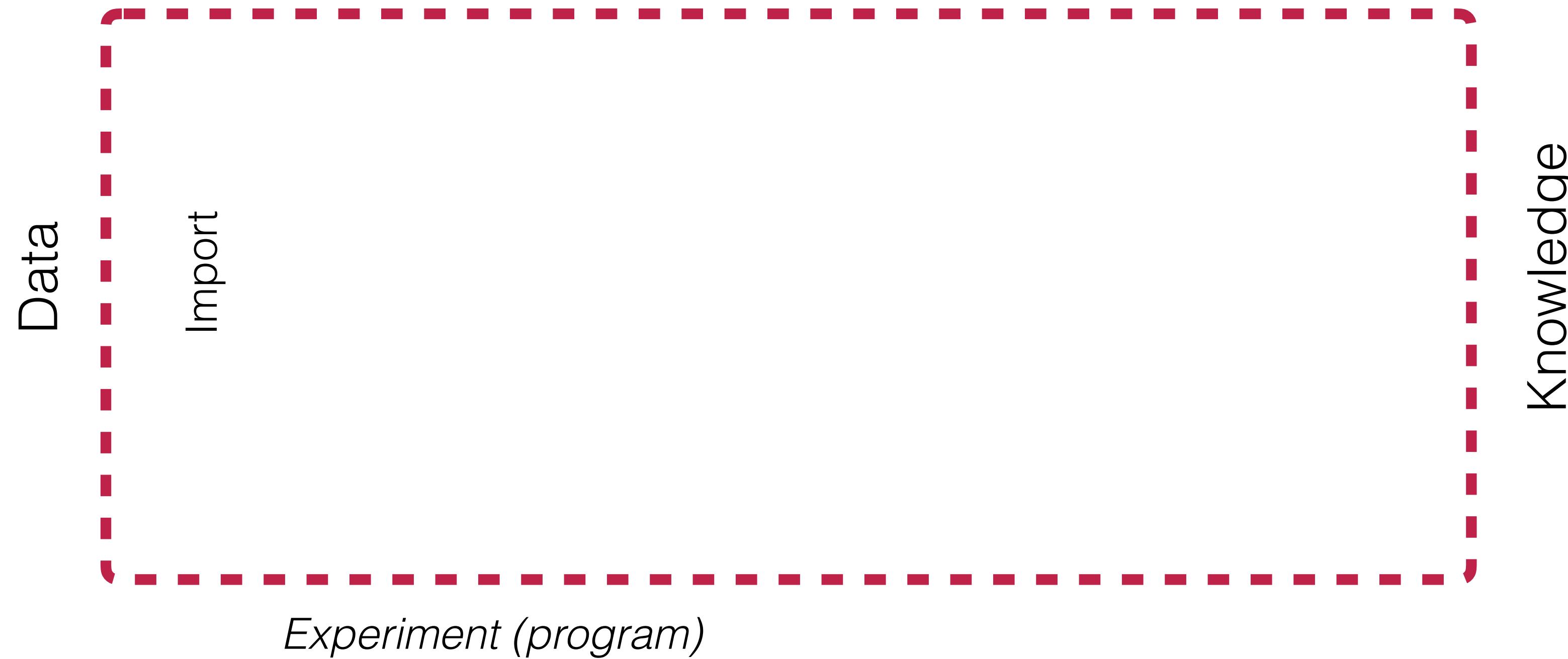
# DS Workflow

---



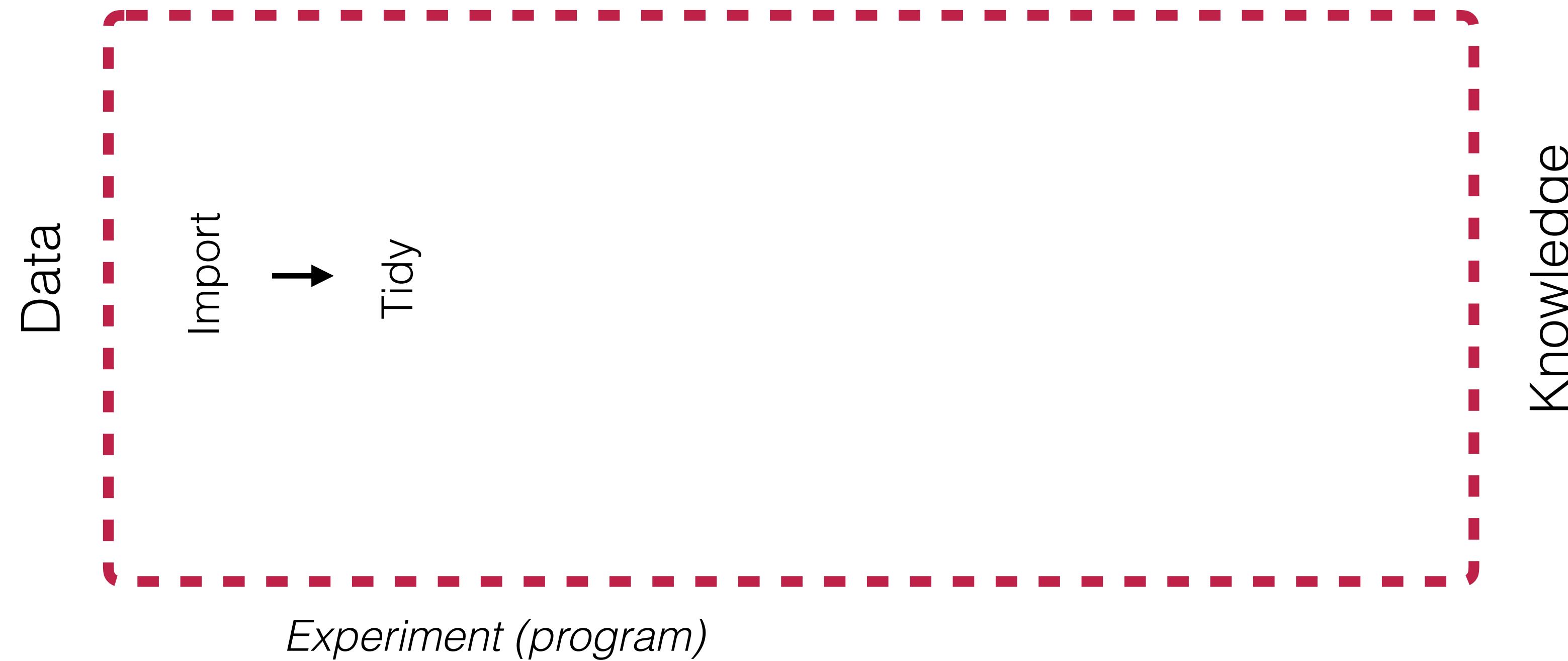
# DS Workflow

---



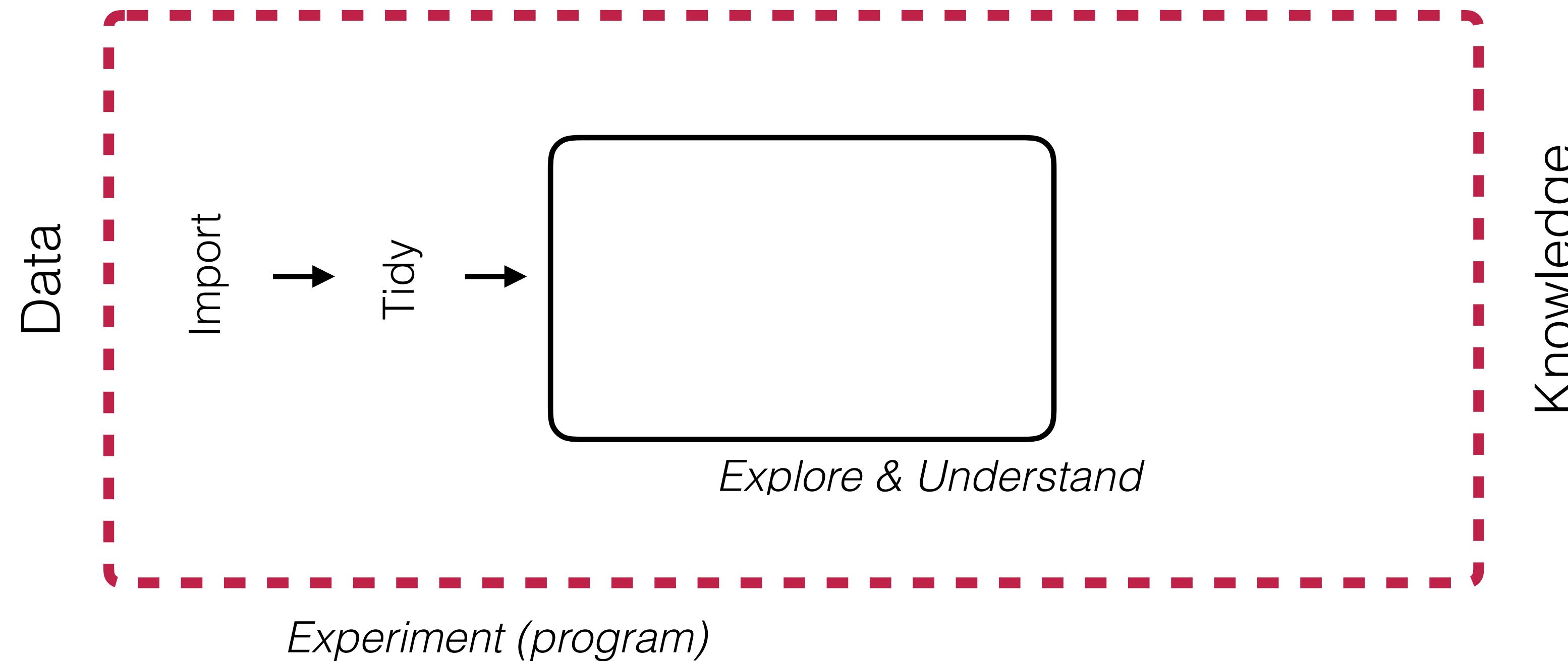
# DS Workflow

---



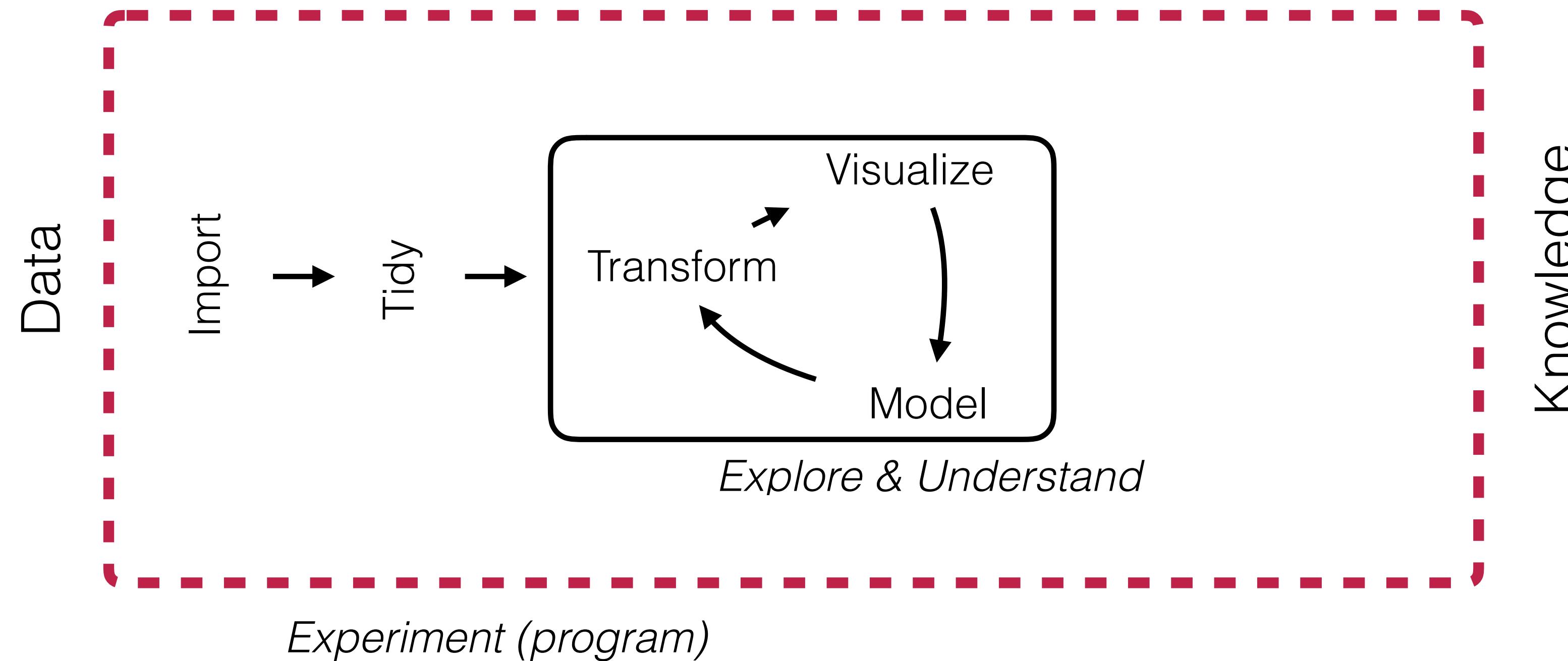
# DS Workflow

---

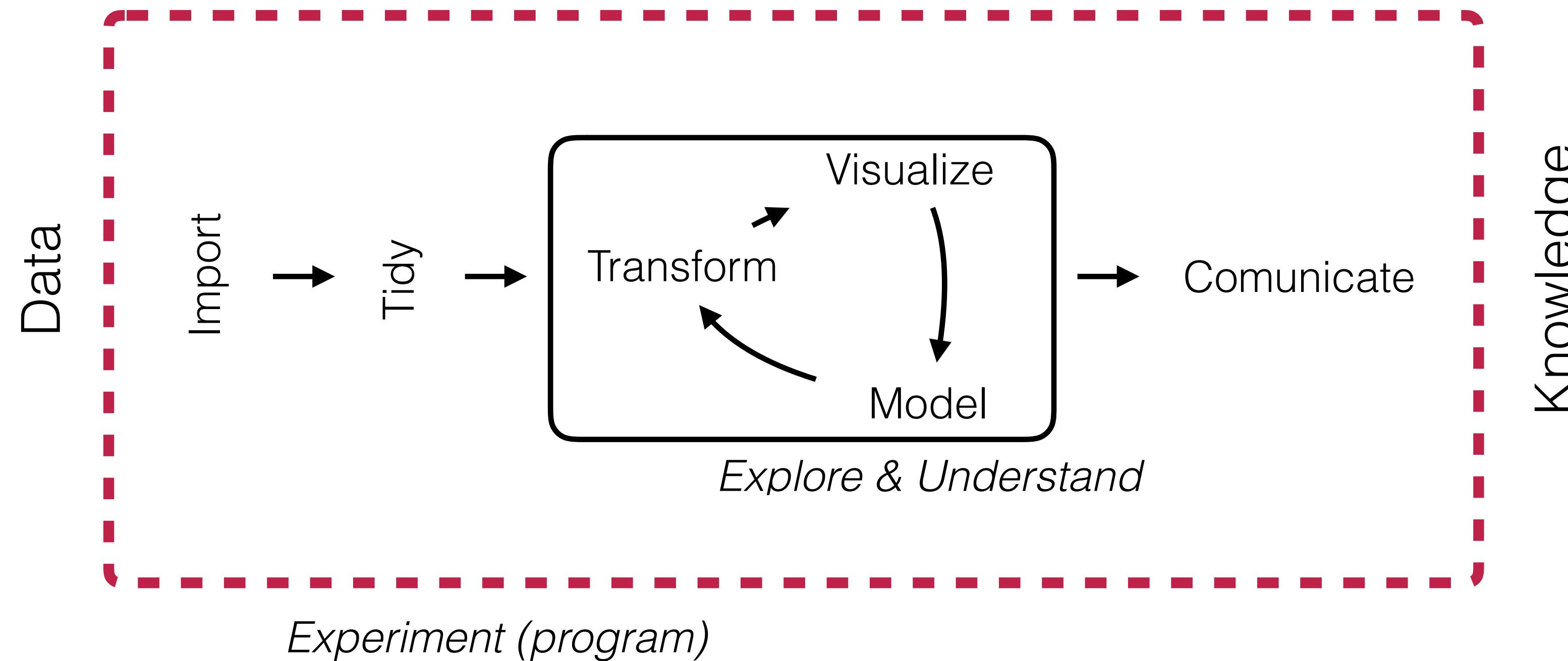


# DS Workflow

---

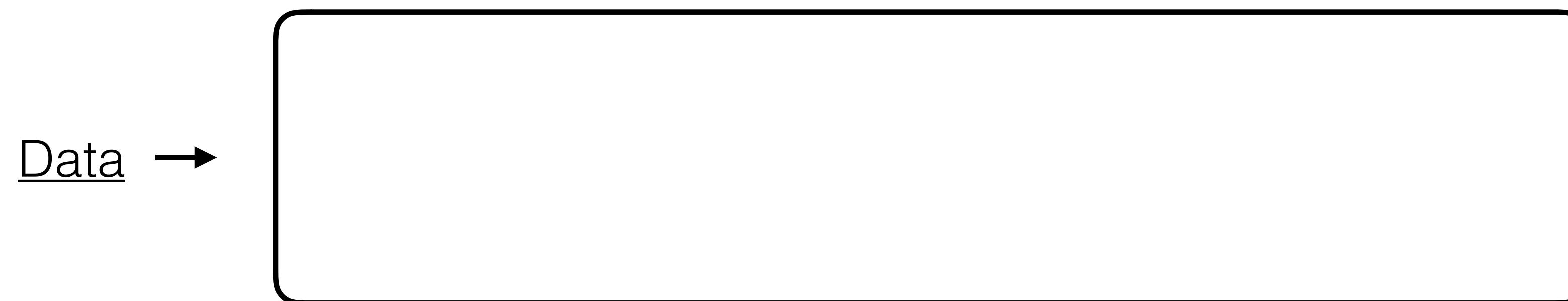


# DS Workflow



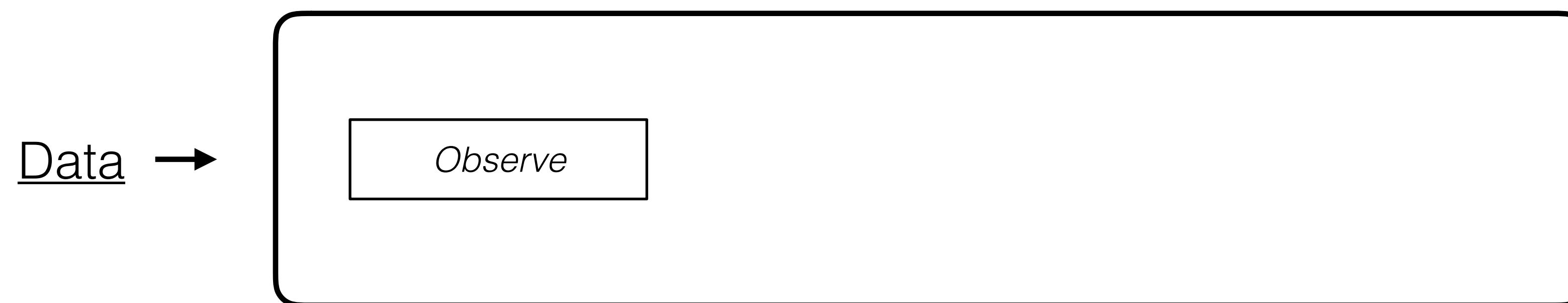
# The Scientific Method

---



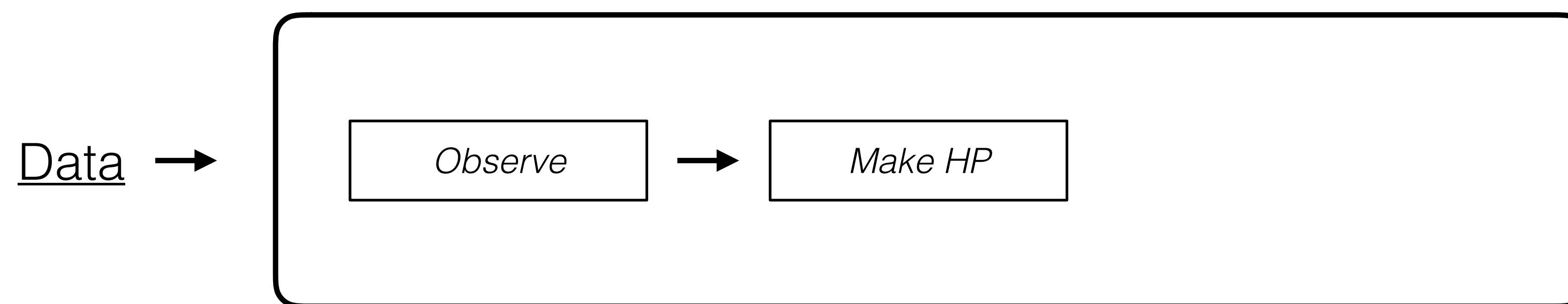
# The Scientific Method

---



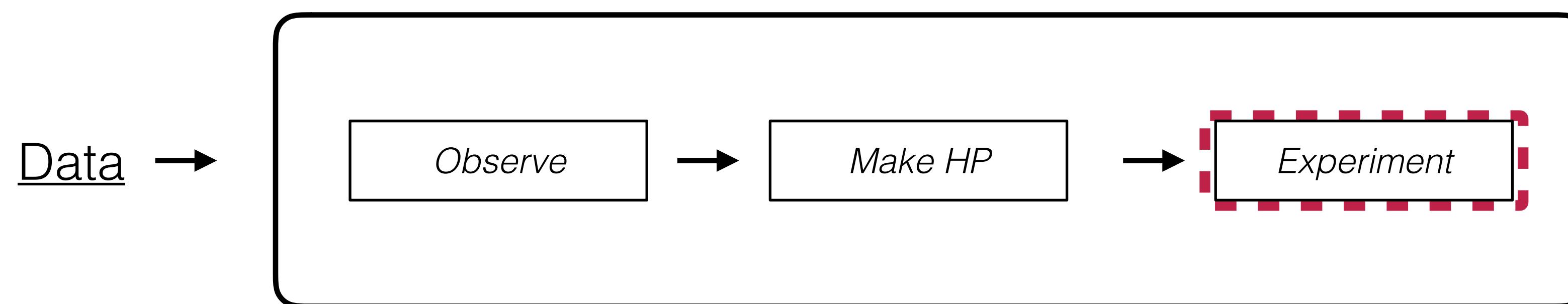
# The Scientific Method

---



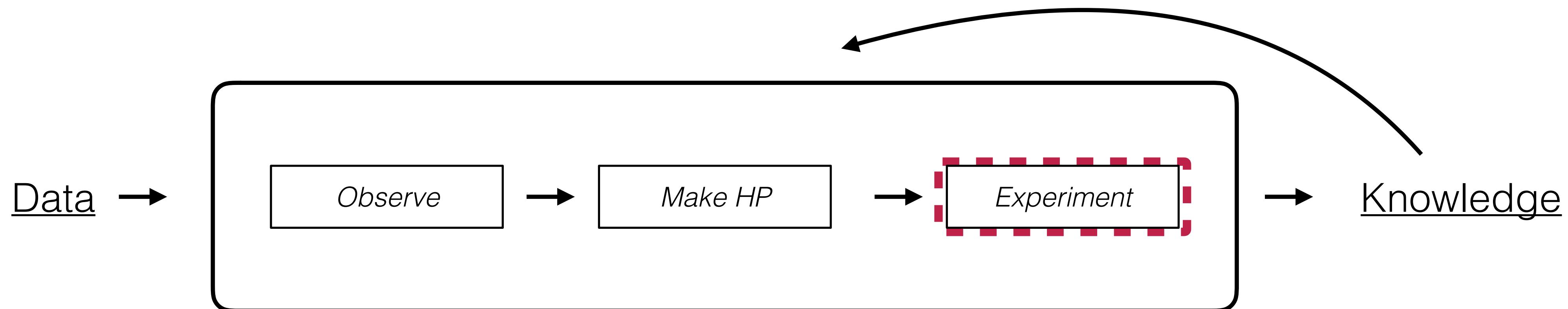
# The Scientific Method

---



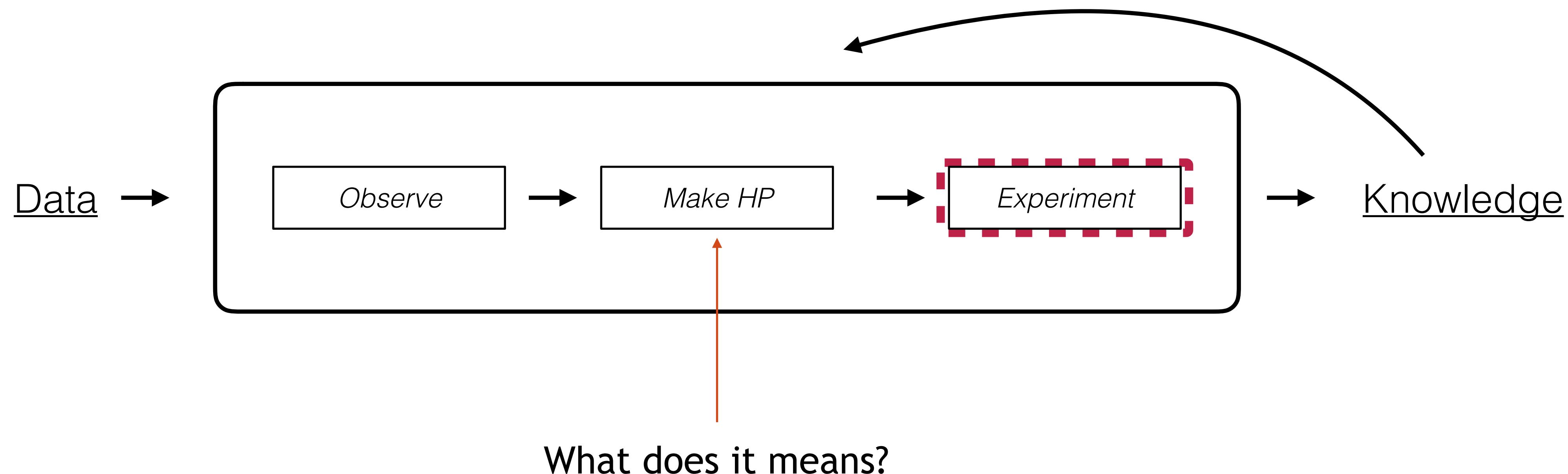
# The Scientific Method

---



# The Scientific Method

---



# Hypothesis

---

A **hypothesis** is a proposed **explanation** for a **phenomenon**. For a hypothesis to be a scientific hypothesis, the scientific method requires that one can **test** it.

Scientists generally base scientific hypotheses on previous **observations** that cannot satisfactorily be explained with the available scientific theories.

WHY DO WE NEED HP TO CREATE NEW KNOWLEDGE?

# Postulation

---

A suggestion or assumption of the existence, fact, or truth of something as a basis for reasoning, discussion, or belief.

WHY DO WE NEED POSTULATION TO CREATE NEW HYPOTHESIS?

# How to Choose an Hypothesis

---

# How to Choose an Hypothesis

---

- Testability (effort to accept or reject the hypothesis)
- Parsimony (discouraging excessive postulation)
- Scope (the apparent application of the hypothesis to multiple cases of phenomena)
- Frustrateness (the prospect that a hypothesis may explain future phenomena)
- Conservatism (the degree of fit with existing recognized knowledge-systems)

# Data Scientist

---

## Skills?

# Top Data Science Skills by Job Role

| Business Manager                        |                  | Developer                               |                  | Creative                                |                  | Researcher                              |                  |
|---|------------------|---|------------------|---|------------------|---|------------------|
| Data Skill                              | % who have skill |
| S - Communication                       | 91%              | T - Managing structured data            | 91%              | S - Communication                       | 87%              | S - Communication                       | 90%              |
| B - Project management                  | 86%              | S - Communication                       | 85%              | T - Managing structured data            | 79%              | S - Data Mining and Viz Tools           | 81%              |
| B - Business development                | 77%              | S - Data Mining and Viz Tools           | 76%              | B - Project management                  | 77%              | M - Math                                | 80%              |
| T - Managing structured data            | 74%              | B - Product Design                      | 75%              | S - Data Mining and Viz Tools           | 77%              | S - Science/Scientific Method           | 78%              |
| B - Budgeting                           | 71%              | M - Math                                | 75%              | M - Math                                | 75%              | S - Statistics and statistical modeling | 75%              |
| B - Product Design and Development      | 70%              | S - Data Management                     | 75%              | B - Product Design and Development      | 68%              | T - Managing structured data            | 73%              |
| M - Math                                | 65%              | B - Project management                  | 74%              | S - Science/Scientific Method           | 68%              | S - Data Management                     | 69%              |
| S - Data Management                     | 64%              | P - Database Administration             | 73%              | S - Data Management                     | 67%              | B - Project management                  | 68%              |
| S - Data Mining and Viz Tools           | 64%              | P - Back-end Programming                | 70%              | S - Statistics and statistical modeling | 63%              | T - Machine Learning                    | 58%              |
| B - Governance and Compliance           | 61%              | P - Systems Administration              | 65%              | B - Business development                | 58%              | M - Optimization                        | 56%              |
| S - Science/Scientific Method           | 59%              | S - Science/Scientific Method           | 64%              | B - Budgeting                           | 58%              | M - Algorithms and Simulations          | 55%              |
| S - Statistics and statistical modeling | 54%              | T - Unstructured data                   | 64%              | P - Database Administration             | 55%              | B - Product Design and Development      | 54%              |
| P - Database Administration             | 50%              | P - Front-end Programming               | 63%              | M - Graphical Models                    | 54%              | M - Graphical Models                    | 51%              |
| M - Optimization                        | 48%              | M - Algorithms                          | 61%              | M - Algorithms and Simulations          | 53%              | M - Bayesian Statistics                 | 50%              |
| M - Algorithms and Simulations          | 46%              | T - Machine Learning                    | 54%              | T - Machine Learning                    | 52%              | T - Managing unstructured data          | 46%              |
| M - Graphical Models                    | 46%              | S - Statistics and statistical modeling | 52%              | T - Managing unstructured data          | 52%              | P - Database Administration             | 46%              |
| T - Managing unstructured data          | 43%              | M - Graphical Models                    | 51%              | M - Optimization                        | 49%              | B - Budgeting                           | 44%              |
| P - Front-end Programming               | 42%              | B - Business development                | 48%              | P - Front-end Programming               | 49%              | B - Business development                | 43%              |
| T - Machine Learning                    | 41%              | T - Big and distributed data            | 48%              | B - Governance and Compliance           | 48%              | T - NLP and text mining                 | 41%              |
| P - Systems Administration and Design   | 41%              | M - Optimization                        | 47%              | P - Systems Administration and Design   | 46%              | P - Systems Administration and Design   | 40%              |
| M - Bayesian Statistics                 | 39%              | B - Budgeting                           | 45%              | M - Bayesian Statistics                 | 45%              | P - Front-end Programming               | 38%              |
| P - Back-end Programming                | 34%              | B - Governance and Compliance           | 44%              | P - Back-end Programming                | 45%              | P - Back-end Programming                | 35%              |
| P - Cloud Management                    | 31%              | M - Bayesian Statistics                 | 44%              | T - NLP and text mining                 | 39%              | B - Governance and Compliance           | 34%              |
| T - NLP and text mining                 | 30%              | T - NLP                                 | 43%              | T - Big and distributed data            | 36%              | T - Big and distributed data            | 32%              |
| T - Big and distributed data            | 29%              | P - Cloud Management                    | 42%              | P - Cloud Management                    | 30%              | P - Cloud Management                    | 21%              |

Note: % who have skill reflects percent of respondents who indicated that they have, at least, an intermediate level of proficiency.

Business Manager (e.g., leader, business person, entrepreneur) N = 250; Developer (e.g., developer, engineer) N = 222; Creative (e.g., Jack of all trades, artist, hacker) N = 221; Researcher (e.g., researcher, scientist, statistician) N = 353

# Data Scientist

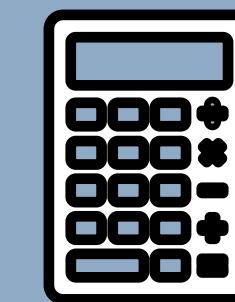
---

1. S – Communication (87% possess skill)
2. T – Managing Structured data (75%)
3. M – Math (71%)
4. B – Project management (71%)
5. S – Data Mining and Viz Tools (71%)
6. S – Science/Scientific Method (65%)
7. S – Data Management (65%)
8. B – Product design and development (59%)
9. S – Statistics and statistical modeling (59%)
10. B – Business development (53%)

# Data Scientist

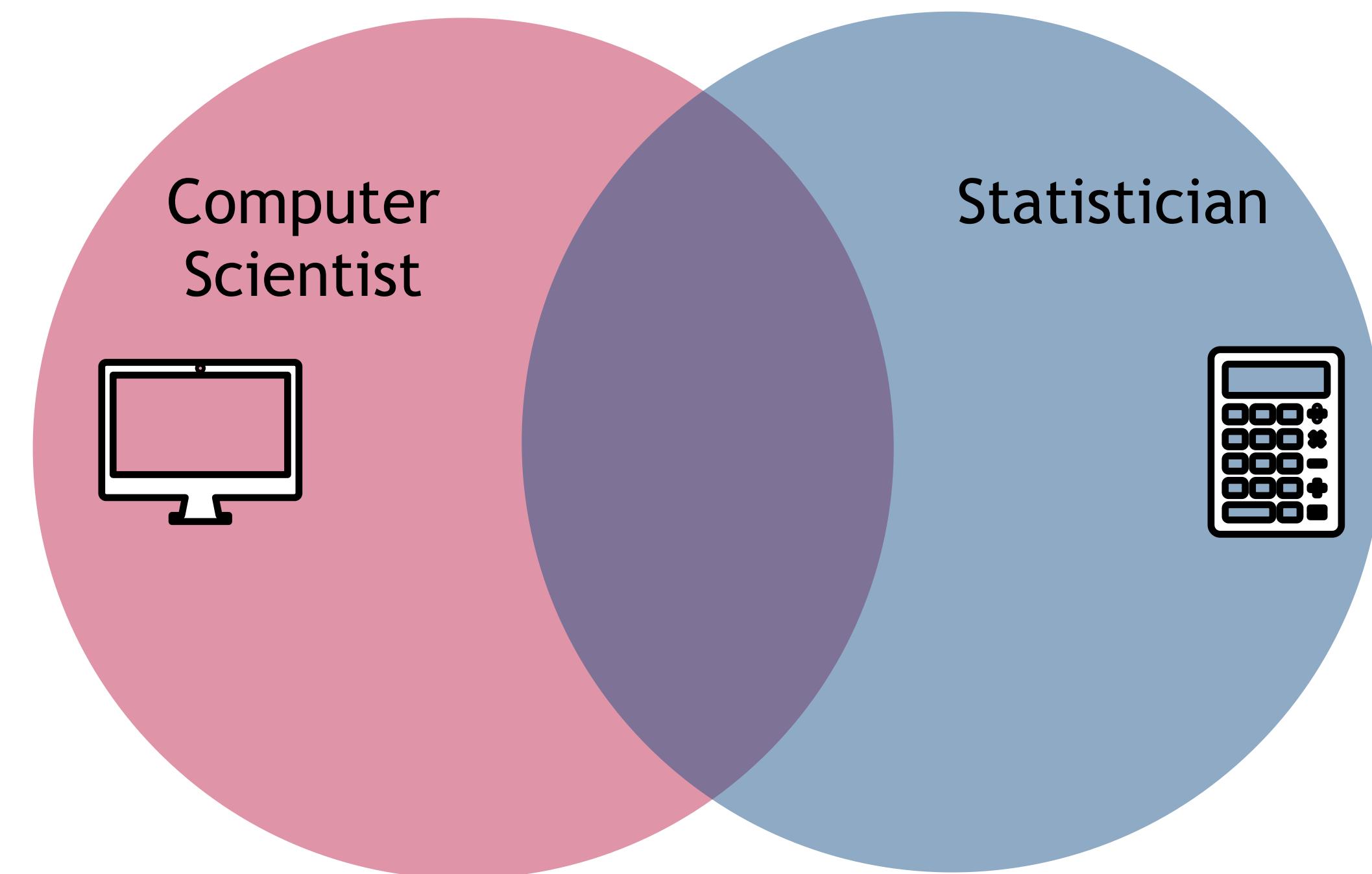
---

Statistician



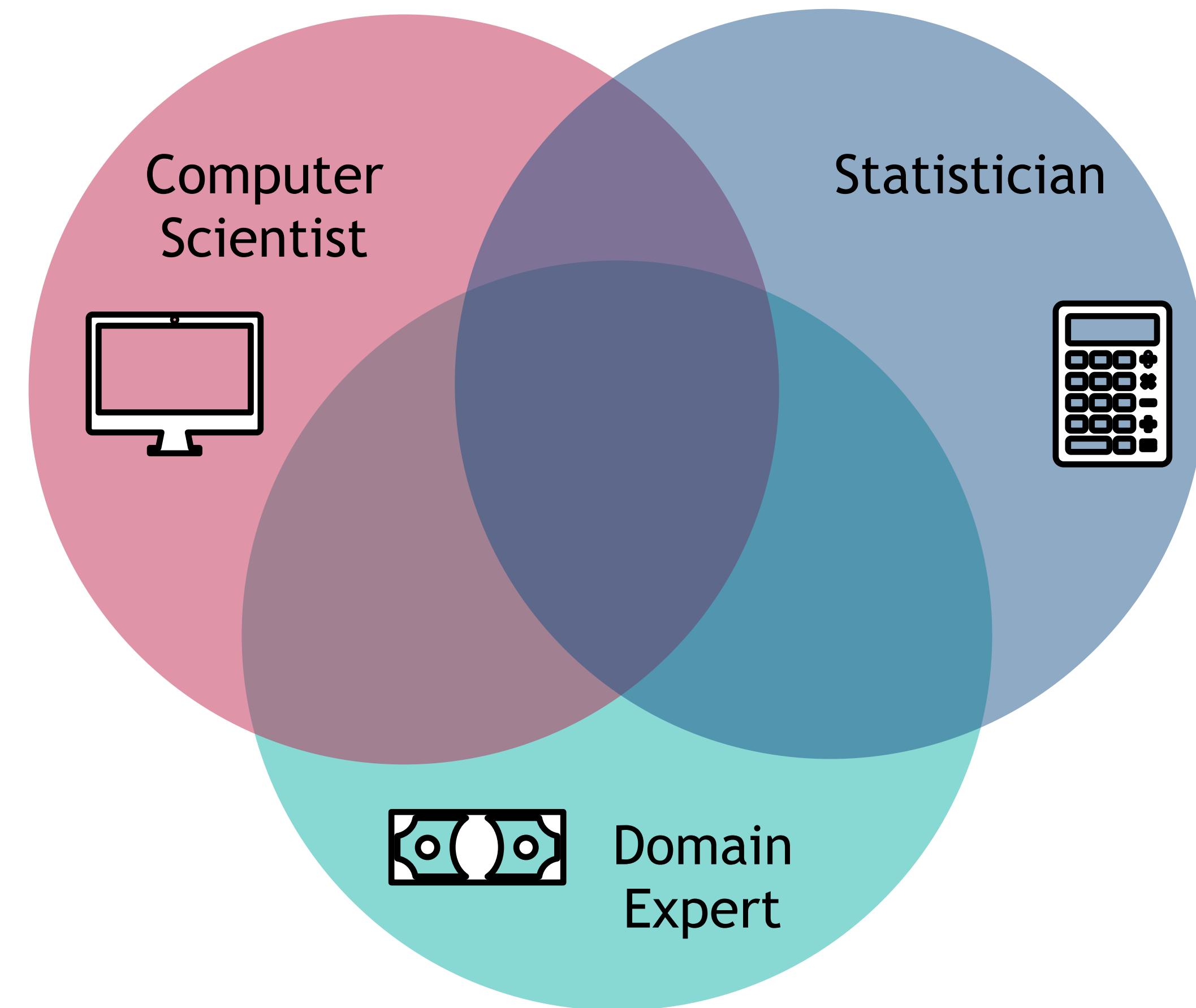
# Data Scientist

---



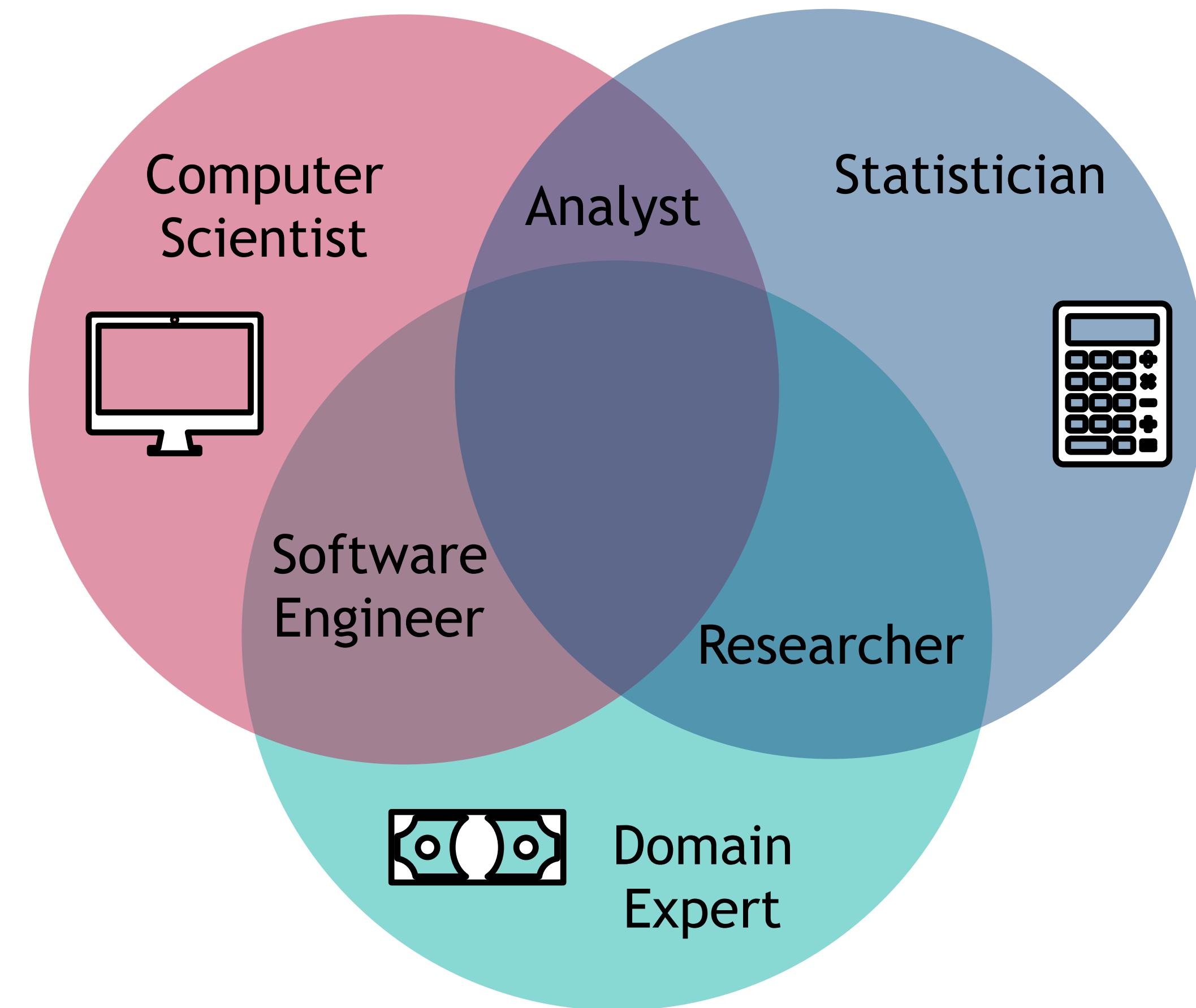
# Data Scientist

---



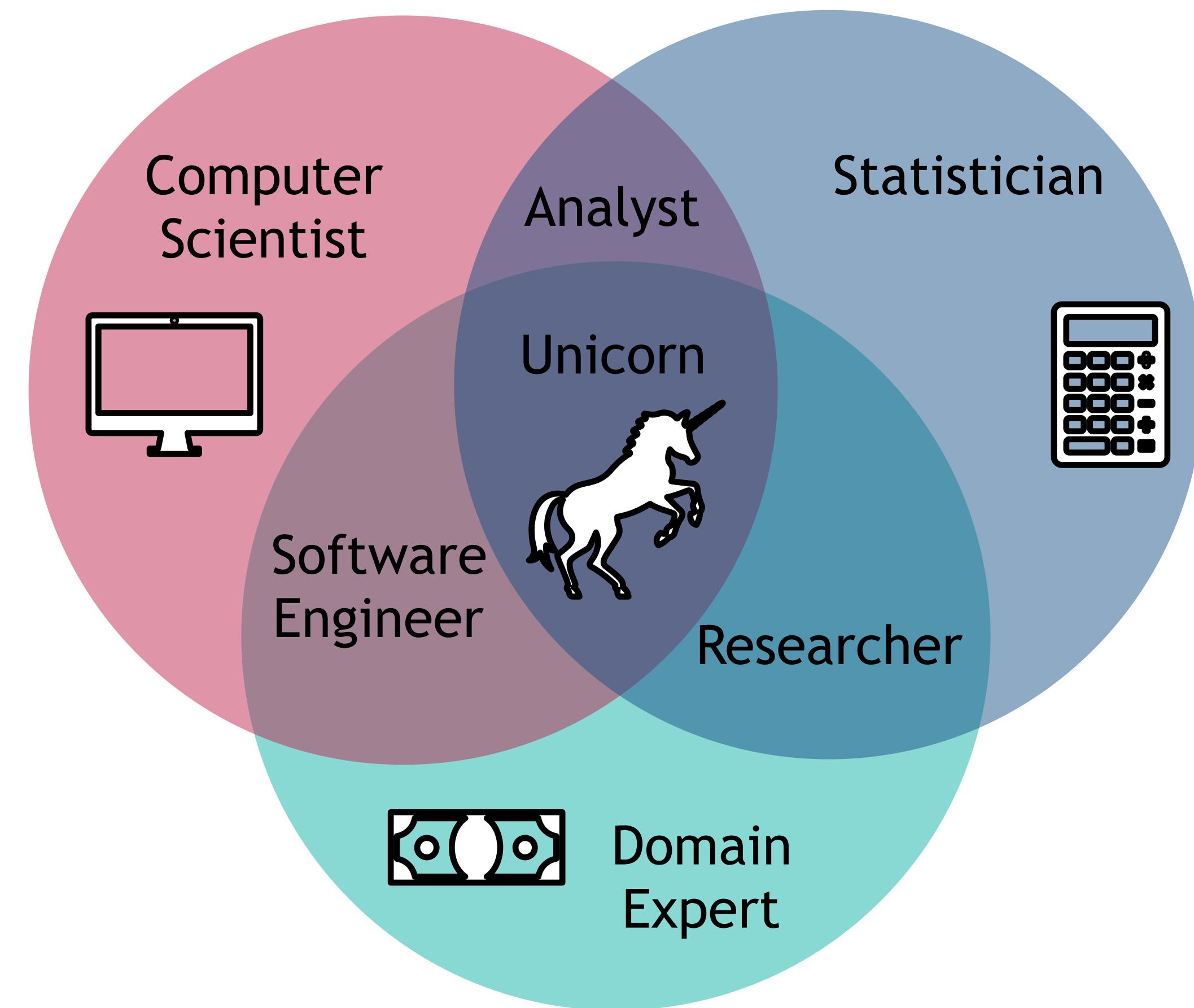
# Data Scientist

---



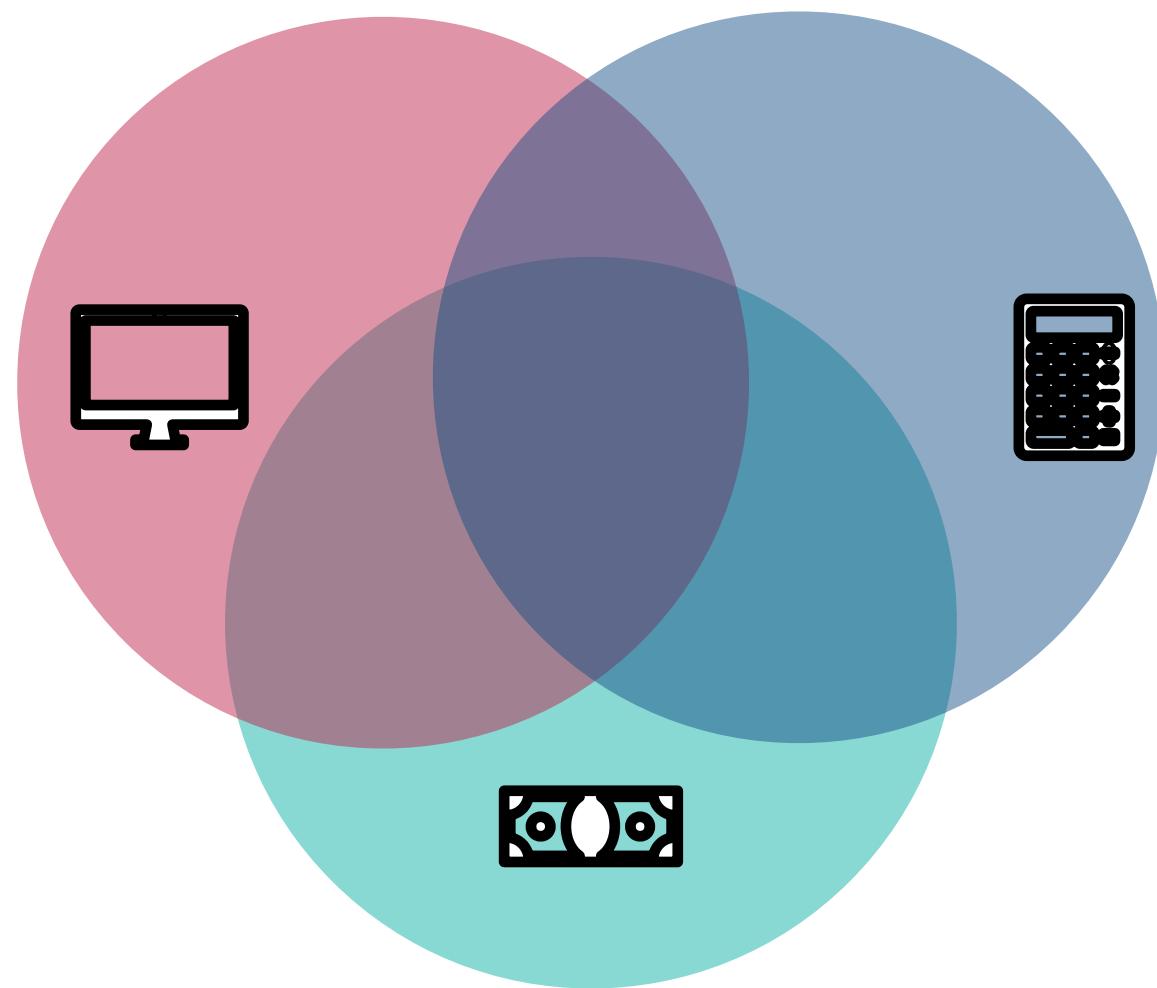
# Data Scientist

---



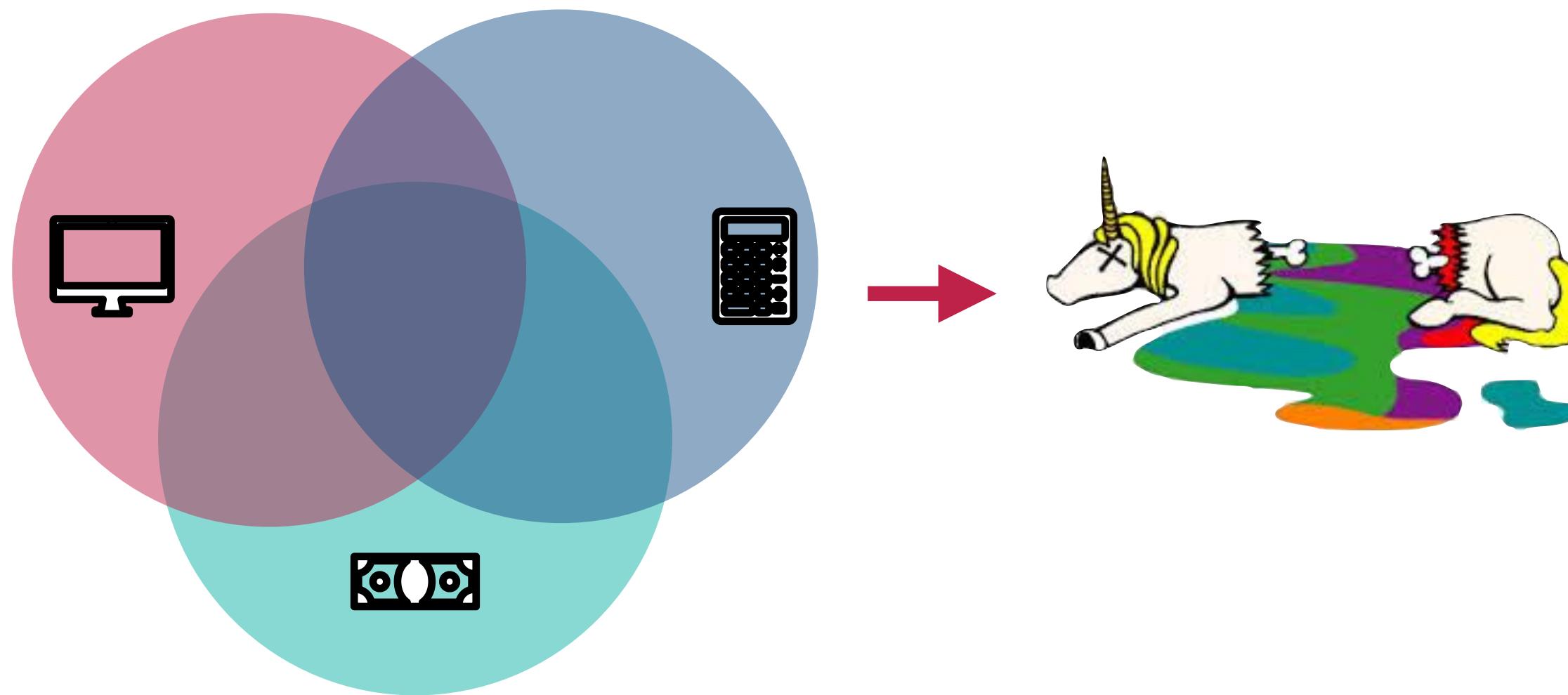
# The importance of Team in Data Science

---



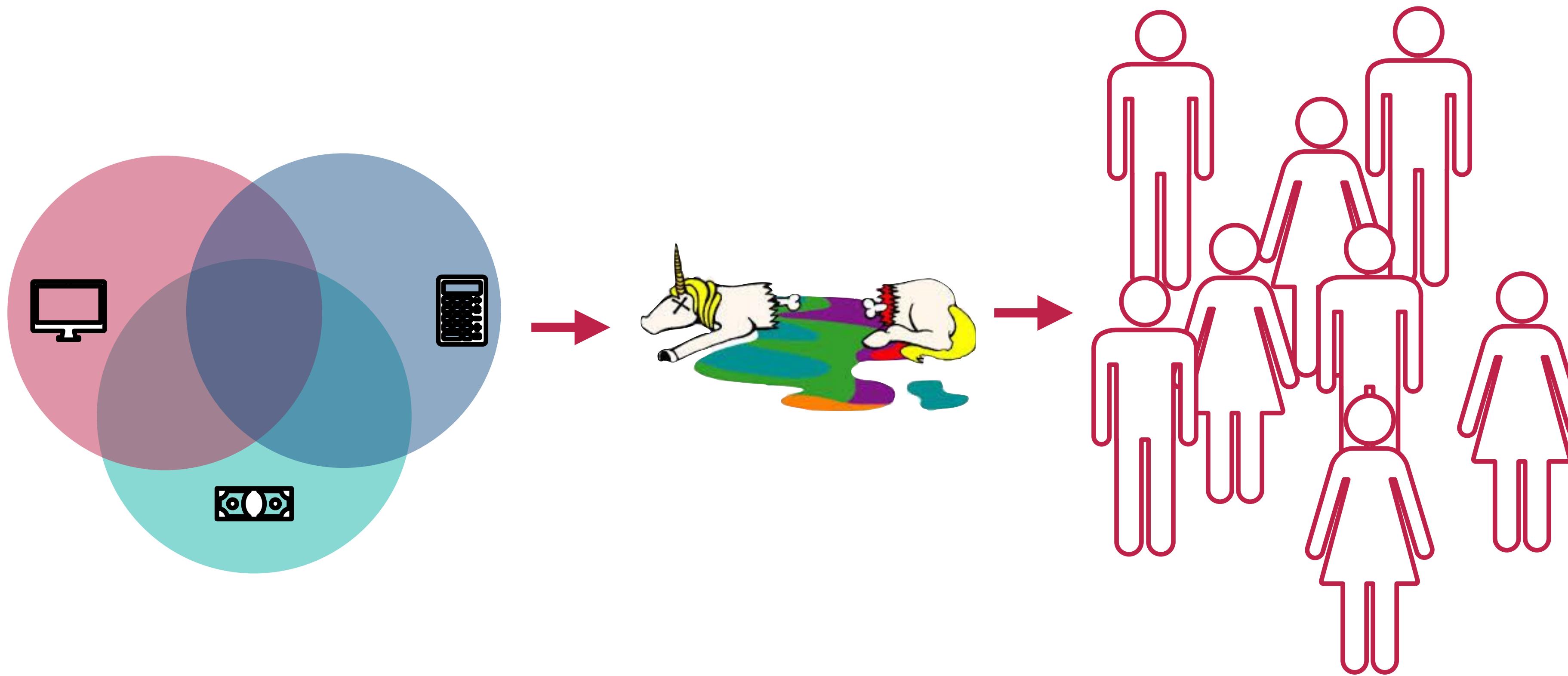
# The importance of Team in Data Science

---



# The importance of Team in Data Science

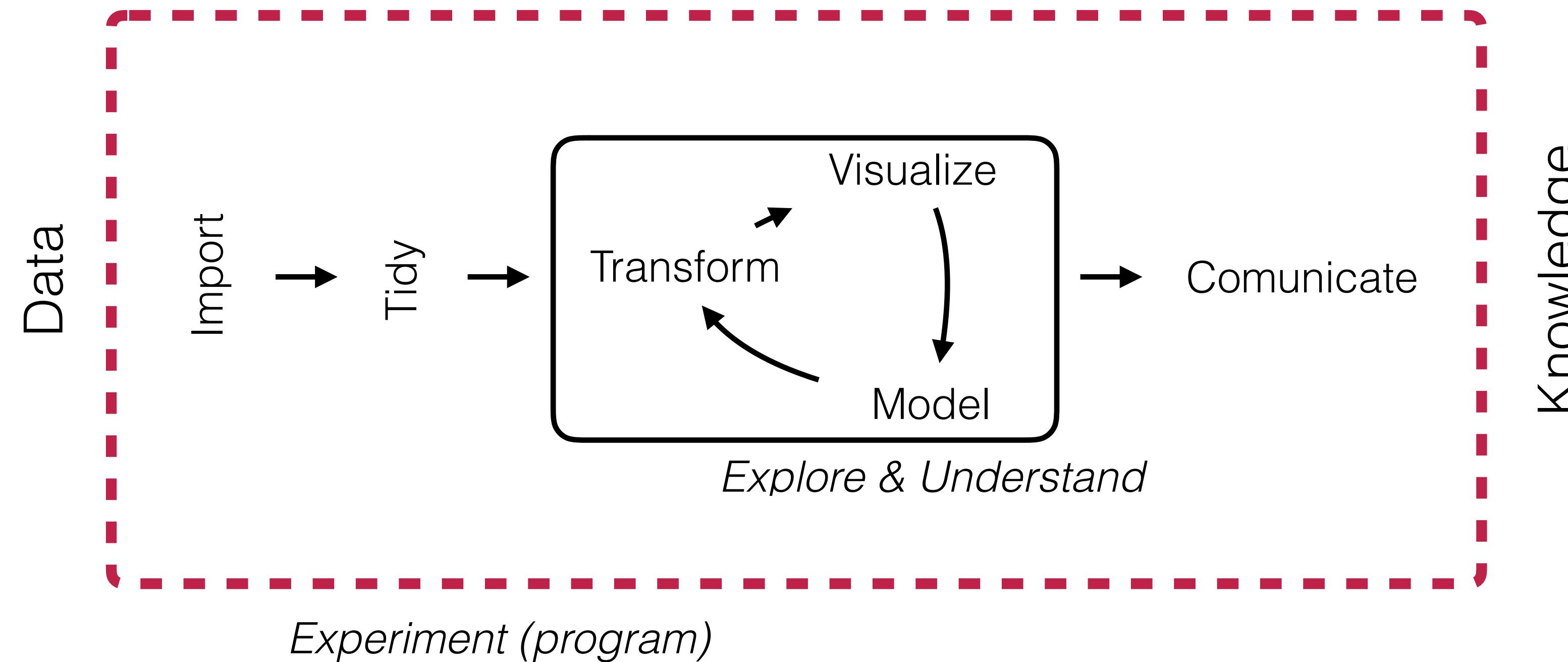
---

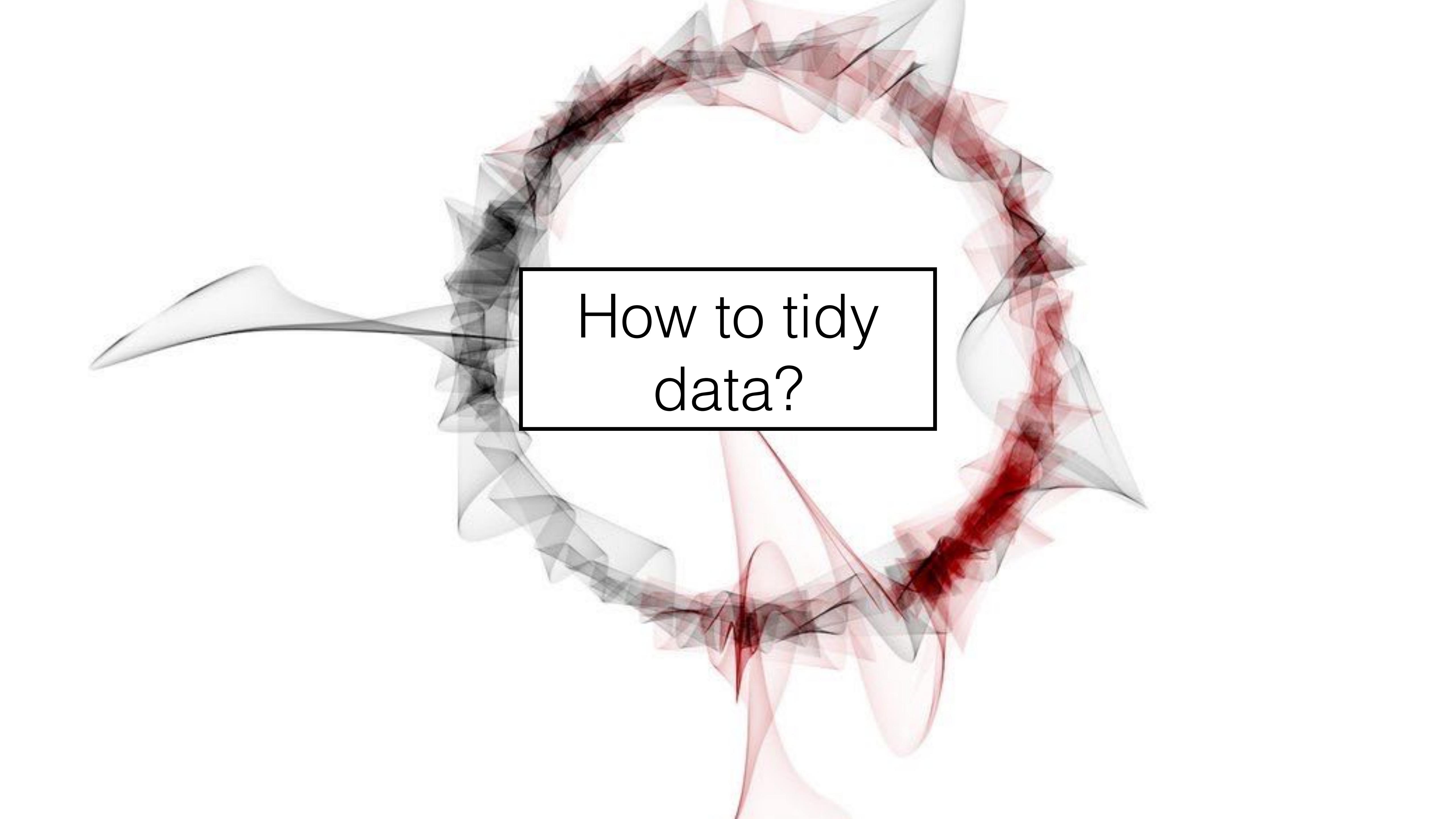


# R for Data Science



# DS Workflow: Which Verb you Didn't expect?





How to tidy  
data?

# The Tidyverse

---

*What is tidy data? Data format that:*

- Makes data analysis easy
- Is easy to model, visualize and transform
- Facilitate the creative process



# The Tydyverse

---

There are three interrelated rules which make a dataset tidy:

- 1- Each variable must have its own column.
- 2- Each observation must have its own row.
- 3- Each value must have its own cell.

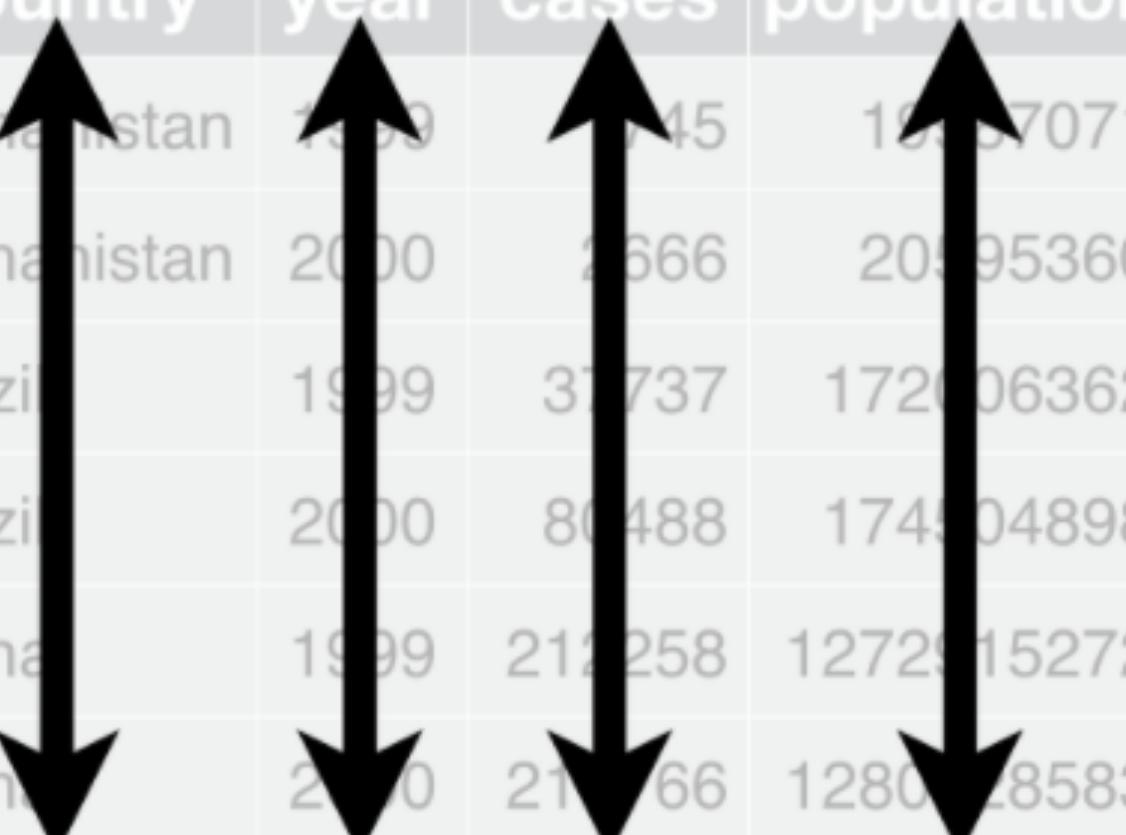


# The Tydyverse

---

| country     | year | cases | population |
|-------------|------|-------|------------|
| Afghanistan | 1999 | 145   | 19081071   |
| Afghanistan | 2000 | 2666  | 20595360   |
| Brazil      | 1999 | 3737  | 172006362  |
| Brazil      | 2000 | 80488 | 174504898  |
| China       | 1999 | 21258 | 1272915272 |
| China       | 2000 | 21666 | 128042583  |

variables



| country     | year | cases | population |
|-------------|------|-------|------------|
| Afghanistan | 1999 | 145   | 19081071   |
| Afghanistan | 2000 | 2666  | 20595360   |
| Brazil      | 1999 | 3737  | 172006362  |
| Brazil      | 2000 | 80488 | 174504898  |
| China       | 1999 | 21258 | 1272915272 |
| China       | 2000 | 21666 | 128042583  |

observations



| country     | year | cases | population |
|-------------|------|-------|------------|
| Afghanistan | 1999 | 145   | 19081071   |
| Afghanistan | 2000 | 2666  | 20595360   |
| Brazil      | 1999 | 3737  | 172006362  |
| Brazil      | 2000 | 80488 | 174504898  |
| China       | 1999 | 21258 | 1272915272 |
| China       | 2000 | 21666 | 128042583  |

values



# Example of Messiness

---

|        | Pregnant | Not Pregnant |
|--------|----------|--------------|
| Male   | 0        | 5            |
| Female | 1        | 4            |

*There are three variables in this data set. What are they?*

# Example of Messiness

---

| Pregnant | Sex    | Freq |
|----------|--------|------|
| no       | female | 4    |
| no       | male   | 5    |
| yes      | female | 1    |
| yes      | male   | 0    |



- R is an **open source** programming language for statistical computing and graphics that is supported by the R Foundation for Statistical Computing.
- Polls, surveys of **data miners**, and studies of scholarly literature databases show that R's popularity has increased substantially in recent years.

# Alternatives?

---

# Alternatives?

---





VS



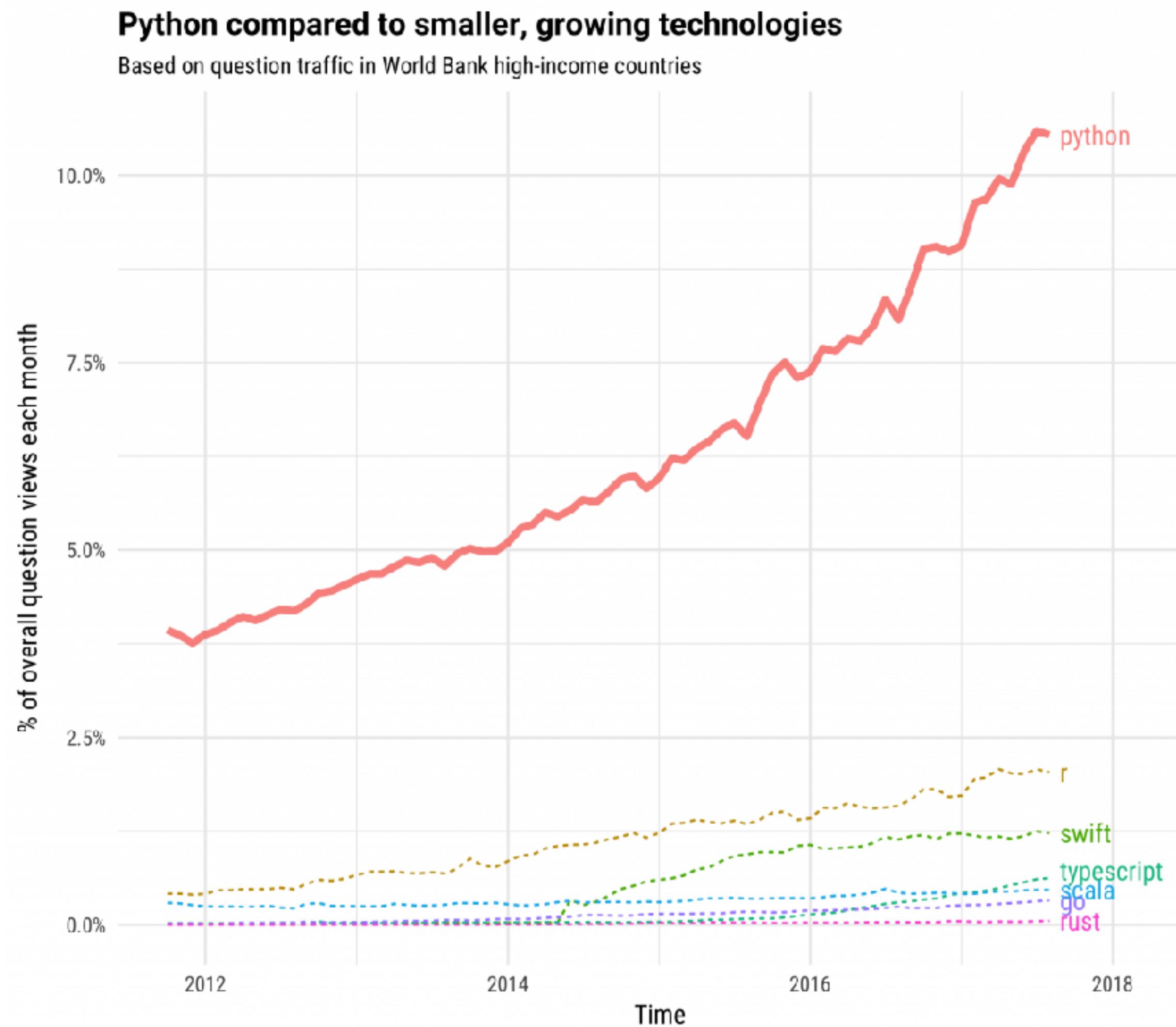


VS



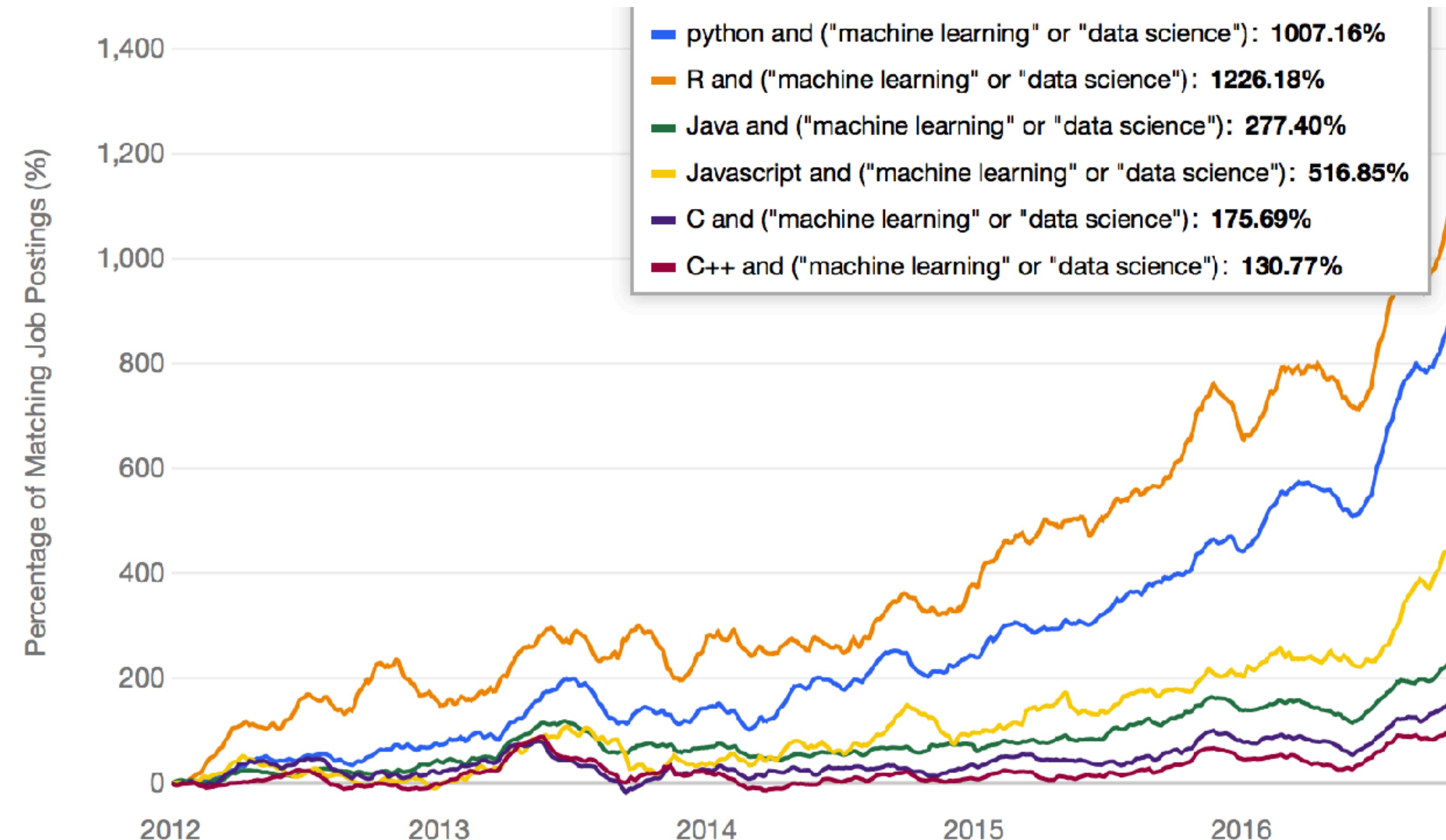
# Why?

---

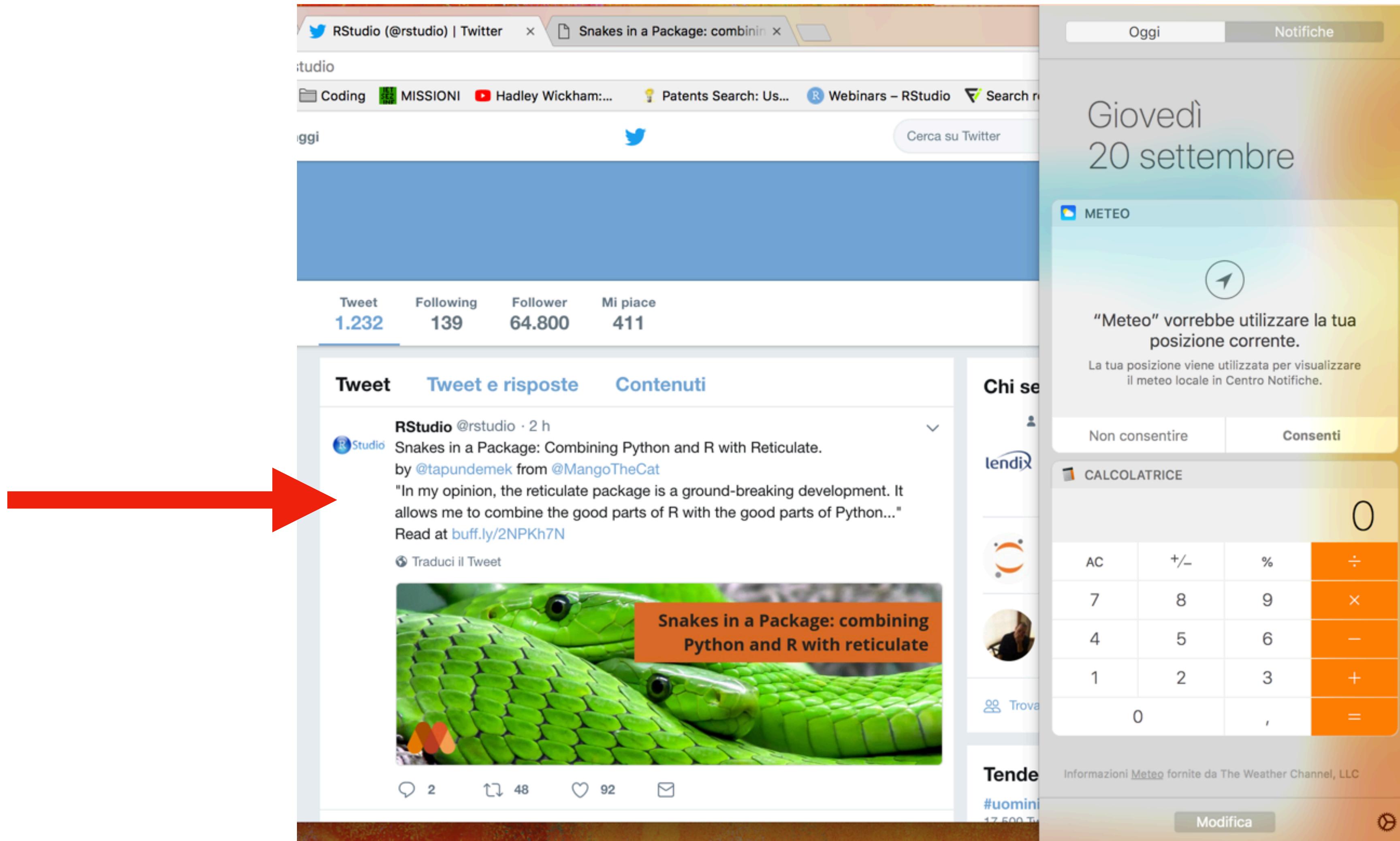


# Why?

---



# Anyway...



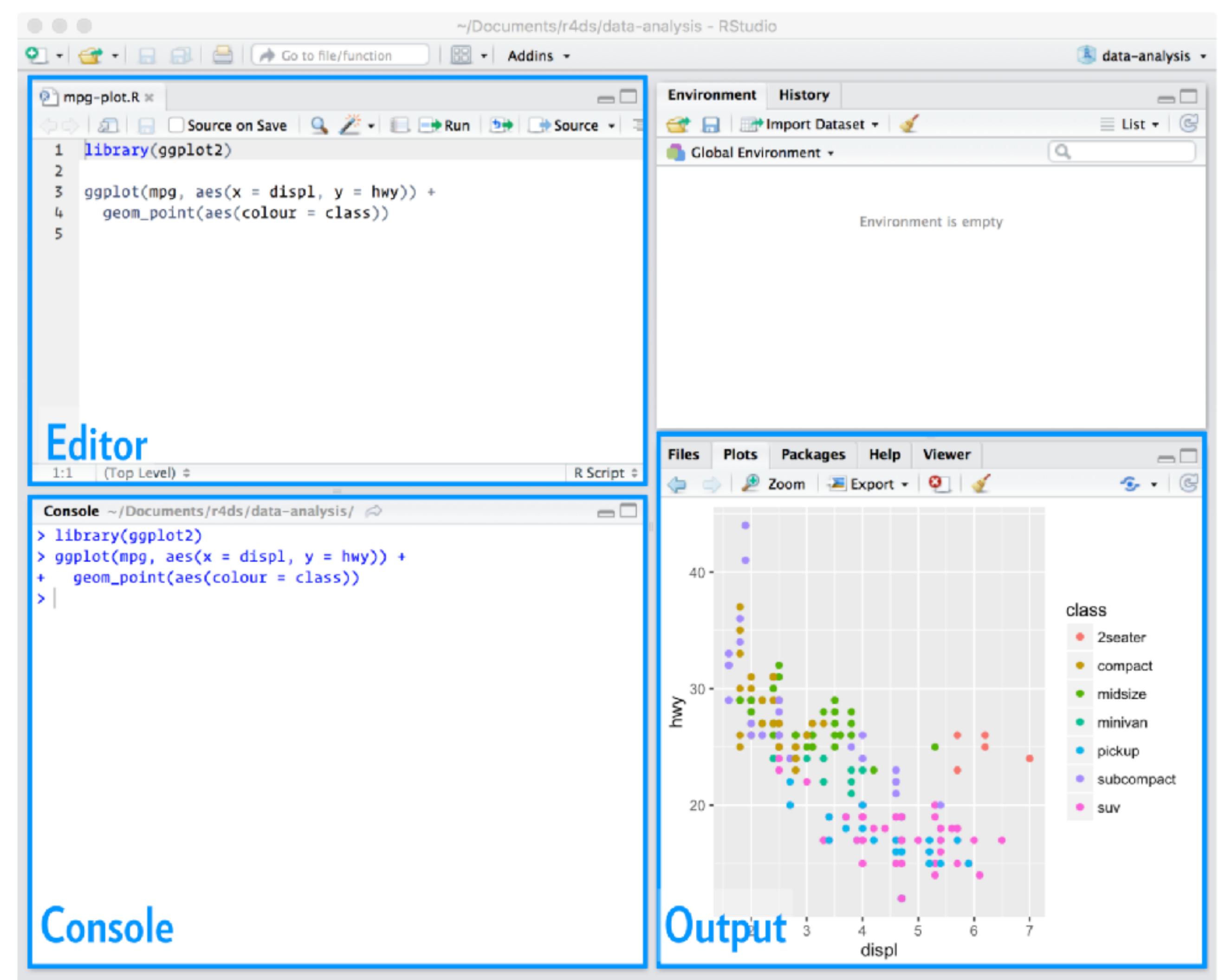
[link](#)

# Here it Comes the King...

---



- Is a free and **open source** integrated development environment (IDE) for R.
- Available in Desktop edition and Server edition, free and commercial
- Written in the **C++** programming language
- Version 1.0 was released on 1 November **2016**





---

**Mac**

---

**PC**

Last Command

↑

↑

Last Matching Command

Command + ↑

Control + ↑

Tab Completion

Tab

Tab



---

|                       | <b>Mac</b>              | <b>PC</b>               |
|-----------------------|-------------------------|-------------------------|
| Run Selection         | Command + Enter         | Control + Enter         |
| Source                | Command + Shift + S     | Control + Shift + S     |
| Source with Echo      | Command + Shift + Enter | Control + Shift + Enter |
| Move focus to console | Control + 2             | Control + 2             |
| Move focus to source  | Control + 1             | Control + 1             |



|                     | <b>Mac</b>          | <b>PC</b>           |
|---------------------|---------------------|---------------------|
| <-                  | Option + -          | Alt + -             |
| %>%                 | Command + Shift + M | Control + Shift + M |
| Comment             | Control + Shift + C | Control + Shift + C |
| Reflow Comment      | Control + Shift + / | Control + Shift + / |
| Undo                | Command + Z         | Control + Z         |
| Redo                | Command + Shift + Z | Control + Shift + Z |
| Shortcuts Reference | Option + Shift + K  | Alt + Shift + K     |



---

**Mac****PC**

---

Multiple cursors

Control + Option + **↑**

Control + Alt + **↑**

Control + Option + **↓**

Control + Alt + **↓**

Control + Option + Click

Control + Alt + Click



|                  | <b>Mac</b>                   | <b>PC</b>           |
|------------------|------------------------------|---------------------|
| Go to line       | Command + Shift + Option + G | Shift + Alt + G     |
| Find and Replace | Command + F                  | Control + F         |
| Find in files    | Command + Shift + F          | Control + Shift + F |
| Close all folds  | Command + Option + O         | Alt + O             |
| Open all folds   | Command + Shift + Option + O | Alt + Shift + O     |



|                  | <b>Mac</b>                   | <b>PC</b>           |
|------------------|------------------------------|---------------------|
| Go to line       | Command + Shift + Option + G | Shift + Alt + G     |
| Find and Replace | Command + F                  | Control + F         |
| Find in files    | Command + Shift + F          | Control + Shift + F |
| Close all folds  | Command + Option + O         | Alt + O             |
| Open all folds   | Command + Shift + Option + O | Alt + Shift + O     |

# An Incredible (Free) Learning Source

---



# Lets get practical..

<https://www.datacamp.com/community/tutorials/tidyverse-tutorial-r>

