# Course overview

### Data visualization (DSC 302) Fall 2022

## Table of Contents

Figure 1: The Red Tower/La tour rouge (Giorgio de Chirico, 1913) Source: Guggenheim

# 1 MUTUAL INTRODUCTIONS

Figure 2: Marc Chagall, Over the town (2018) Source: Wikiart

1. Why are you here?
2. What would delight you?
3. What would disappoint you?
4. Where are you headed?

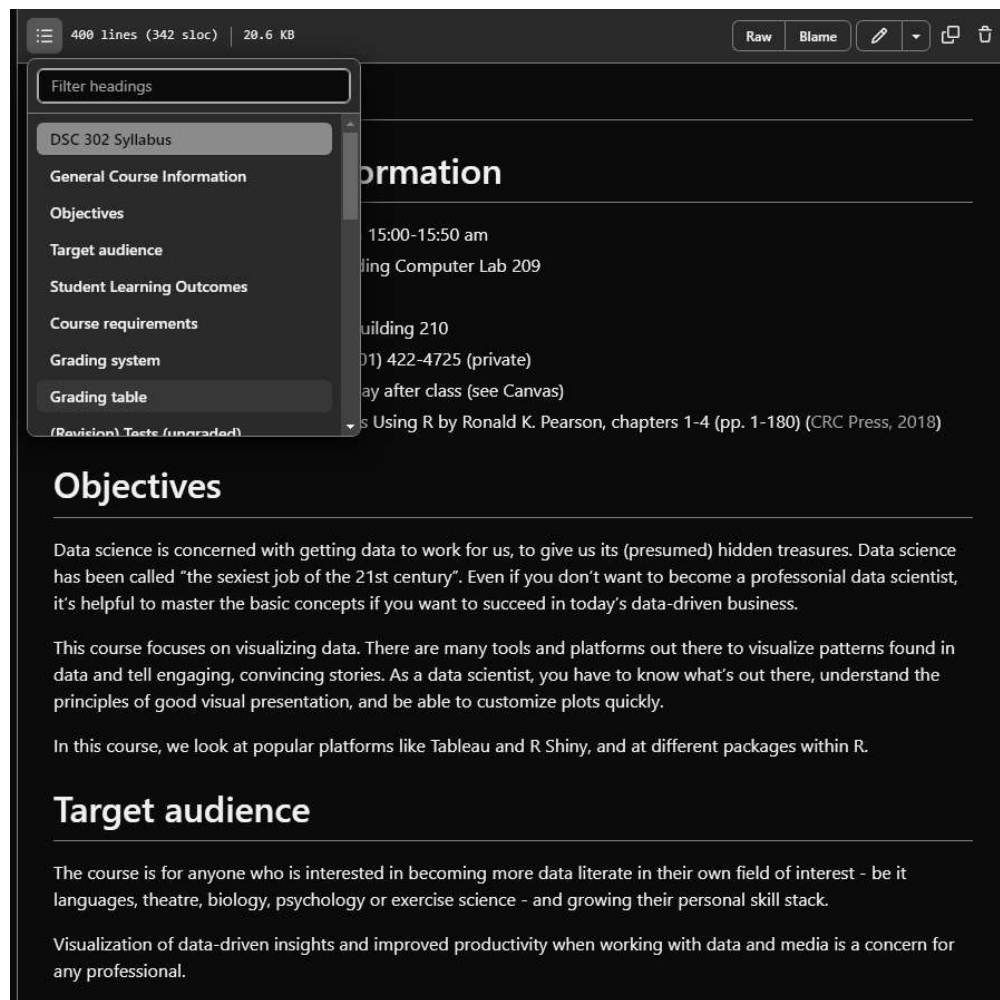# 2 COURSE SYLLABUS (on GitHub and on Canvas)

Figure 3: DSC 302 Syllabus on GitHub

- General information & standard policies
- Course information (grading, attendance)
- Schedule with dates of tests and assignments
- The GitHub repo contains course material
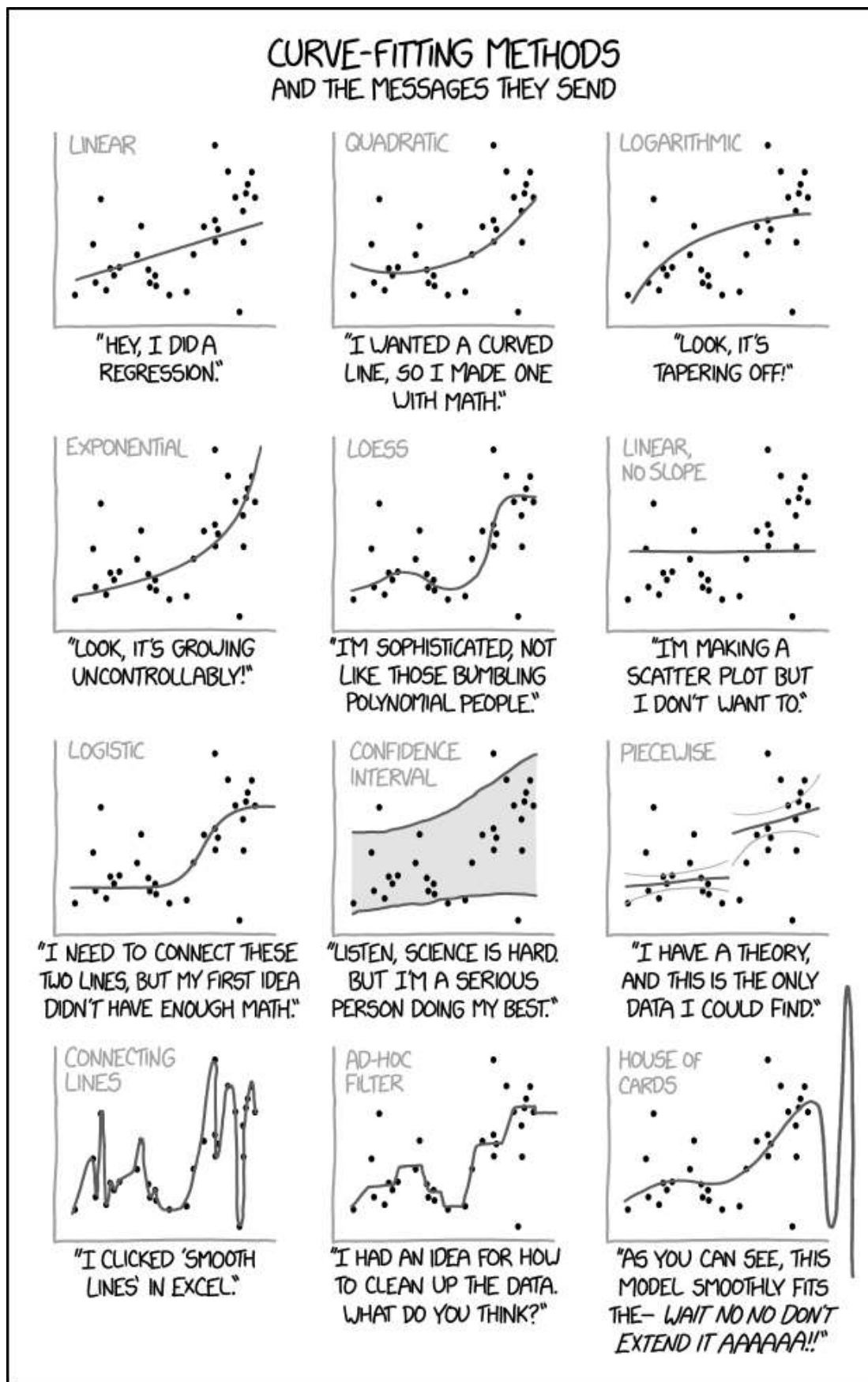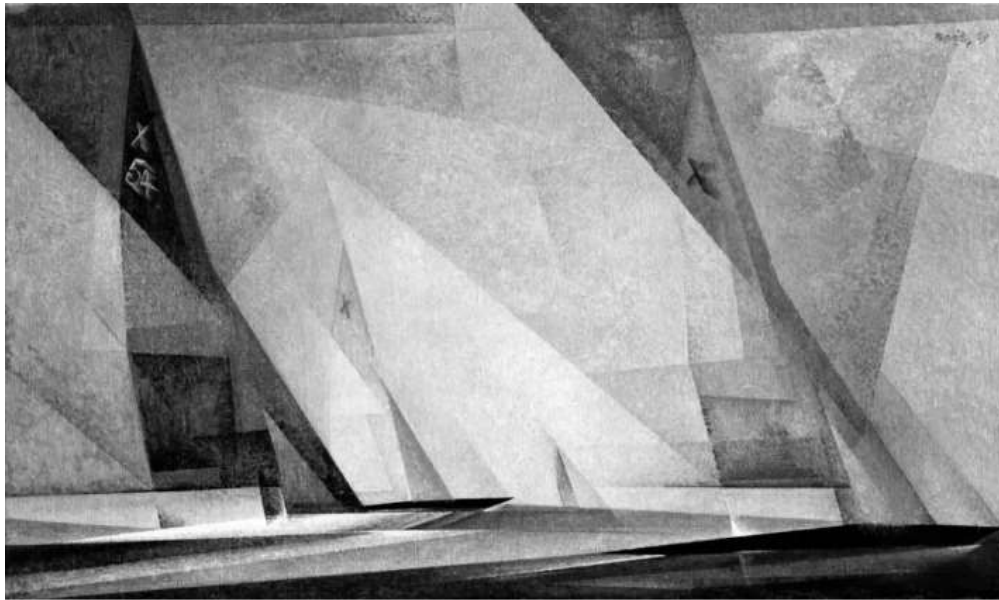
# 3 COURSE TOPICS (ILLUSTRATED)

Figure 4: Course topics

# 4 COURSE TOPICS (SPELLED OUT)



Lyonel Feininger, *Sailboats*, 1929, Detroit Institute of Arts, Detroit, MI, USA.

Figure 5: Lyonel Feininger, Sailing Boats (1929)

1. Exploratory Data Analysis (EDA) using R
2. Graphics in base R with applications
3. Working with external data (critically)
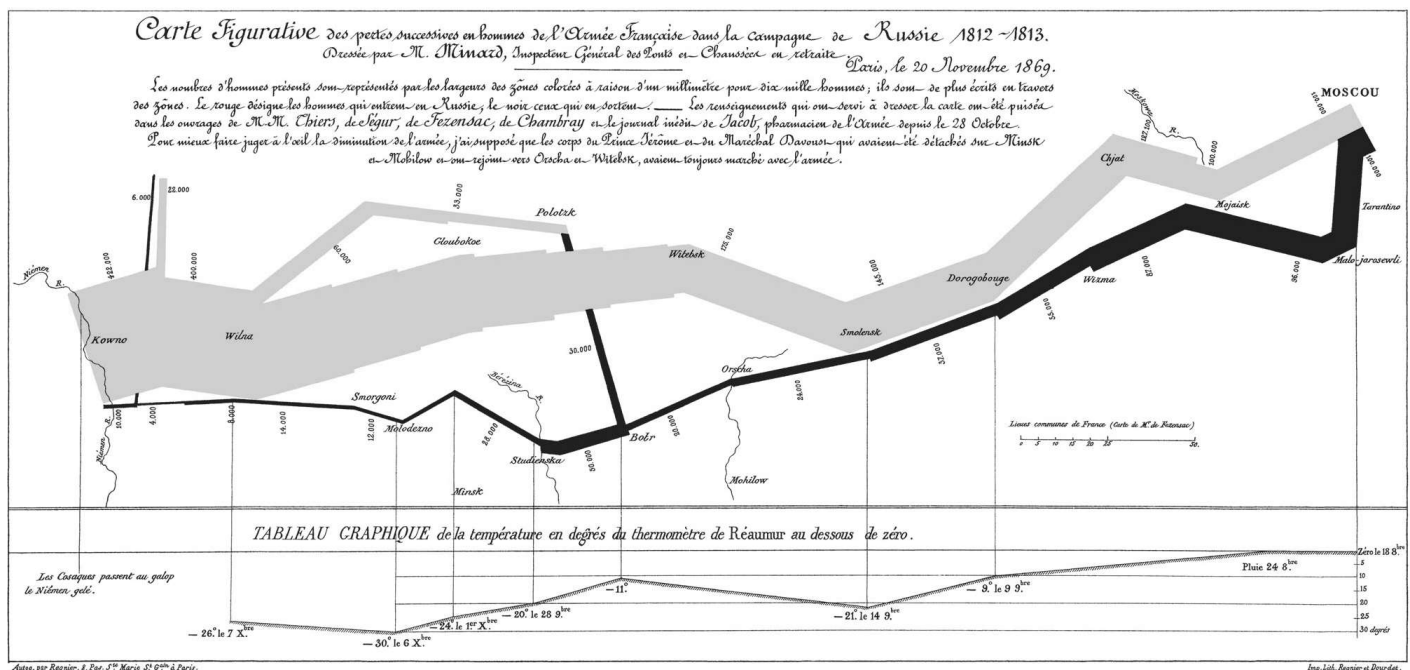
# 5 WHY "DATA VISUALIZATION"?

Figure 6: Charles Minard, Napoleon's Russian campaign 1812

- The purpose of data science is *pattern identification*
- *Visualization* happens in the head of the researcher first
- *Graphing* happens throughout, *storytelling* happens last

- The diagram by Charles Minard (1869) tells the story of Napoleon's disastrous Russian campaign in 1812 (datavizblog.com, 2013)
- Variables: army location, temperature, size over time
- Diagram type: Sankey flow diagram (many examples)
- Data type: time series (an object class, ts, in R)
- The story of this campaign is also the backstory for Tolstoy's novel "WAR AND PEACE" (Война и мир, 1867)

# 6 GET THE STORY BEHIND THE STATS

Even *The Fayetteville Observer* is trying to catch readers with data visualization / data story offers:
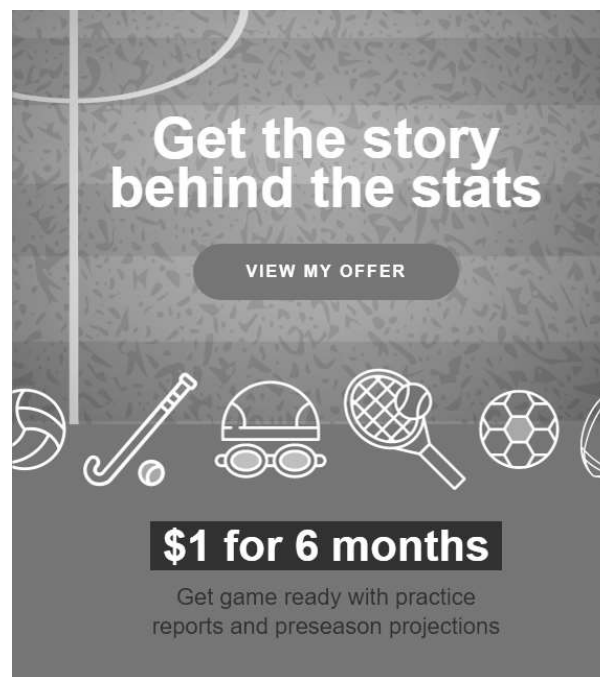


Figure 7: The Fayetteville Observer ad (Aug 5, 2022)
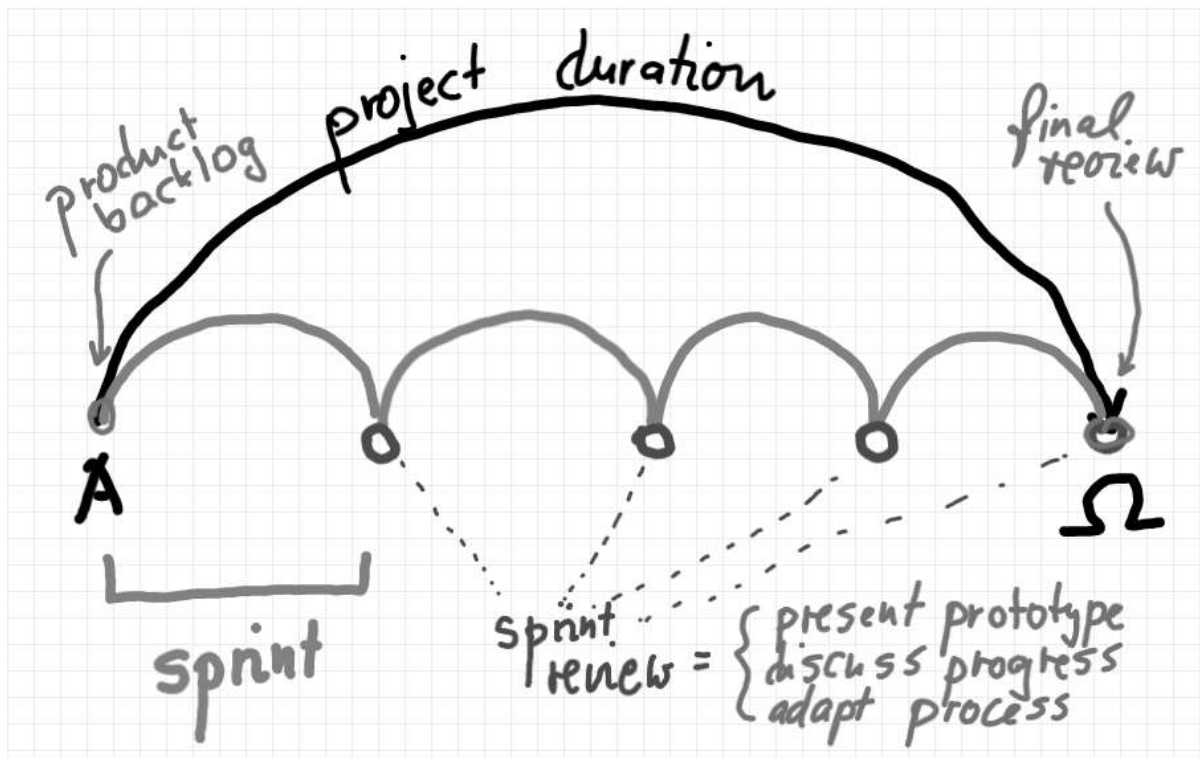
# 7 AGILE TEAM PROJECT

Figure 8: Agile (Scrum) project

The team project makes up 20% of your final grade for this course.

- What is a team project? (FAQ)
- Do you have examples for data science projects? (FAQ)
- Can you do a project as an absolute beginner? (FAQ)

**Note:** the first *sprint review* is on August 31. Use it to present your initial results (see FAQ on who to deliver, and 1st sprint review).
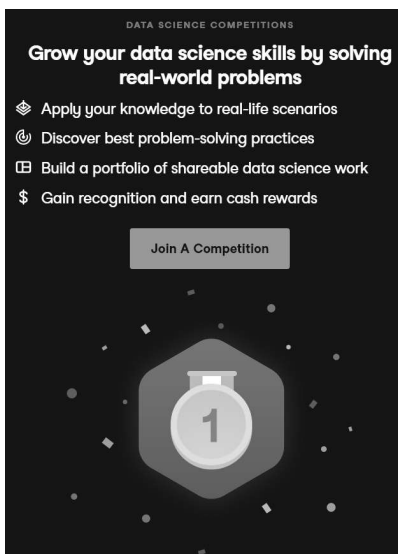
# 8 MANY PROJECT OPPORTUNITIES

Figure 9: DataCamp competition announcement

- Create an interesting data visualization
- Explore a graphics or animation package
- Solve a real-world problem
- Analyse existing visualizations
- See DataCamp projects for examples
- Explore a data visualization tool
- Visualize whale song / double up between 2 or 3 courses

- Explore any of these graphics solutions (`base`, `ggplot2` and `Shiny` are covered in this course already):



Figure 10: Source: Modern Data Visualization with R (Kabacoff, 2021)
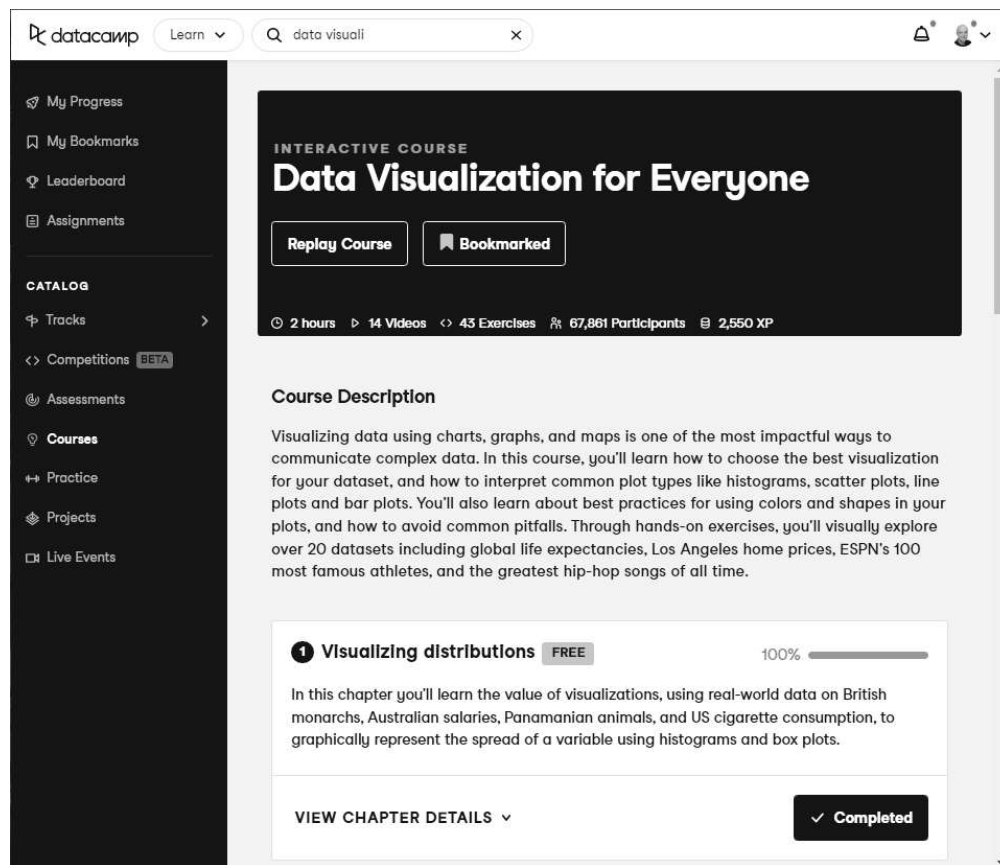
# 9 INTRODUCTION TO DataCamp

Figure 11: DataCamp course "Data Visualization For Everyone" start page

- DataCamp is a data science learning platform
- Access for you is free (classroom license)
- 9/15 assignments are DataCamp assignments
- Assignments are drawn from 5 courses
    1. Data visualization for everyone
    2. Data visualization with R
    3. Introduction to data visualization with ggplot2
    4. Building web applications with Shiny in R
    5. Introduction to Tableau
- Complete them on time to get full points
- Completed DataCamp courses can support your resume
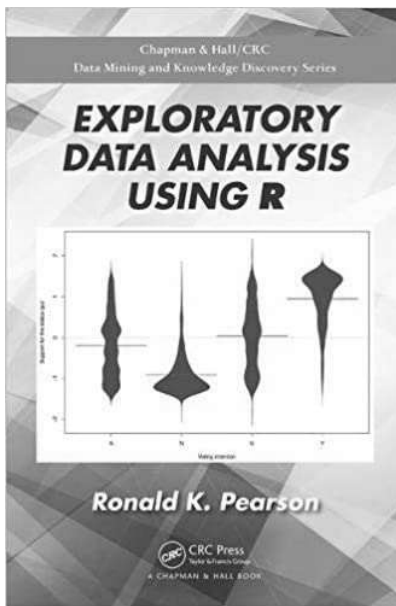
# 10 INTRODUCTION TO THE TEXTBOOK

Figure 12: Cover of EDA Using R (Pearson, 2018)

- R is *FOSS* with focus on stats and graphics
- Pearson's "EDA Using R" is extensive (563 pp.)
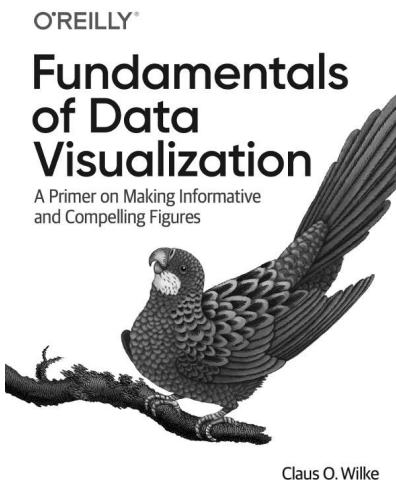- You don't have to read along but it might help

# 11 OTHER SOURCES

Figure 13: Cover of Fundamentals of Data Visualization (2019) by Claus Wilke

- Introduction to data visualization: Wilke (2019) **- in library**
- Many other tutorials and textbooks available
- The best (free) short online tutorial: Matloff's "fasteR"

- The best complete textbook: Davies' "Book of R" **- in library**
- Beware of ideologies (cp. Matloff's "TidyverseSceptic")
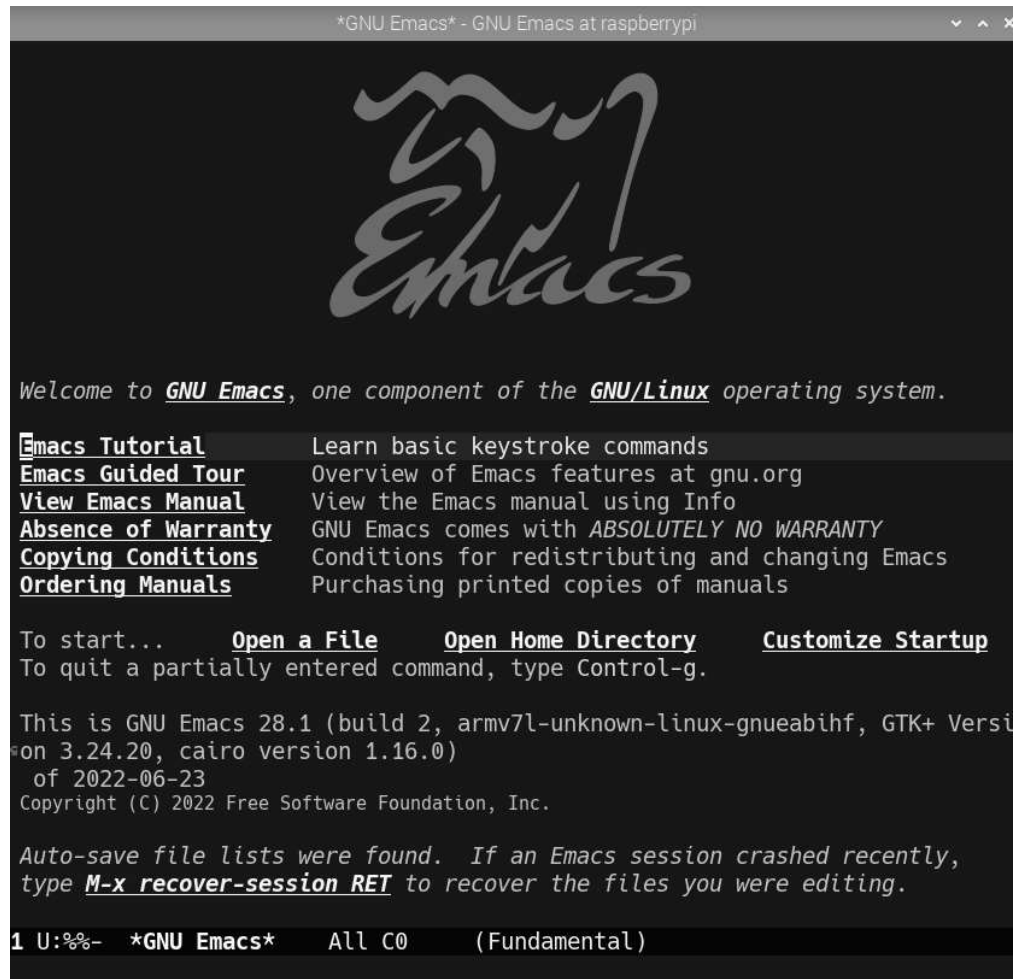
# 12 INTRODUCTION to GNU Emacs + ESS + Org-mode

Figure 14: GNU Emacs start page

- Emacs: self-documenting, extensible *FOSS* text editor
- Process, file and package management (like an OS)
- *Literate programming* environment for 43 languages
- *IDE* for R programming and *REPL* for interactive coding
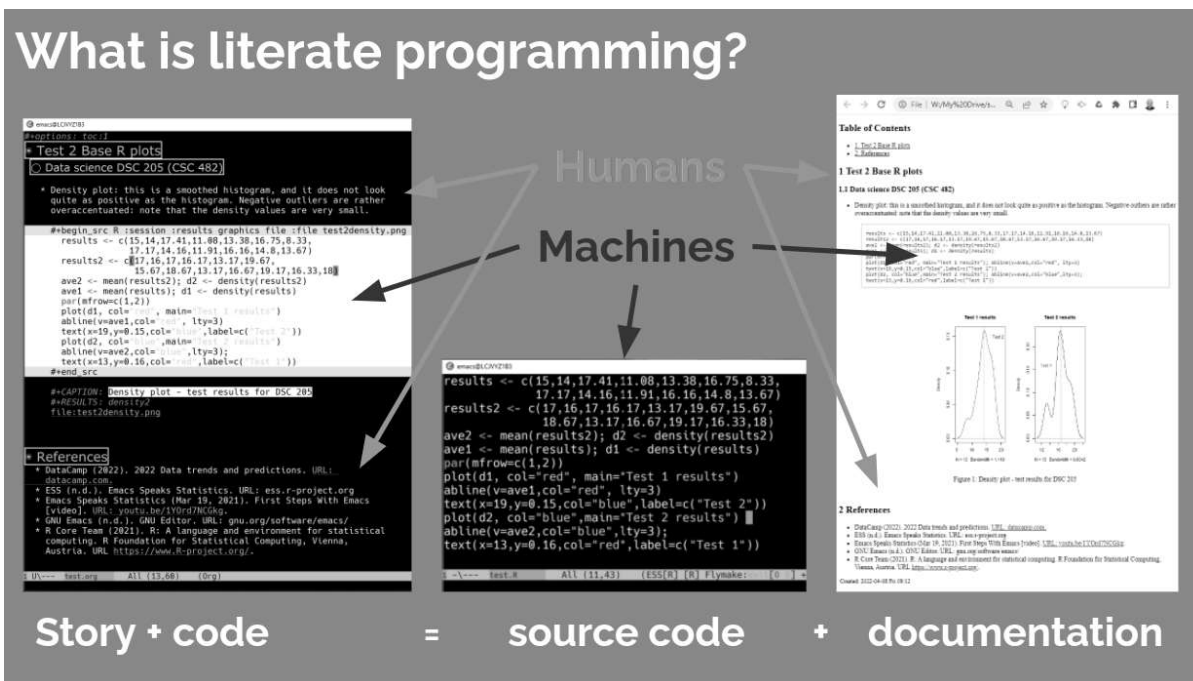
# 13 LITERATE PROGRAMMING

Figure 15: What is literate programming?

Source: "Teaching data science with hacker tools" (2022)

- Common practice among data scientists
- *Paradigm* behind interactive computing notebooks
- Useful when learning any programming language

# 14 HOME ASSIGNMENTS

- There are 15 programming assignments altogether = 10 points each, or 30% of your final grade.
- Register with DataCamp and complete the DataCamp chapter "Visualizing distributions" from the course "Data visualization for everyone" by Monday, 22 August at 3 pm (ca. 20 min).
    - Motivating visualization of data
    - Continous vs. categorical variables
    - Plot types: histograms and box plots
- Complete the Emacs on-board tutorial and upload an edited copy to Canvas by Friday, 26 August at 11 am (ca. 60 min).
    - Get comfortable with Emacs keyboard bindings
    - Learn how to create, view, edit, save files
    - Learn how to insert a time stamp automatically

# 15 TESTS (NOT GRADED)

Figure 16: Start page of the entry quiz on Canvas

- Tests have to be completed online, are timed, and have a deadline; after the deadline, you can play them an unlimited number of times
- There will be a revision quiz on Canvas every week, consisting of 5-10 multiple choice, matching and true/false questions.
- A subset of the test questions will form the final exam (20% of your final grade) - we will practice in the last week before the exam.

# 16 PRACTICE - COURSE INFRASTRUCTURE

**Useful:** take notes! Practice leads to mastery and the practice exercises will often come back to haunt you in the tests.

1. Open a browser
2. Find the GitHub repos (birkenkrahe/ds1 and /org)
3. Open the command line terminal
4. Open/close R
5. Open Emacs
6. Find the Emacs tutorial
7. Open/close R inside Emacs
8. Run R in an Org-mode file
9. Close Emacs
10. Close the command line terminal

**Note:** Class room practice completion = 10 points each for active participation.

# 17 GLOSSARY

| TERM | MEANING |
|---|---|
| Command line | aka terminal/shell to talk to the OS |
| Emacs | GNU self-extensible text editor |
| FOSS | Free and Open Source Software |
| GitHub | Software development platform |
| Git | Version control software |
| GNU | GNU's not Unix |
| IDE | Integrated Development Environment |
| "Literate Programming" | Story + code => source code + doc |
| Paradigm | A standard way of looking at things |
| R | FOSS statistical programming language |
| REPL | Read-Eval-Print-Loop |
| Repo | Code repository |
| "Tidyverse" | Popular R package bundle |
| Scrum | Agile project management method |
| Sprint review | Period to complete a prototype |
| Prototype | Intermediate (not perfect) solution |

# 18 REFERENCES

- datavizblog.com (May 26, 2013).DataViz History: Charles Minard's Flow Map of Napoleon's Russian Campaign of 1812. Online: datavizblog.com
- Davies T D (2016). The Book of R. NoStarch Press.
- Pearson R K (2018). Exploratory Data Analysis Using R. CRC Press.
- Wilke C (2019). Fundamentals of Data Visualization. O'Reilly Media. Online: clauswilke.com

Author: Marcus Birkenkrahe

Created: 2022-08-07 Sun 16:45