

DS Agenda

Agenda for CSC482/DSC205 Introduction to Advanced Data Science Spring 2022

Table of Contents

- [1. README](#)
- [2. Course introduction - w1s1 \(12-Jan\)](#)
- [3. Installing R / Windows PATH - w1s2 \(14-Jan\)](#)
- [4. Installing and setting up GNU Emacs - w2s3 \(19-Jan\)](#)
- [5. Understand Emacs Org-mode - w2s4 \(21-Jan\)](#)
- [6. Customizing Emacs \(init file\) - w3s5 \(24-Jan\)](#)
- [7. Running code in Org-mode 1 - w3s6 \(26-Jan\)](#)
- [8. Running code in Org-mode 2 - w3s7 \(28-Jan\)](#)
- [9. Org-mode lab session - w4s8 \(31-Jan\)](#)
- [10. 2022 Data Trends - w4s9 \(2-Feb\)](#)
- [11. Studying with DataCamp - w5s10 \(7-Feb\)](#)
- [12. Installing packages, using index vectors - w5s11 \(9-Feb\)](#)
- [13. Writing functions 1 - w6s13 - \(14-Feb\)](#)
- [14. Reviewing test 1, xkcd, plots - w6s14 \(16-Feb\)](#)
- [15. Guest talk - Stone Ward - w6s15 \(18-Feb\)](#)
- [16. Guest talk - Post mortem - w7s16 \(21-Feb\)](#)
- [17. Emacs recent files, Writing functions 2 - w8s17 \(28-Feb\)](#)
- [18. Change R download repo - w8s18 + 19 \(4-Mar\)](#)
- [19. Graphical devices dev.list,.Library - w9s19 \(7-Mar\)](#)
- [20. Gapminder, Pomodoro timer - Tour of ggplot - w9s20 \(9-Mar\)](#)
- [21. Test preparation / quiz 4-6 - w9s21 \(11-Mar\)](#)
- [22. Test 2 with ggplot2, org-skeleton - w10s23 \(16-Mar\)](#)
- [23. ggplot2 boxplots - w10s24 \(18-Mar\)](#)
- [24. Analysis COVID-19 DataCamp project - w11s25/26 \(28/30-Mar\)](#)
- [25. UpdateR, scatterplots - w11s27 \(1-Apr\)](#)
- [26. Faceting - w12s28 \(4-Apr\)](#)
- [27. Time series plots - w12s29 \(6-Apr\)](#)
- [28. Scale and value transformations - w13s30 \(11-Apr\)](#)
- [29. Boxplots and ridge plots - w13s31 \(13-Apr\)](#)
- [30. EDA with categorical data - w14s32 \(20-Apr + 22-Apr\)](#)
- [31. Review quiz 7-9 - w15s34 \(25-Apr\)](#)
- [32. C++ and R \(Wyatt\) - w15s34 \(27-Apr\)](#)
- [33. Data science on the command line - w15s35 \(2-May\)](#)
- [34. Last Rites / Final exam review - w16s36 \(4-May\)](#)
- [35. References](#)

1 README

This file contains the agenda overview (what I had planned), the objectives (what we managed to do) and (much of the) content of each taught session of the course. I want to avoid splitting the content up over many files - so that you have to navigate as little as possible (like a book)!

The companion file to this file, less structured and with the captain's log, is the [notes.org](#) file.

2 Course introduction - w1s1 (12-Jan)

2.1 Welcome



- Aspirations - changes spring 2022
- Ambitions - program 2021-2023
- Antagonization - new data science credo
- Syllabus - this course
- DataCamp assignments
- GNU Emacs Org-mode

"After a course is launched, we don't consider it to be complete: the launch is just the start of data collection." Richie Cotton, DataCamp

2.2 Syllabus



- [Syllabus in Schoology](#)
- [Syllabus in GitHub](#)
- [Schedule in GitHub](#)

2.3 Aspirations (Changes in Spring 2022)

Cp. Good-bye fall 2021

FALL 2021	SPRINT 2021
Base R (stick shift) instead of "TidyVerse" (automatic)	Adding the "Tidyverse"
Use of interactive notebooks (literate programming!)	Intro to RStudio IDE and Emacs
Use GitHub as a code and materials repository	GitHub repo
Create lots of (ungraded) tests	Graded quizzes and tests
Use of DataCamp assignments	DataCamp assignments
Avoid mathematics as much as possible	No math
Reuse tests for the final exam	Reuse quizzes for final exam
Let students pick their own projects	No projects (only optional)

2.4 Ambitions (DS program 2021-2023)

CLASS	CODE	TERM	Topics
Data Science Tools and Methods	DSC 101	Fall 2021	R, Basic EDA, Base R
Introduction to Advanced Data Science	DSC 205	Spring 2022	R, Advanced EDA, Tidyverse, shell
Database Theory and Applications	CSC 330	Spring 2022	SQL, SQLite
Operating Systems	CSC 420	Spring 2022	Bash, awk, sed, regular expressions
Applied Math for Data Science	DSC 482/MTH 360	Fall 2022	Probability, Statistics + R
Data Visualization	DSC 302	Fall 2022	D3, Processing, Javascript, Bokeh
Machine Learning	DSC 305	Spring 2023	Predictive algorithms, neural nets
Digital Humanities	CSC 105	Spring 2023	Data science applications

2.5 DataCamp

 Intermediate R Conditions and Control Flow Chapter	Team	Active	Jan 31, 15:00 CST	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Intermediate R Loops Chapter	Team	Active	Feb 7, 15:00 CST	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Intermediate R Functions Chapter	Team	Active	Feb 14, 15:00 CST	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Intermediate R The apply family Chapter	Team	Active	Feb 21, 15:00 CST	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Intermediate R Utilities Chapter	Team	Active	Feb 28, 15:00 CST	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Introduction to the Tidyverse Data Wrangling Chapter	Team	Active	Mar 7, 15:00 CST	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Introduction to the Tidyverse Data visualization Chapter	Team	Active	Mar 14, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Introduction to the Tidyverse Grouping and summarizing Chapter	Team	Active	Mar 28, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Introduction to the Tidyverse Types of visualizations Chapter	Team	Active	Apr 4, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Exploratory Data Analysis in R Exploring Categorical Data Chapter	Team	Active	Apr 11, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Exploratory Data Analysis in R Exploring Numerical Data Chapter	Team	Active	Apr 20, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Exploratory Data Analysis in R Numerical Summaries Chapter	Team	Active	Apr 25, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View
 Exploratory Data Analysis in R Case Study Chapter	Team	Active	May 2, 15:00 CDT	<div style="width: 0%; height: 10px; background-color: #ccc;"></div>	0	0	0%	View

- Why are we using it?
- How are we using it?
- What will you have to do?

2.6 Antagonization

A new credo.

“Getting it right is crucial when people’s lives are affected.” -Jonathan Steinhart



Figure 4: Lego fencing (Source: Unsplash)

2.7 What's next?



- See schedule:
 - install R / Emacs IDE - may do this together
 - Entry quiz (by Tue 18 Jan) - you should get > 50%
- Watch online lecture on "Systems" (to be published)
- Online followup notes ([notes.org](#) in GitHub)
- See you Friday 14-Jan online!
- Hopefully Wednesday 19-Jan in class!

3 Installing R / Windows PATH - w1s2 (14-Jan)

3.1 Overview

HOW	WHAT
Practice	Install R from CRAN
	Set PATH environment variable
	Test R in terminal and GUI
Install GNU Emacs + ESS (FAQ)	
	Set PATH environment variable
	Test R in Emacs
	Set .emacs init file
	Create Org file
	Run R code blocks in an Org file

3.2 Objectives

- [X] Install R
- [X] Set PATH environment
- [X] Test R in terminal and GUI
- [] Install GNU Emacs
- [] Test R in Emacs

4 Installing and setting up GNU Emacs - w2s3 (19-Jan)

4.1 I'm back

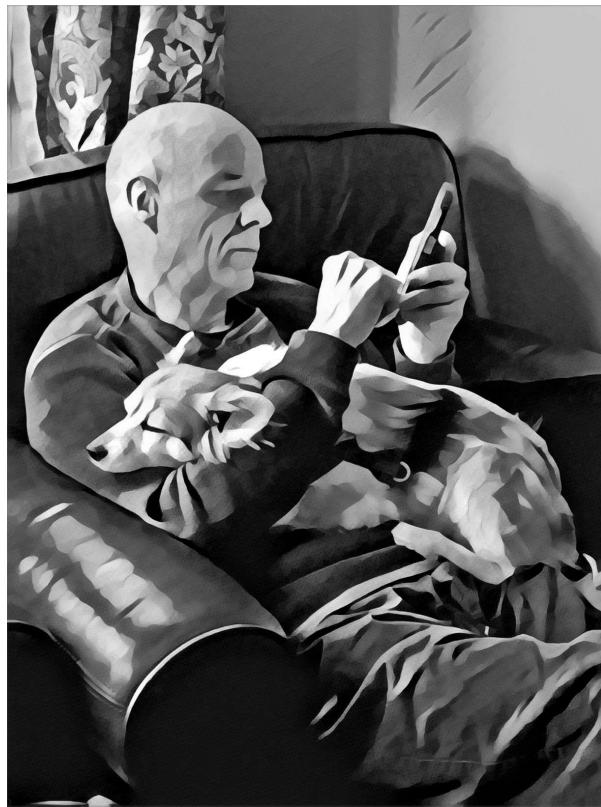


Figure 6: "I'm back, baby."

4.2 Overview

HOW	WHAT
Review	Entry quiz Quiz 1 + feedback + discussion
Practice	Install GNU Emacs + ESS (FAQ) Set PATH environment variable

HOW	WHAT
Test R in Emacs	
Set .emacs init file	

4.3 Objectives

- [X] Install GNU Emacs + ESS
- [X] Set PATH environment to run R in Emacs
- [X] Test R in Emacs (however, see course FAQ)
- [] Configure Emacs

4.4 Next

- Create Emacs Org file
- Run R code blocks in an Org file
- DataCamp assignments beginning soon!

5 Understand Emacs Org-mode - w2s4 (21-Jan)

5.1 Overview

HOW	WHAT
Lecture/Demo	GNU Emacs <u>Org-mode</u>
Practice	GNU Emacs Tutorial (gh)
Homework	Set <code>emacs</code> init file
	Create <code>.org</code> file
	Run code in an <code>.org</code> file

5.2 Objectives

- [X] Understand what Org-mode is and what it's for
- [] Create an `.emacs` init file for GNU Emacs
- [] Create an Org file
- [] Run a code block in your Org file

5.3 Next

- Create Emacs Org file
- Run R code blocks in an Org file
- DataCamp assignments beginning soon

6 Customizing Emacs (init file) - w3s5 (24-Jan)

6.1 Overview

HOW	WHAT
Review	Quiz 2
Lecture/Demo	GNU Emacs <u>Org-mode</u> (Part 2) <u>New: video playlist</u>
Practice	GNU Emacs Tutorial cont'd (gh)
- Package manager	M-x package-list-packages RET
- Start R shell in Emacs	M-x R (R must be installed & in the PATH)
- Add init file	.emacs sample file (GitHub)
<u>Assignment</u> ¹	Set <code>emacs</code> init file
<u>Assignment</u>	Read 2022 Data trends and predictions Put your summary thoughts in an .org file Check the FAQ "How should you read?"

6.2 Objectives

- [X] Create an .emacs init file for GNU Emacs
- [] Create an Org file
- [] Run am R code block in your Org file

6.3 Reading assignment

- [Read "2022 Data trends and predictions"](#) (DataCamp, 2022).
- Prepare for discussion in class:
 - Which quantitative and which qualitative predictions were made?
 - What do you think how valid these predictions are?
 - Put your thoughts in an Org-mode file (filename = YourName.org)
 - Upload your submission to [assignment/2022_predictions](#) on GitHub

To identify yourself, use the #+AUTHOR: option. You can see how

this works from the options in the header of this README.org file.

There is no upper or lower limit on the number of words. The main point is to create a proper Org-mode file.

6.4 Next

- Create Org-mode file with R code in it and run it
- Org-mode assignment
- DataCamp assignments beginning soon (due Jan 31)

Assignments / DSC 205 Introduction to Advanced Data Science ▾

+ Create Assignment

ACTIVE PAST DUE ARCHIVED

Active Assignments

Filter By Type ▾

Search assignments...

TITLE ▾	ASSIGNEES ▾	STATUS	DUE BY ▾	C ▾	A ▾	CR ▾	DESCROLL >
Intermediate R Conditionals and Control Flow Chapter	Team	Active	Jan 31, 15:00 CST	0	8	0%	View
Intermediate R Loops Chapter	Team	Active	Feb 7, 15:00 CST	0	8	0%	View
Intermediate R Functions Chapter	Team	Active	Feb 14, 15:00 CST	0	8	0%	View
Intermediate R The apply family Chapter	Team	Active	Feb 21, 15:00 CST	0	8	0%	View
Intermediate R Utilities Chapter	Team	Active	Feb 28, 15:00 CST	0	8	0%	View

Figure 7: DataCamp assignments

7 Running code in Org-mode 1 - w3s6 (26-Jan)

7.1 Overview

HOW	WHAT	Link
Preview	DataCamp course "Intermediate R"	datacamp.com
Demo	Creating an Emacs Org-mode file with code and run it	README.org
Practice	Create Org-mode file with an R code block	

7.2 Objectives

- [X] Understand DataCamp assignment 1
- [X] Create an Org file
- [X] Run an R code block in your Org file

7.3 Next

- Submit Org-mode assignment in Schoology
- DataCamp assignments due Jan 31

The screenshot shows a web-based assignment management system. At the top, there's a header bar with the text "Assignments / DSC 205 Introduction to Advanced Data Science" and a "Create Assignment" button. Below the header are three tabs: "ACTIVE" (which is selected), "PAST DUE", and "ARCHIVED". A search bar labeled "Search assignments..." is located below the tabs. To the right of the search bar is a "Filter By Type" dropdown menu. The main area displays a table of assignments:

TITLE	ASSIGNEES	STATUS	DUE BY	C	A	CR	DESCROLL >
Intermediate R Conditionals and Control Flow Chapter	Team	Active	Jan 31, 15:00 CST	0	8	0%	<button>View</button>
Intermediate R Loops Chapter	Team	Active	Feb 7, 15:00 CST	0	8	0%	<button>View</button>
Intermediate R Functions Chapter	Team	Active	Feb 14, 15:00 CST	0	8	0%	<button>View</button>
Intermediate R The apply family Chapter	Team	Active	Feb 21, 15:00 CST	0	8	0%	<button>View</button>
Intermediate R Utilities Chapter	Team	Active	Feb 28, 15:00 CST	0	8	0%	<button>View</button>

Figure 8: DataCamp assignments

8 Running code in Org-mode 2 - w3s7 (28-Jan)

1. We continue where we left it last Wednesday
2. Fixing the .emacs problem on Windows lab computers
3. Change of some deadlines - to finish basic Emacs training

Upcoming · 26

Add Event

Friday, January 28, 2022

 CREATE AND RUN R IN AN
EMACS ORG FILE AND IN THE
SHELL (in-class exercise) 11:59
pm

Monday, January 31, 2022

 Read trend report, put your
thoughts in an Emacs Org-
mode file 3:00 pm

Wednesday, February 2, 2022

 DataCamp assignment 1 3:00
pm

Figure 9: deadline changes in Schoology

4. Finish (expanded) Org-mode assignment
5. Submit results to Schoology.

9 Org-mode lab session - w4s8 (31-Jan)



Figure 10: Teaching Emacs on Dagobah

We will hold a special lab session tomorrow, Monday 31 January 3-3.50 PM, to sort out any issues related to Emacs and R. Bring your own PC to the session, or work on a lab desktop. I will spend the time going round to make sure that you can

- Install/ open / use the Emacs editor
- Create, run and tangle Org-mode files with R code
- Install / use the R programming language
- Understand the recent program assignments

The necessary steps are also demonstrated [in this tutorial video playlist](#).

We will continue with our regular program on Wednesday, 2nd February at 3 PM - a short quiz will be available before.

For those who know or can do all of this already: here's a [second challenge](#) (with solution) to practice while I sort others out.

9.1 What's next

- Deadline for 1st DataCamp assignment is looming ([Wed 2 Feb 3pm](#))
- Scenario building for "Data Trends and Predictions 2022" report ([assignment](#)) - think about the 2 most important dimensions & watch this video about [scenario planning](#)
- Complete **quiz 3** including a **poll** on the prediction report before class
- Check out the [webinar recording](#) with DataCamp luminaries (panel)
- Use the breathing space to complete the Emacs tutorial (c-h t)

10 2022 Data Trends - w4s9 (2-Feb)

We meet today at 3-3.5- PM in the seminar room Lyon 106 - this room is directly adjacent to 104, our usual lab. We'll discuss the DataCamp 2022 trend report. The quiz will be available before end of the week. The planned first test (in class) will take place next Wednesday instead. ([Schoolology Update](#))

10.1 Overview

HOW	WHAT
Discussion	DataCamp 2022 report on Data Trends
Groupwork	Data science scenario planning (video)

10.2 Objectives

- [X] Understand the implications of the 2022 DataCamp trend report
- [X] Understand and apply the scenario planning technique

10.3 Next

- Quiz 3 - Conditionals and Control Workflow (DataCamp review)
- Test 1 (Friday 11 Feb 3 PM)
- Interactive R notebook - Writing functions

11 Studying with DataCamp - w5s10 (7-Feb)

11.1 Overview

HOW	WHAT
Review	Quiz 3 - Relational and logical operators How to study R with DataCamp
Preview	While and For Loops
Lecture	Writing functions in R
Test info	Test 1 on Friday 11 Feb 3.05-3.50 pm

11.2 Objectives

- [X] Review quiz 3 & how to study with DataCamp
- [X] Understand test conditions (Friday 11 Feb)
- [] Understand how to write functions in R (lecture)

11.3 Test 1 info

- Online in Schoology
- Entry quiz and Quiz 1-3 are not visible during the test
- The 10 hardest questions of entry quiy + quiz 1-3 (< 50%)
- 10 new questions
- Maximum time = 45 min

11.4 Next

- Interactive R notebook - loop problems
- Test 1 (Friday 11 Feb 3 PM)

12 Installing packages, using index vectors - w5s11 (9-Feb)

12.1 Overview

HOW	WHAT
Review	While and For loops
Lecture	Writing functions in R (part 1)

12.2 Objectives

- [X] Org-mode PROPERTY "shebang" stuff (meta data)
- [X] Review: install packages and loading datasets
- [X] Understanding and using index vectors

12.3 Next

- Test 1 (Friday 11 Feb 3-3.50 PM)
- Matthew Stewart, Stone Ward (Friday 18 Feb 3-3.50 PM)

13 Writing functions 1- w6s13 - (14-Feb)

13.1 News

- [2022 Data analytics competition \(accounting data\)](#)
- Matthew Stewart, Stone Ward (Fri 18 Feb 3-3.50 PM) in Derby 209

13.2 Overview

HOW	WHAT
Class assignments	How do they work?
Practice Class assignments	Write a hello world function Installing loading packages .Rprofile configuration file
Review	Writing functions (DataCamp)
Interactive Lecture	Writing functions in R (part 2) Statistical functions in R

13.3 Objectives

- [X] Mark guest talk in your calendar (Fri 18-Feb) Derby 209
- [X] Understand how "class assignments" work
- [X] Complete a couple of class assignments
- [] Practice: install packages and loading datasets
- [] Review DataCamp chapter on writing functions

13.4 How do class assignments work?

- In-class assignments are **10%** of your total grade
- They are labeled **class assignments** in the Schoology gradebook
- You get the points if you attend and participate **actively**
- If you check your phone instead, you're **not** active
- If you could not attend (with a good excuse), submit **late**
- Submit an **Org-mode file**, not a screenshot

13.5 Next

- Wednesday: Review of test 1
- See some fun plotting techniques

14 Reviewing test 1, xkcd, plots - w6s14 (16-Feb)

14.1 News

- Eliminated some DataCamp assignments
- Remaining assignments mostly bi-weekly
- Emacs package of the week: xkcd

14.2 xkcd - life is too serious sometimes

- Package is pre-installed (list: `M-x package-list-packages`)
- `M-x xkcd` opens current comic
- `o` in xkcd mode opens browser with current topic
- `C-h ? m` opens full mode description

14.3 Overview

HOW	WHAT
Review	Hello function
	Test 1 - first month of class
How to make up for bad test results	Complete a mini-project

14.4 Objectives

- [] Review: Hello function
- [] Review: results of test 1
- [] Learn how to plot a density distribution and the mean
- [] Understand factor vectors
- [] Master Vector element extraction
- [] Understand the difference: Emacs Org-mode, ESS, and Base R
- [] Understand R comments
- [] Understand NA
- [] Understand the difference: object, storage class, data type
- [] Understand the help available in and outside of R
- [] Understand print and paste
- [] Understand vectorization
- [] Understand purpose and properties of interactive notebooks

14.5 CHALLENGE: Write a hello function with your name as an argument

- You already learnt how to write a `hello()` function without arguments. Write a function that takes your name as an argument and prints "Hello, [your name]". Write and test the function in the same code block.

```
hello <- function(name) {
  print(paste("Hello, ", name))
}
hello(name="Marcus")
```

[1] "Hello, Marcus"

- Another solution, this time with two arguments.

```
hello2 <- function(fname, lname) {
  print(paste("Hello, ", fname, lname, "!"))
}
hello2(fname="Marcus", lname="Birkenkrahe")
```

[1] "Hello, Marcus Birkenkrahe !"

14.6 Lab 104 Emacs check

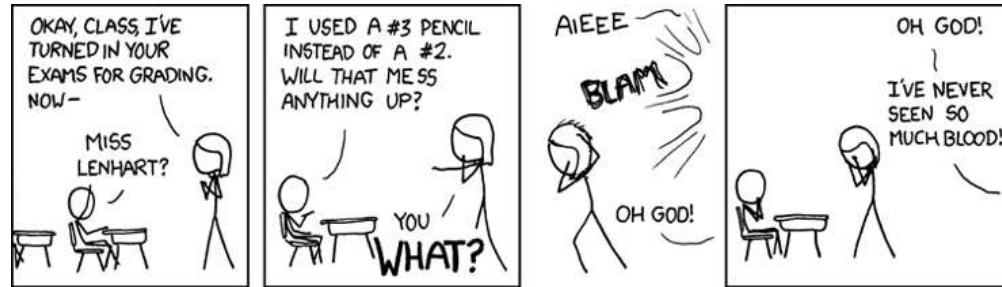
- First thing, when you sit down at your desktop in the computer lab, open Emacs, write a code block in an Org-mode file (`test.org`), and try to run it:

```
str(mtcars)
```

- If it does not work but instead complains about missing `org-babel` whatever, you need to install a `.emacs` file in the `$HOME` directory.
- Download the file or its content from <https://tinyurl.com/lyonemacs>. Make sure the file has the right name, then restart Emacs and run the code block again.
- You unfortunately need to do this any time you sit at a computer in the lab you have not sat at before.
- To make things easier, you could also put a `.emacs` file in your GDrive and download it in one go.

14.7 Test review

14.7.1 Paper vs Screen



Never again! Preparing such a test on paper and grading it while allowing for partial credit is a nightmare: future tests will be online in Schoology!

14.7.2 Test 1 results

- The test results are OK (average 70%). Better next time!

Statistics

# of Grades	14	Average	13.06 (65.29%)
Max Points	20	Standard Deviation	4.3 (21.5%)
Highest Grade	17.17 (85.85%)	Median	14.08 (70.4%)
Lowest Grade	0 (0%)	Mode	N/A (N/A)

Figure 12: Test 1 results (Schoology)

```
results <- c(15,14,17.41,11.08,13.38,16.75,8.33,
           17.17,14.16,11.91,16.16,14.8,13.67)
```

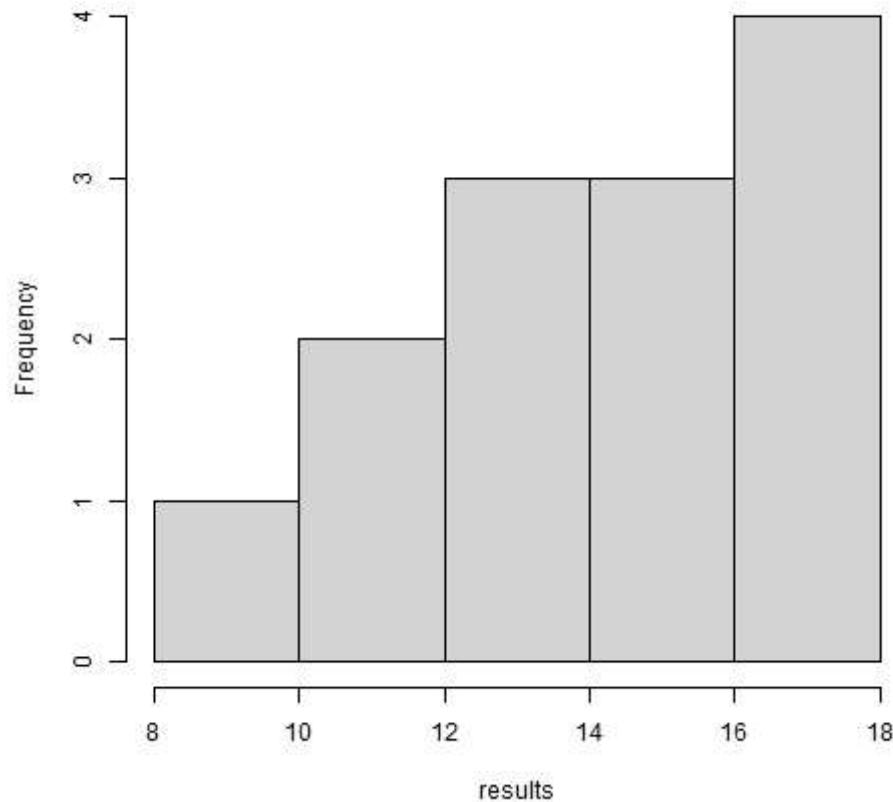
- When checking the stats with R, I find different results. Why?²

```
paste("Sample:",length(results))
paste("Standard deviation:", sd(results))
paste("Average:", 100*mean(results)/20)
summary(results)
```

```
[1] "Sample: 13"
[1] "Standard deviation: 2.59571120632991"
[1] "Average: 70.7"
Min. 1st Qu. Median Mean 3rd Qu. Max.
 8.33   13.38   14.16   14.14   16.16   17.41
```

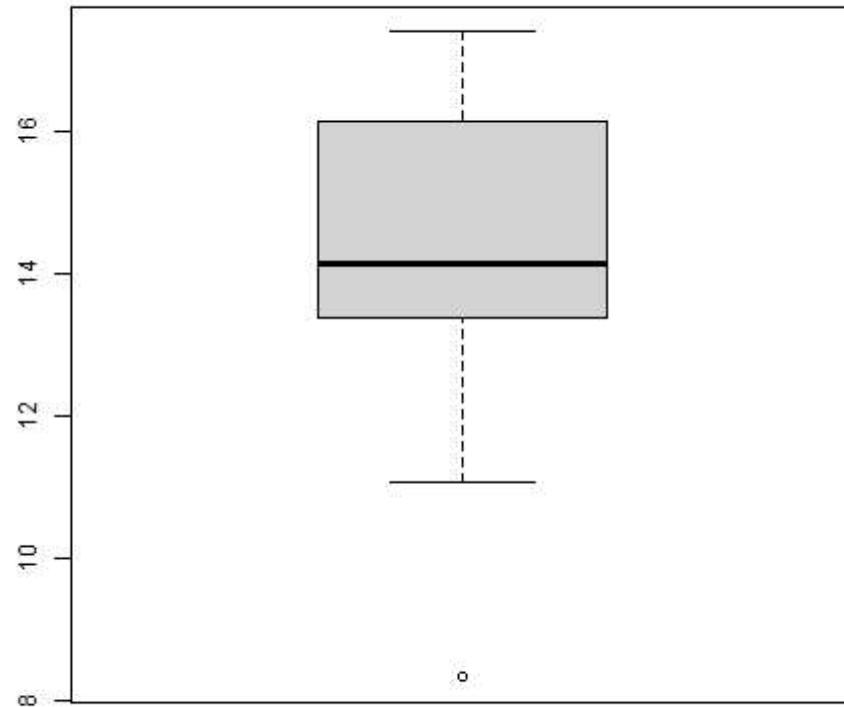
- Let's make some plots: histogram, boxplot and density plot.
- Fetch the vector from GitHub and run the code in Emacs.
- Histogram. Demonstrates the fact that almost the entire course but one is above 50% (= pass). Looks more positive than the whole truth, because the x-axis ends with the maximum result achieved, and not with the maximum points available (20).

```
hist(results, main="Test 1 results, DSC 205 Spring 2022")
```

Test 1 results, DSC 205 Spring 2022

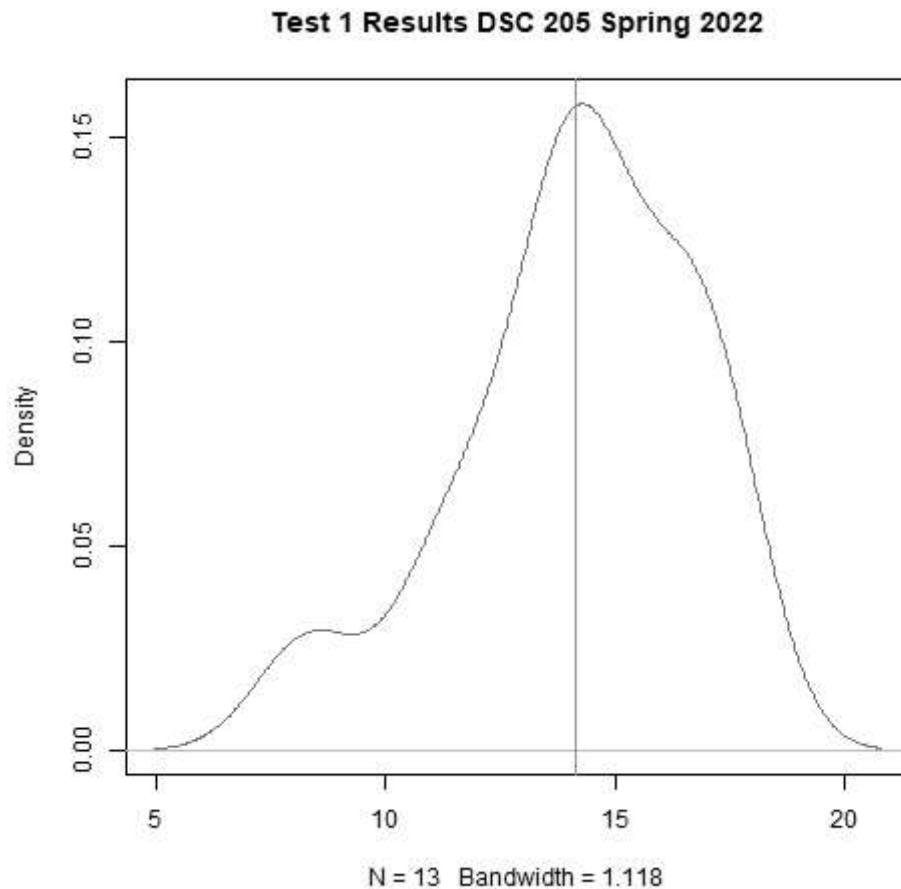
- Boxplot: this graph is deceptively positive, because it doesn't show the maximum points (20) but only the maximum achieved points. The "whiskers" correspond to the outliers, and the thick black line is the median (the middle value).

```
boxplot(results, main="Test 1 results, DSC 205 Spring 2022")
```

Test 1 results, DSC 205 Spring 2022

- Density plot: this is a smoothed histogram, and it does not look quite as positive as the histogram. Negative outliers are rather overaccentuated.

```
ave <- mean(results)
med <- median(results)
d <- density(results)
plot(d, col="steelblue", main="Test 1 Results DSC 205 Spring 2022")
abline(v=ave, col="red")
abline(v=med, col="green")
```



14.7.3 Analysis - feedback and action points

- Test 1 can now be played an unlimited number of times. I will add feedback to all new questions by the end of today.
- If you didn't play the other quizzes until you reached 100%, you had it coming. (My question: why wouldn't you do that?)
- What surprised me most was that many of you did not use the available time. However, I have not (yet) been able to correlate test time and test success (it's a project).
- Plots: I'd like the histogram and the density plot (a smoothed histogram) to peak more to the right, and for the boxplot to be smaller and higher up.
- See also: "I can teach it to you but I cannot learn it for you"
- Questions:
 - How did you study for this test?
 - If you didn't perform well, what will you change?
 - What can I do to help you help yourself?
- Changes to be applied in future quizzes/tests:
 - Fewer multiple choices (max. 4)
 - Announce if a question has > 1 answer (and/or how many)
 - Try to avoid having > 1 test on the same day

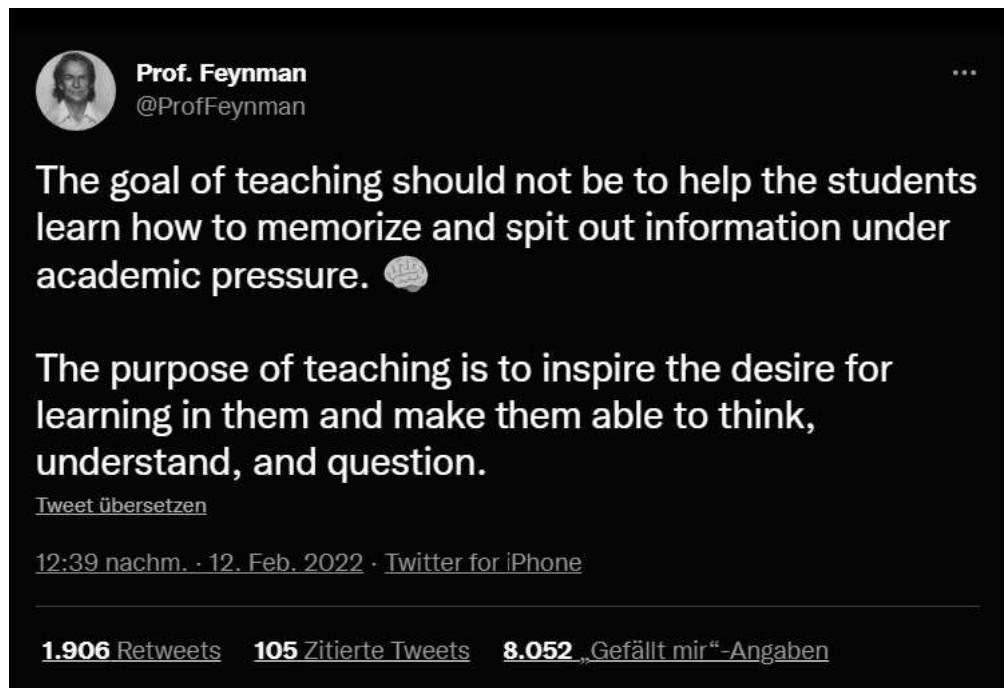


Figure 16: Feynman (via Twitter)

14.8 Next (topical)

- Writing R system functions
- Statistical functions
- Reading tables with `read.table`

15 Guest talk - Stone Ward - w6s15 (18-Feb)

15.1 Potential questions:

These are my questions informed e.g. by the 2022 data trends report.

1. What do your clients typically expect from you with regard to data science?
2. In the 2022 data trends report, we read that "upskilling [with data literacy skills] becomes a mandate". What is the level of data literacy (with examples) at Stone Ward? Where would you like it to be?
3. How well did your studies prepare you for what you're doing now as a data scientist?
4. What should undergraduates at Lyon know before they decide to embark on a potential career as data scientists or data analysts?
5. How important is machine learning in 2022 - and where is it going?
6. If you compare data science from an industry perspective 5 years ago, now, and 5 years from now - what's different?
7. What should students know before they approach you/Stone Ward for internships? What if they approach Stone Ward for a job?
8. What about a data science minor/major: important? Useful? Relevant?
9. Which projects would you like students to have attempted or completed? Is project experience important at all?
10. Which soft skills are most relevant at Stone Ward?

15.2 Presentation questions:

These are some of my questions after leafing through a pre-view of Matthew's presentation "[Data in Business](#)":

1. Why do clients want analysis? What do they do with the results? (Example)
2. Are clients typically more interested in descriptive (historic), prescriptive (normative) or predictive (future) analyses?
3. How much time do you still spend coding? Reading about R, new packages etc. How important do you think this is?
4. Tidyverse or base R?
5. How important is Excel to your work? How important is it to your clients still? (Compared to R or Python, or platforms like Tableau or Power BI)
6. What's with Plato's cave!?
7. Clients only remember "1-3 numbers" - which numbers are these (example)? How would I know what's important to them?
8. What if I screw up as a data analyst (example)?
9. How did you learn to talk about data and data science?
10. Do clients ever ask you for helicopter presentations like these, or only data analysis presentations (close to the result)?
11. What is a "non-data minded person"? (What are they missing?)
12. Who is on the analytics team?
13. Have you had interns or employees from Lyon College yet?
14. Can you tell us more about the scope of the problem or problems to be tackled in a mini-internship? How much does a student have to know?
15. How large are the data sets that you encounter at clients?³
- 16.

16 Guest talk - Post mortem - w7s16 (21-Feb)

16.1 News

- If you answered TRUE for question 18 on vectorization, contact me and you'll get an extra point for your test. My question was too confusing because the comparison could be seen as vectorization: check with `is.vector("hello")` - scalars and characters are internally represented as vectors, hence the simultaneous application of an operator (`<`) to all its elements could be called vectorization!

16.2 Objectives

[Link to the slides](#)

- [] Summarize talk (small group discussion + presentation)
- [] Identify what is most relevant to you
- [] Critically review claims and recommendations
- [] Apply the presentation to your own learning and career
- [] Learn more about the Google Data Analytics certificate
- [] Understand the problems with Excel
- [] Understand the difference between the "Tidyverse" and Base R
- [] Learn more about ggplot2

16.3 Overview

WHAT	HOW
Discussion in small groups	What did you think?
Overview	Summarize main messages
	Google Data Analytics Certificate
	Excel Errors
	Tidyverse vs. Base R
	The ggplot2 package

16.4 What did you think of the talk?

- Summarize the main messages of the presentation?
- What were your personal takeaways?
- Is there anything you'd like to know from me?
- Do you want to have more presentations like this one?
- Are you interested in the mini-internships at all?

16.5 Next

- Writing functions
- Reading tables
- apply family of functions

17 Emacs recent files, Writing functions 2 - w8s17 (28-Feb)

17.1 News

- Updated schedule: Rcpp, bash, Excel and SQLite
- Interesting [GitHub issues](#) for "forward studying"
- The quizzes are hard(er): I uploaded [PDF versions to GitHub](#)
- Emacs package(s) of the week: recentf (and ace-window)

17.2 Objectives

- [X] Get started with interactive Emacs notebooks
- [X] Save and load user-defined functions
- [X] Practice writing functions (system, stats)
- [] Understand stats functions in R
- [X] Learn a new Emacs package (or two)

17.3 Emacs package of the week: how to display recent files

- Package is actually built in. You call it with `M-x recentf-open-files`, which, on my Windows Emacs, leads to this buffer right now (the buffer below shows the buffer list - `C-x C-b`).
- It is nice that this also works when you kill Emacs by mistake. It is very handy to just be able to continue your work after you've opened and worked in 5-10 files!

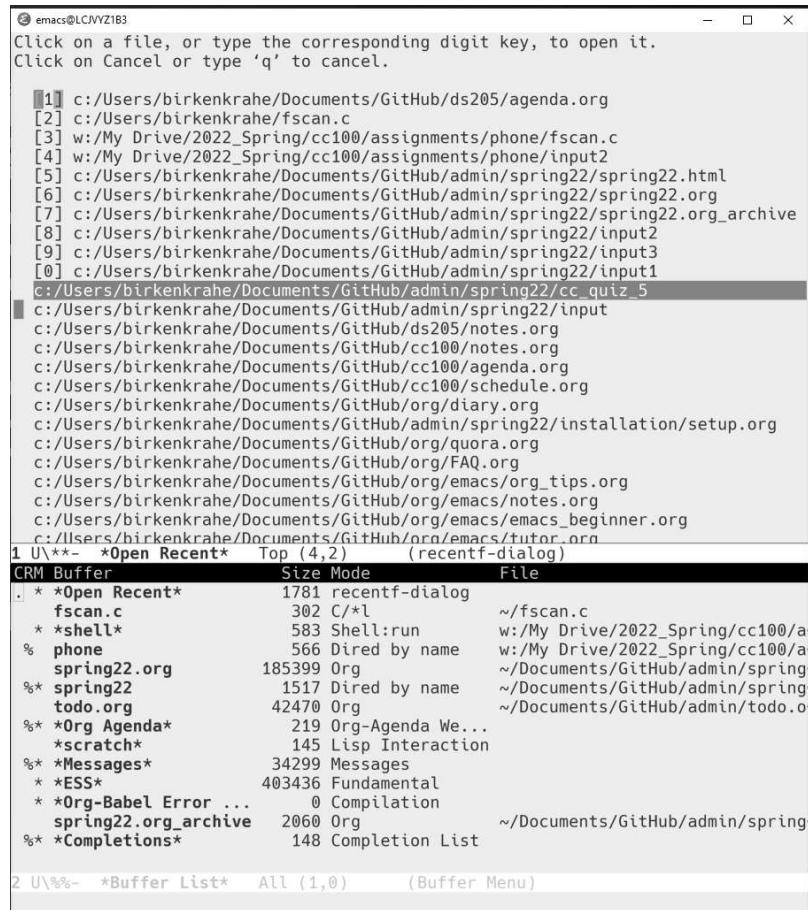


Figure 17: recent files on Emacs

- You can put the following code into your `.emacs` file to bind the command to `C-x r e` (or any other available combination you choose), and to enable it.

```
;; enable recentf mode and bind it to
(recentf-mode 1)
(global-set-key (kbd "C-x rf") 'recentf-open-files)
```

- More information:
 - [Details on Emacs keybindings \(Petersen, 2019\)](#)
 - [Customizing Key Bindings: GNU Emacs manual](#) (also available inside Emacs: `C-h i` opens the Emacs info reader)

17.4 Interactive notebook practice

- Download [save_nb.org from GDrive](#) and work through it in class.

17.5 Next

- Statistical functions (lecture)
- Tour of `apply` (notebook)
- `ggplot2` (lecture + notebook)
- Visualizing COVID-19 (project)

18 Change R download repo - w8s18 + 19 (4-Mar)

18.1 News

- [X] Mid-term grades - Improve your grade with a project ([FAQ](#))
- [X] Waking up in AppData/Roaming? [Change your Emacs HOME now.](#)

18.2 Writing functions: set default R download repo

- Download repos_nb.org from GDrive and get cracking!
- You can find the solution in repos_solution_nb.org in GDrive.

18.3 Next

- Statistical functions
- Tour of apply
- ggplot2 graphics
- Visualizing COVID-19 (DataCapmp project)

19 Graphical devices dev.list, .Library - w9s19 (7-Mar)

19.1 News

- [X] Waking up in AppData/Roaming? [Change your Emacs HOME now.](#)

19.2 Writing functions: set default R download repo

- Download repos_nb.org from GDrive and get cracking!
- You can find the solution in repos_solution_nb.org in GDrive.

19.3 Next

- Tour of apply
- ggplot2 graphics
- Visualizing COVID-19 (DataCapmp project)

20 Gapminder, Pomodoro timer - Tour of ggplot - w9s20 (9-Mar)

20.1 Preparations for test 2 (Mon 14-Mar)

- Test 2 will only cover questions from quiz 4-6 + new questions.
- You can find quiz 4-6 with solutions + feedback as PDF ([in /quiz](#))
- I will create an update of content Org files ([in pdf/](#))

20.2 Emacs package of the week: Pomorodo timer

- Auto-complete org-timer-
- org-set-timer
- org-timer-start
- org-timer-stop

- Lines for .emacs:

```
;; pomodoro timer
(require 'org)
(setq org-clock-sound "c:/Users/birkenkrahe/Documents/Sounds/ding.mp3")
```

20.3 Multiple device control

Device control is quite important in all graphical program. In R, there are multiple functions that achieve this - look up `help(dev.list)` for the documentation, especially `example(dev.list)` (examples at the end of the help file).

Even though the `_` is supposedly Unix-specific, it runs under Windows, too. `x11` is the windows manager. In the example, a series of devices is opened, used, and finally closed again.

```
x11()
plot(1:10)
x11()
plot(rnorm(10))
dev.set(dev.prev())
abline(0, 1) # through the 1:10 points
dev.set(dev.next())
abline(h = 0, col = "gray") # for the residual plot
dev.set(dev.prev())
dev.off(); dev.off() #- close the two X devices
```

```
windows
 3
windows
 4
windows
 3
windows
 4
windows
 2
```

```
windows
 2
windows
 3
windows
 2
windows
 3
null device
 1
```

- Tour of `apply`
- Tour of statistical functions

20.4 ggplot2 and dplyr

- [X] In winter 2020 (DSC101) I gave an introductory lecture to `ggplot2` using [this Google Colaboratory workbook](#). You can copy the workbook to your GDrive and work through it if you like.
- [X] This workbook came from an Org-mode file that I have uploaded in the practice directory in GDrive as `ggplot2_2021.zip` with all plots and images required to view and run it if you so wish. There is also a 24-page PDF copy in the pdf directory.
- [X] In this advanced introductory class, we will instead dive right into an interesting EDA example, the gapminder dataset, to explore the possibilities of `ggplot2` and `dplyr`.
- [X] I say it once at the outset (and hopefully not too many times). I am not a friend of the "Tidyverse" ideology. In practice, I prefer `data.table` over `dplyr`, and base R graphics over `ggplot2`. The reasons are manifold - [see here for details](#).
- [X] See also my notes on the "Tidyverse" and `ggplot2` from the post mortem session of our guest talk by Matthew Stewart.

20.5 Next

- Quiz 4-6 check and test preparation
- Visualizing COVID-19 (project)
- `ggplot2` (lecture + notebook)

21 Test preparation / quiz 4-6 - w9s21 (11-Mar)

- [X] You find the materials of the past month in GitHub ([pdf/](#))
- [X] Review of quiz 4-6
- [X] Review of DataCamp chapters "apply" and "utilities"
 - `lapply`, `sapply`, `vapply` and their uses
 - Useful mathematical functions: `abs`, `sum`, `mean`, `round`
 - Pattern matching: `grep`, `grepl`, `sub`, `gsub`
 - Time and data manipulation and calculation
- [X] Best way to prepare: practice "Intermediate R" in the DataCamp app - my questions are likely to come from that source!

22 Test 2 with ggplot2, org-skeleton - w10s23 (16-Mar)

22.1 News

- Cleaning function with `gsub` in The Economist ([issue](#))
- Emacs snippet of the week: `org-skeleton`
- Test grades plotting with `ggplot2`

22.2 Emacs snippet of the week

- `Org-skeleton` is a built-in Org function but you need to customize it.
- Put the code below into your `.emacs` file, create a new Org file, and enter `M-x org-skeleton` to automatically add the meta data.
- You can also enter `CTRL + SHIFT + F4` to activate the function

```
(define-skeleton org-skeleton
  "Header info for a emacs-org file."
  "Title: "
  "#+TITLE:" str " \n"
  "#+AUTHOR:" str " \n"
```

```

"#+SUBTITLE: " str " \n"
"#+STARTUP:overview hideblocks\n"
"#+OPTIONS: toc:nil num:nil ^:nil\n"
;; "#+email: your-email@server.com\n"
;; "#+INFOJS_OPT: :view:info \n"
;; "#+BABEL: :session *R* :cache yes :results output graphics :exports both :tangle yes
")
(global-set-key [C-S-f4] 'org-skeleton)

```

23 ggplot2 boxplots - w10s24 (18-Mar)

- COVID-10 project deadline moved
- Finish ggplot2 boxplots
- Retrieve updated file from GDrive

24 Analysis COVID-19 DataCamp project - w11s25/26 (28/30-Mar)

24.1 Tentative plan for the rest of the term

- Week 11: Finish gapminder intro to ggplot2 and dplyr packages
- Week 12: Base R bar charts, DataCamp review (categorical data)
 - No class on Friday, April 8
- Week 13: Tableau, DataCamp review (numerical data)
- Week 14: SQLite and R, Excel and R
- Week 15: R and C++, data cleaning on the command line (Wyatt/Ben?)
- Week 16: Webscraping and final project
- 5 quizzes, 1 test, 1 final exam
- 4 DataCamp assignments (EDA using ggplot2 and dplyr)

24.2 Review: COVID-19 project

Download the whole directory ["covid"](#) from GDrive - it includes the data sets in the directory data, and a location for the images in img.

Extract the file anywhere on your computer and go through the Org-mode workbook.

There is a [30-page PDF on GitHub](#) with the whole workbook filled in.

- [] Aspects of the project
- [] Getting the data
- [] Reading the data (readr vs utils)
- [] Plotting with ggplot2
- [] Filtering with dplyr (`filter`)
- [] Grouping with dplyr (`group`)

25 UpdateR, scatterplots - w11s27 (1-Apr)

25.1 Updating R on Windows in 6 steps

```
#+begin_src R :exports both :session :results output
  library(ggplot2)
  library(openintro)
#+end_src

#+RESULTS:
: Loading required package: airports
: Loading required package: cherryblossom
: Loading required package: usdata
: Warning messages:
: 1: package 'openintro' was built under R version 4.1.3
: 2: package 'airports' was built under R version 4.1.3
: 3: package 'cherryblossom' was built under R version 4.1.3
: 4: package 'usdata' was built under R version 4.1.3
```

Figure 18: Screen message (R version update)

1. [] Install and load the installr package
2. [] Run updateR()
3. [] Keep the old version of R (4.1.2 in my case)
4. [] Add new Rterm location to your PATH environment variable
5. [] Set Emacs variable inferior-ess-r-program to Rterm path
6. [] Restart Emacs and M-x R

25.2 Preview: DataCamp course EDA in R

- [X] You find the datasets at the end of the course page
- [X] Read them into data frames using `read_csv` or `read.csv`
- [X] Introduction of the contingency table function `table`
- [X] Uses `factor` and `levels` (both function and attribute)
- [X] Presents different formats of bar charts with `geom_bar`
- [X] Faceting with `facet_wrap`

25.3 Gapminder - scatterplots

- [] Download the gapminder directory from [GDrive](#)
- [] Open gapminder1.org

26 Faceting - w12s28 (4-Apr)

26.1 Homekeeping

- [] Add `(setq-default org-hide-emphasis-markers t)` to `.emacs`, or customize the variable by typing `C-h v org-hide-emphasis-markers` and choose `toggle` and `Apply` and `Save` in the menu⁴.
- [] R 4.1.3 is missing `Internet.dll` - fix⁵

26.2 Review: last time (gapminder 1)

Getting the data

- Knowing alternative paths in R is not a waste
- Installing and loading R packages

- Updating R packages
- Tibble format for data frames
- Fixing an `Internet.dll` issue on Windows

Checking and getting to know the data

- Reviewing structure checking commands
- Changing the display width option
- Printing a data frame as a tibble
- Pipes to pass data to functions
- Pipeline concept

Filtering the data

- Review data findings for: quality, sample, context
- `dplyr` commands resemble SQL: `select`, `filter`, `%in`
- The pipe operator in `dplyr` is `%>%`

26.3 Topics

- [] Assignment: writing `dplyr` function
- [] Topics: facetting with gapminder
- [] Practice workbook [gapminder2.org](#) in [GDrive](#)

27 Time series plots - w12s29 (6-Apr)

27.1 Emacs tip: hide emphatic characters

To **not** see the emphatic characters like `~` or `*` or `/` in the Org file text, run the following code chunk (or put the code in your `/.emacs` file): if successful, you should see "t" in the minibuffer.

```
(setq-default org-hide-emphasis-markers t)
```

Note: If you don't put it in your `/.emacs` file, the command will only work for the current Emacs session, and for all new buffers (not for buffers already opened).

See also: [articles on becoming more productive with Emacs](#).

27.2 Emacs tip: change theme to "Leuven"

- In Emacs, open the theme selector with `M-x custom-themes`
- Select Leuven and click on `Save theme settings`
- Your code blocks are now more clearly visible

27.3 Package tip: overviewR ([issue #43](#))

- [Link to tweet and cheatsheet](#) & [vignette at CRAN](#)
- This could be a cool project - exploring pros and cons

27.4 Topics

- [] Finish facetting: `facet_wrap` in [gapminder2.org](#)
- [] Topic: time series plots with [gapminder](#)
- [] Practice workbook [gapminder3.org](#) in [GDrive](#)
- [] Quiz 8 = DataCamp review + time series + facetting

27.5 Reading and writing assignment for Friday April 8

There will be no class on **Friday, April 8**. Instead, please complete the following assignment:

- [] Read the **introduction** to the essay "[Teaching R in a Kinder, Gentler, More Effective Manner: Use Base-R, Not the Tidyverse](#)"
- [] Pick one of the 10 **case studies** included in the essay and study it so that you can answer the following questions
 1. **Message:** What's the main message of the case study?
 2. **Content:** In your own words, what is the author trying to say?
 3. **Opinion:** What are your personal views (if any) on the case study?
 4. And anything else you wish to share!
- [] Post your thoughts as an [issue to GitHub](#) (with your name and the title of the case study in the title of the issue).
- [] You should not spend less than an hour on this assignment. This may include looking stuff up or trying things out (in R).

Background to the assignment: this essay (from 2019) by R legend Norman Matloff has recently been [greatly revised](#). It is important for any user of R because of the way the 'Tidyverse' is seen (as an alternative R universe). It is important for you as non-teachers, because you're likely to have to teach, or instruct other users of data (who aren't happy enough with Excel). Lastly, it's interesting because it's political, and to see politics, history & business so clearly in computing is rare and will help you grow professionally.

28 Scale and value transformations - w13s30 (11-Apr)

28.1 Housekeeping

- [X] Review: Tidyverse Sceptic essay reading + writing assignment
 - RStudio is not evil, just mistaken
 - Check out [Programming Games with R Shiny](#) (for DataViz? Here?)
- [X] Quiz 8 = DataCamp review + time series + facetting
- [X] Test 3 will be available for 24 hours (we need the time)
- [X] All quizzes are now randomized (don't learn the sequence)x

28.2 Emacs tips - keyboard macros

- Keyboard macros are powerful time savers
- Example: when creating quizzes, I need to surround expressions with the ~ symbol to emphasize code
- Keyboard macro: `C-x ([key sequence] C-x)`
- Repeat with
- E.g. add [] after each bullet point of this list
 - Type `C-x (` to start the macro
 - Enter exact sequence that you want repeated
 - Type `C-x)` to close macro definition
- To run the macro: `C-x e`.
 - Repeat manually with `e`
 - Repeat N times with `C-u N C-x e`
- Bind the macro to a key sequence (for the session):

- C-x C-k b
- Choose 0 to 9
- Save it for later use:
 - Give it a *name* with C-x C-k n
 - Go to /.emacs
 - Save the definition with M-x insert-kbd-macro RET /name/ RET
- In my Emacs file, my macro looks like this:

```
;; surround region from point to EOL with ~
(fset 'tt
  (kmacro-lambda-form [?~ ?\C-e ?~] 0 "%d"))
```

(See [GNU Emacs Manual](#))

28.3 Conceptual [data science] stuff

- This assignment completes our monthly "conceptual" series:
 - 2022 data trends: knowing what the future might hold
 - Stone Ward guest talk: knowing what the industry wants
 - Tidyverse Sceptic: knowing that there are movements
- You do not need to deeply engage with this side of things but you should be aware that there are serious (*ideological*) issues.
- Computer science and data science are also *battlefields* where you may have to pick your sides and your heroes (based on *values*).
- In an *open source* landscape, these battles are more open and visible but of course they do exist in the commercial world, too.
- We met a fourth issue (literate programming) through our Emacs practice - topic of my recent faculty colloquium talk ([link](#)).
- What is programming? It has a coding and a meta part:

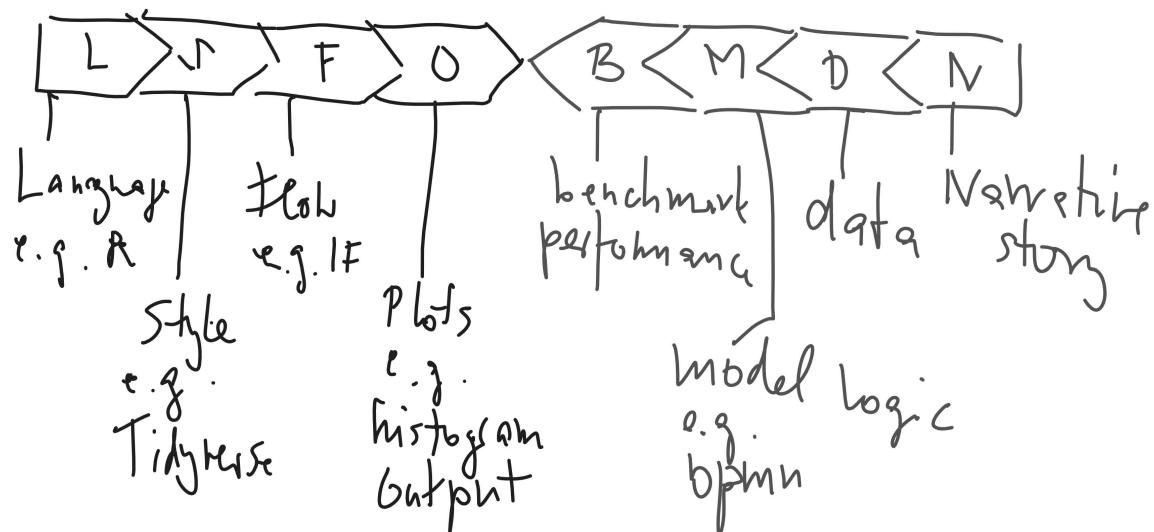


Figure 19: Programming as coding (left) and meta activity (right)

CODING EXAMPLE	META-CODING	EXAMPLE
Language & paradigm	R, C, Java OOP, FP	Story LitProg
Style & Package	"Tidyverse" <code>data.table</code>	Plotting Emacs + Org <code>gapminder</code>
Flow	<code>if, goto</code>	Model logic
Output:	Histogram	Benchmark
		Performance

28.4 Scattermore demo

- This is a cool new package. "Bazillion of points without wait."
- Check out their [GitHub repo](#). Demo code below.

```
library(scattermore)

## create 10 million 2D datapoints
data <- cbind(rnorm(1e7), rnorm(1e7))

## prepare empty plot
par(mar=rep(0,4))

## plot the datapoints and see how long it takes
system.time(plot(
  scattermore(data,
              rgba=c(64,128,192,10),
              xlim=c(-3,3),
              ylim=c(-3,3))))
```

28.5 Topics

- Remember: `ggplot2` is actually not really part of the "Tidyverse"
- [X] Fun package: `scattermore` - "more points = more power"
- [X] Scale and value transformations (`gapminder4.org` in GDrive)
- [] Review: DataCamp - EDA with R - categorical variables
- [] Practice: Base-R vs. "Tidyverse" (`practice/tidyverse`)

29 Boxplots and ridge plots - w13s31 (13-Apr)

29.1 Emacs tip

- With Emacs, you're now part of the free software world
- With great[er] power comes great[er] responsibility
- You need to update your Emacs packages:
 - M-x `package-list-packages` (will try to refresh the list)
 - In `*Packages*`, type U to update
 - If there are packages to update, type x to execute

Package	Version	Status ▾	Archive	Description
tildify	4.6.1	built-in		adding hard space mode for keeping transparent
timeclock	2.6.1	built-in		Transparent Rem
tramp	2.4.5.27.2	built-in		major mode for
vera-mode	2.28	built-in		major mode for
verilog-mode	2019.12.17...	built-in		A full-featured
viper	3.14.1	built-in		Rename files ed
wdired	2.0	built-in		minor mode to v
whitespace	13.2.2	built-in		browse UN*X man
woman	0.551	built-in		Emacs Speaks St
ess	18.10.2	obsolete		Growl Notificat
gntp	0.1	obsolete		The missing has
ht	2.3	obsolete		Make bindings t
hydra	0.15.0	obsolete		provide logging
log4e	0.3.3	obsolete		Other echo area
lv	0.15.0	obsolete		A color theme f
night-owl-theme	0.1.0	obsolete		Support library
pdf-tools	0.91	obsolete		a simple wrappe
pfuture	1.10.2	obsolete		Pop a posframe
posframe	1.1.7	obsolete		The long lost E
s	1.12.0	obsolete		Extended tabula
tablist	1.0	obsolete		View xkcd from
xkcd	1.1	obsolete		Yet another sni
yasnippet	0.14.0	obsolete		Display content
comint-mime	0.1	incompat	gnu	Cycle (rotate)
emacsql-sqlite-...	20220406.1340	incompat	melpa	EmacSQL back-en
flymake-rest	20220409.1233	incompat	melpa	Core features f
fontsloth	20211118.2018	incompat	melpa	Elisp otf/ttf f
helm-gitignore	20170211.8	incompat	melpa	Generate .gitig
magit-commit-mark	20220406.2314	incompat	melpa	Support marking
mode-line-idle	20220406.2322	incompat	melpa	Evaluate mode l
poetry	0.1.0	incompat	melpa-s...	poetry in Emacs
project-tab-groups	20220331.918	incompat	melpa	Support a "one
undo-fu-session	20220412.1212	incompat	melpa	Persistent undo
xref-rst	20220406.2311	incompat	melpa	Lookup reStruct

1 U:%*- *Packages* Bot (8263,0) (Package Menu)

Figure 21: End of the **Packages** buffer

29.2 Topics and resources

- [] Topic: boxplots and ridge plots with [gapminder](#)
- [] [ggplot2: practice workbook gapminder5.org in GDrive](#)
- [] Base-R: [How to make a boxplot in R](#) (Negoita, April 2022)
- []

ggplot2 and dplyr summary

How to make and modify boxplots

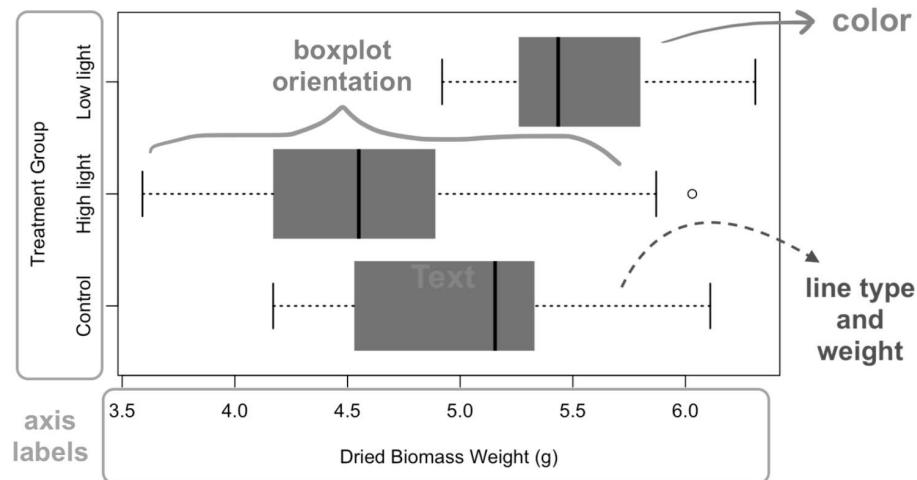


Figure 22: How to make and modify boxplots (Source: R-blogger)

- [] Another resource: '[apply' functions in R](#) (Cheng, 2022)

29.3 ggplot2 and dplyr summary

29.3.1 Concepts

- Aesthetics (ggplot2) = mapping of data on a plot
- Piping (dplyr) = A form of summarizing / passing data on to other functions
- Faceting = displaying plots in different panels (to show plots next to one another)
- Other forms of showing information: stack, ridge, fill, identify
- Object-orientation is part of R (via "S3") - everything is an object - in ggplot2, plots become objects
- Functional programming: analysis is done with functions, which can be nested
- Layout is controlled via geometry and suitable `geom_` functions

29.3.2 Code

30 EDA with categorical data - w14s32 (20-Apr + 22-Apr)

- [X] In software development "[Lightweight is the right weight](#)"
- [X] Followup from the math test lecture yesterday: [igraph](#)
- [X]

Book recommendation for the holidays! "[How to solve it](#)" by George Pólya (1990 - see also [application example](#)).

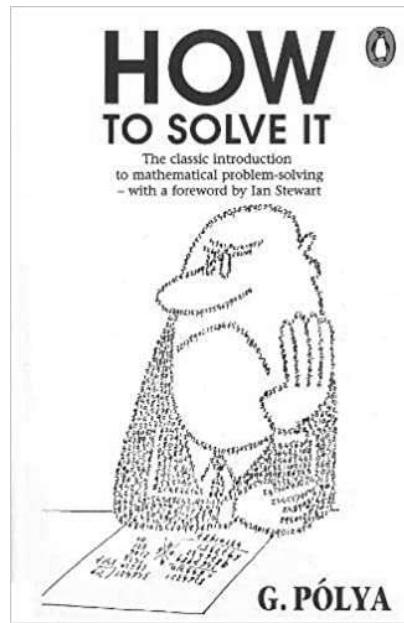


Figure 23: Perennial classic: Pólya's HOW TO SOLVE IT

- [X] Quiz 9 on EDA (last quiz!) online from Fri until Monday
- [X] Exploring Categorical Data workbook ([eda_1.org in GDrive](#)) - solutions in the pdf/ directory in GitHub
- [X] Complete 2nd workbook on categorical variables, [eda_2.org](#).

31 Review quiz 7-9 - w15s34 (25-Apr)

- FREE PIZZA + ORGANIC CHEMISTRY!

Interesting Topics in Organic Chemistry Seminars - On April 27th at noon in Derby 16. Pizza will be provided. The short seminars will be conducted by CHM 220 students and the content is geared to freshmen-sophomore and/or a general audience. (Prof. Narawathne)

- [X] Who is interested? (headcount = 0)
- [X] *How are you?*
- [X] Review quiz 7-9
- [] **Test 3 will run from Thursday 6 pm to Saturday 8 am:**
 - You have 30 minutes for 15 questions
 - **You cannot resume an incomplete submission**
 - Let me know if there are any difficulties during that time!

32 C++ and R (Wyatt) - w15s34 (27-Apr)

32.1 Evaluate

- []

Evaluate this course! Help Lyon decide if I stay or if I go!

Extra credit for evaluating! - Don't forget to click the box.

Enable Send Proof: Allow Students to send proof of survey completion. [?](#)

Participation Incentives: Enable Contest Incentives for this survey.

 Allow Faculty to leave This option allows Students to send an email receipt to Responsible Faculty from their completed survey list. This is for anonymous surveys only.

32.2 R and C++

- []

The amazing Rcpp package - introduced by *Wyatt Frerichs*



Figure 25: Wyatt hails from Oklahoma (Photo: R Sanner on Unsplash)

33 Data science on the command line - w15s35 (2-May)

- [X] What is the command line
- [X] Why data science at the command line
- [X] How to get a command line that works for data science
- [X] Interactive demo: csvkit Python library (YouTube video)
- [] Ben Grafton: Data processing with shell (Wed)
- [] Beyond: csvkit

33.1 Final exam (Schoology update) - Sunday May 8, 3.30pm

COURSE NO	MEETS	TOPIC	TIME	ROOM
DSC 205	MWF 3.00p	Data science	Sun 8-May 3.30p - 5.30p	Lyon 104

The final exam will be **in room 104 of the Lyon building** (our usual classroom) on **Sunday, May 8 from 3.30 pm to 5.30 pm** ([see exam schedule](#)).

The exam will consist of 40 questions drawn at random from the pool of quiz and test questions. Completing the various quizzes and tests until you've reached 100%, and revisiting your past tests should enable you to pass this exam with flying colors!



33.2 csvkit demonstration (video)

Here is the 2nd part of "data science on the command line", the demo of the csvkit library, [as a video on YouTube \(14 min\)](#). The left hand side of the screen shows the Docker container (see FAQ), the right hand side shows the script ([GitHub](#)).

A screenshot of a video player interface. The main content area shows a terminal window with several lines of text. The text includes:

```
state_id country_fips state_name quantity_of_acquisition
state_code total_cost ship_date federal_supply_category federal_supply_category_name
federal_supply_class federal_supply_class_name
```

id	country_fips	state_name	quantity	of	acquisition
1	ADAMS	33,000	1,000	94-585-1273	RIFLE,7.62 MILLIMETER
2	ADAMS	33,000	1,000	97-11	WEAPONS
3	ADAMS	33,000	1,000	through 20 mm	
4	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
5	ADAMS	33,000	1,000	99-585-1273	WEAPONS
6	ADAMS	33,000	1,000	99-585-1273	through 20 mm
7	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
8	ADAMS	33,000	1,000	99-585-1273	WEAPONS
9	ADAMS	33,000	1,000	99-585-1273	through 20 mm
10	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
11	ADAMS	33,000	1,000	99-585-1273	WEAPONS
12	ADAMS	33,000	1,000	99-585-1273	through 20 mm
13	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
14	ADAMS	33,000	1,000	99-585-1273	WEAPONS
15	ADAMS	33,000	1,000	99-585-1273	through 20 mm
16	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
17	ADAMS	33,000	1,000	99-585-1273	WEAPONS
18	ADAMS	33,000	1,000	99-585-1273	through 20 mm
19	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
20	ADAMS	33,000	1,000	99-585-1273	WEAPONS
21	ADAMS	33,000	1,000	99-585-1273	through 20 mm
22	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
23	ADAMS	33,000	1,000	99-585-1273	WEAPONS
24	ADAMS	33,000	1,000	99-585-1273	through 20 mm
25	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
26	ADAMS	33,000	1,000	99-585-1273	WEAPONS
27	ADAMS	33,000	1,000	99-585-1273	through 20 mm
28	ADAMS	33,000	1,000	99-585-1273	RIFLE,7.62 MILLIMETER
29	ADAMS	33,000	1,000	99-585-1273	WEAPONS
30	ADAMS	33,000	1,000	99-585-1273	through 20 mm
31	STATE				
32	COUNTY				
33	FIPS				
34	NAME				
35	STATE_NAME				
36	STATE_CODE				
37	ACQUISITION_DATE				
38	TOTAL_COST				
39	DATA_DATE				

Watch later Share

Examine data... Beyond covid... References...

34 Last Rites / Final exam review - w16s36 (4-May)

34.1 cvskit demonstration (video)

Here is the 2nd part of "data science on the command line", the demo of the csvkit library, [as a video on YouTube \(14 min\)](#). The left hand side of the screen shows the Docker container (see FAQ), the right hand side shows the script ([GitHub](#)).

Data science on the command line: csv...

Watch later Share

CSV file loaded

- should be re-written as CSV file
- can use to get a table, not just raw data
- column provides a tabular view of the data

CSV file loaded
CSV file -> data.csv
CSV file -> data2.csv

• CSV TO CUT OUT columns from a CSV File

- CSV TO is a version of cut for CSV Files
- The --all option shows all columns.
- The --columns option shows specific columns

CSV file loaded
CSV file -> data.csv
CSV file -> data2.csv
CSV file -> data3.csv

• output columns can be called by name, too

CSV file -> data2.csv

- the pipe prints the first 3 rows of the respective columns
- I want to use some of the output later so I put it into a file

CSV file loaded
CSV file -> data2.csv

- ALL of the previous operations can be put together in one pipe.

CSV file -> data2.csv

• Examining data
• Beyond csv...
• References...

Watch on YouTube

34.2 Agenda

- [X] Ben Grafton - Data processing with the shell (offline)
- [X] Your final exam questions
- [X] Going through the exam questions
- [X] Last rites and outlook - exciting data science trends
- [X] **Any outstanding work must be finished by Monday May 9**
- [X] **Please take 10 minutes to complete the course evaluation now**
- [X]

Please be specific: what did you like, what could be changed?



34.3 What I would (probably) change - (video proto script)

34.3.1 What's going to be interesting in the next few years?

In data science, I mean. Some things I became aware of via Twitter:

1. AUTOMATIC THEOREM PROVING (see [Tim Gowers, April 29, 2022](#))

Coincidentally, this was sent to me by Prof Barry Gehm over the w/e:

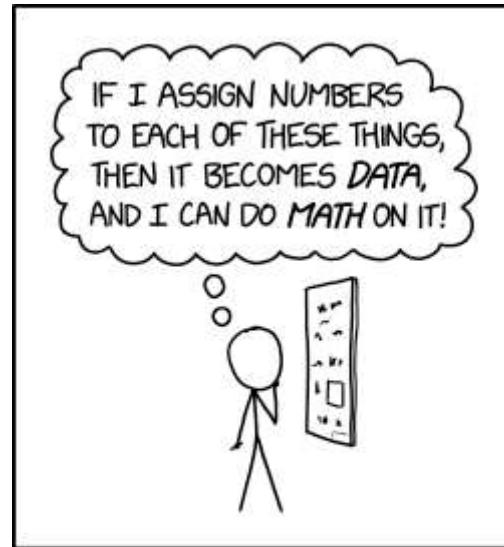


Figure 30: Gödel and data science (Source: xkcd)

Hovertxt: "Gödel should do an article on which branches of math have the lowest average theorem number. (Source: xkcd)

2. FACE RECOGNITION APPLICATIONS OF DEEP LEARNING

Not much new happens technically but there are philosophical and political implications that have powerful economical consequences.

"Tying the technology to accusations of racism has made the technology toxic for large, responsible technology companies, driving them out of the market. IBM has dropped its research entirely. Facebook has eliminated its most prominent use of face recognition. And Microsoft and Amazon have both suspended face recognition sales to law enforcement. These departures have left the market mainly to Chinese and Russian companies. In fact, on a 2019 NIST test for one-to-one searches, Chinese and Russian companies scored higher than any Western competitors, occupying the top six positions. In December 2021, NIST again reported that Russian and Chinese companies dominated its rankings. The top-ranked U.S. company is Clearview AI, whose business practices have been widely sanctioned in Western countries.

Given the network effects in this business, the United States may have permanently ceded the face recognition market to companies it can't really trust. That's a heavy price to pay for indulging journalists and academics eager to prematurely impose a moral framework on a developing technology" [Baker \(2022\)](#).

34.3.2 What I would change about this course

- I wouldn't spend as much time on `ggplot2` and `dplyr`. After teaching it for the second term in a row, I am more, not less convinced that beginner courses should stay as far away from the "Tidyverse" as possible. I

knew this, you see, but after the talk by Matthew Stewart, I decided to respond directly to his request for "More Tidyverse, more ggplot2". I'll be smarter next time.

- I really, really missed projects in this class. I think you'd have been great at completing an analysis by yourself (it worked out very well last term). I was reluctant to try this because so many course participants had not completed DSC 105, and I knew that DataCamp would not be a proper substitute.
- I'd spend more time doing work on the command line. Getting back to the command line interface (**CLI**) after many years was one of my personal re-discoveries in this term - not just in this course. I'd also look more closely at packages like and **Rcpp**. Fortunately, Wyatt covered **Rcpp** and Ben motivated me to at least present the **CLI**. Other packages are just as important, for example **RSQLite** and **data.table**. Spending a week or two on all of these would have been good.
- I would force everyone to complete the (1 hr) Emacs tutorial **in class**. To the very end, I had the impression that 1/2 the class was not really mastering this tool and spent too much time trying to do simple things (like changing buffers, opening, saving and writing files, finding directories, etc.)
- My learning: I need to spend more effort trying to get you to take your **tools** seriously. Without mastery of your tools you are nothing in any discipline but this is especially true for programming, and even more for statistical programming and data science. Many tools (like Emacs, liked editors, compilers, IDEs) don't look like tools because clever software engineers have created an environment for you where the food falls from the tree straight into your mouth most of the time. But that's not the real world! In the real world you have to **learn** to use a hammer to drive a nail into wood without hurting yourself!
- To improve the test results: more, not less drills. More assignments not from DataCamp.
- I cannot decide if I should rather have assigned book chapters or sections of one or more books to you to read from week to week. What do you think? I don't work with one, but with many books all the time, and I really like learning this way.
- Having said all that, I really, really enjoyed the course - if only because it was a constant challenge for me to challenge myself in many different ways: to structure the sessions, to come up with interesting notebooks, to work out quizzes and drills, to complement the DataCamp assignments, etc.
- At the end, I am still inspired by your presence and your willingness to learn, and I am dying to get to teach this 2nd introductory course again so that I can make it better!
- And for the seniors: I'd really like to hear from you!

35 References

- Birkenkrahe (Jan 11, 2022). Interactive shell vs. interactive notebook (literate programming demo). [URL: youtu.be/8HJGz3IYoHI](https://youtu.be/8HJGz3IYoHI).
- Cheng (Mar 8, 2022). Complete tutorial on using 'apply' functions in R - How to use apply(), lapply(), sapply(), and tapply(), and how they compare to using dplyr [blog]. [URL: rforecology.com](https://rforecology.com).
- Cotton (Oct 25, 2018). How DataCamp Handles Course Quality [blog]. [URL: www.datacamp.com](https://www.datacamp.com).
- DataCamp (2022). 2022 Data trends and predictions. [URL: datacamp.com](https://datacamp.com).
- ESS (n.d.). Emacs Speaks Statistics. [URL: ess.r-project.org](https://ess.r-project.org)
- Emacs Speaks Statistics (Mar 19, 2021). First Steps With Emacs [video]. [URL: youtu.be/1YOrd7NCGkg](https://youtu.be/1YOrd7NCGkg).
- GNU Emacs (n.d.). GNU Editor. [URL: gnu.org/software/emacs/](https://gnu.org/software/emacs/)
- Negoita (Apr 6, 2022). How to make a boxplot in R - The complete beginner's tutorial on boxplots in R [blog].
- Petersen (2019). Mastering Key Bindings in Emacs [website]. [URL: masteringemacs.org](https://masteringemacs.org).
- R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.r-project.org/>.
- System Crafters (Aug 1, 2021). Emacs Has a Built-in Pomodoro Timer?? [video]. [URL: youtu.be/JbHE819kVGQ](https://youtu.be/JbHE819kVGQ).

35.1 Image references

- Oklahoma [Photo by Raychel Sanner on Unsplash](#)

Footnotes:

¹ Submission of the assignment by Monday 24 January 3pm gives 10 extra credit points.

² Somehow Schoology counts 14, not 13 participants.

³ This was answered in the talk later. MS said that he had been asked to analyze time series data sets containing no more than 3 months of data. Depending on the number of observations, this could mean that the data set consists of 90 lines only, which is very small indeed.

⁴ And if you wish for the emphasis markers to appear when you hover over an expression with the mouse, install the `org-appear` package and add this code to your `.emacs` file, too:

```
;; Show hidden emphasis markers
(use-package org-appear
  :hook (org-mode . org-appear-mode))
```

⁵ On my machine, R version 4.1.3 cannot open documentation in a browser. This seems to be a Windows 10 issue (`internet.dll` is missing). I simply copied the file from `R-4.1.2/modules/x64` to `R-4.1.3/modules/x64` and that fixed it.

Author: Marcus Birkenkrahe

Created: 2022-05-04 Wed 21:30