

Applied Data Analysis (CS401)



Lecture 1 Intro to ADA 10 Sep 2025

EPFL

Maria Brbić

1

Important websites



<http://ada.epfl.ch>

Your main entry point. All materials linked from there.



<https://edstem.org/eu/courses/2502/discussion>



Main communication channel. Sign in with your EPFL email address (or simply access via Moodle).

<https://github.com/epfl-ada/2025>

Used for exercises, homework, project, and final exam.

2

Credits

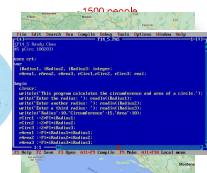
- Robert West
- Now your ML instructor



3

About your instructor

- Born in Tučepi, Croatia



- Education:
University of Zagreb, Croatia
Stanford University, USA

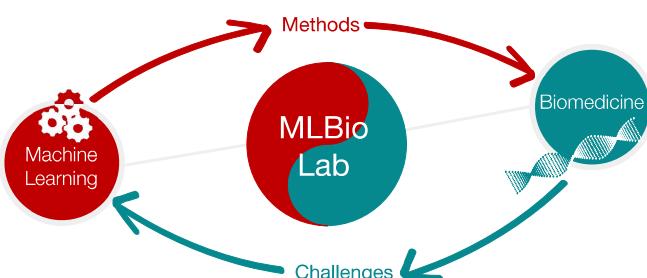


- Assistant Professor at EPFL since Sep '22
Machine Learning for Biomedicine (MLBio) lab

EPFL

4

Our research @ MLBio



5

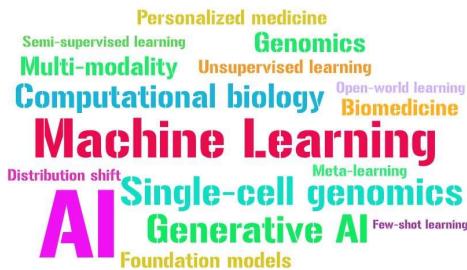
Our research @ MLBio

- Develop new AI methods
 - Generative AI, foundation models, multi-modal models
- Collaborate with biologists and medical researchers and have new “untouched datasets” or collect datasets
- Gain new insights from these datasets ◻ what are interesting questions that need new AI algorithms to be answered?
- Apply AI algorithms we develop to advance biomedical research and drive new discoveries in biology and medicine



6

Our research



7

Data analysis

“... the process of **inspecting, cleaning, transforming, and modeling data** with the goal of **discovering useful information**, suggesting conclusions, and supporting decision-making.”



“Data analysis has multiple facets and approaches, encompassing **diverse techniques** under a variety of names, **in different business, science, and social science domains.**”

8

Applied data analysis

- This course is about **breadth**, not depth
- “**What methods, principles, and tools are out there?**”, rather than “**How can I become an expert in deep learning for computer vision applied to images of cats?**”
- Data science is a fast-paced, shifting field
- Obsessing on one tool or technique won’t pay off in a few years
- Be ready to explore and keep learning on your own

Goal of this class: Enable you to conduct a full-fledged data science project from start to finish

That being said, depth matters too...

Complementary courses:
[Machine learning](#)
[NLP](#)
[DIS](#)
[Data viz](#)

9

Let's abbreviate this course as **Ada**, not A-D-A, in honor of **Ada Lovelace**, “the world’s first computer programmer.”

https://en.wikipedia.org/wiki/Ada_Lovelace



10



Syllabus

- Handling data
 - “Slicing and dicing”: obtaining, preparing, juggling data
- Visualizing data
 - Exploration of data, communication of results
- Describing data
 - How to support (and be suspicious of) claims about data
- Regression analysis for disentangling data
 - How to disentangle datasets with correlated variables
- Causal analysis of observational data
 - How to deal with “found data”
 - Correlation != causation

11

Syllabus (cont'd)

- Learning from data
 - Supervised learning
 - Unsupervised learning
 - Applied aspects of machine learning
- Handling specific types of data
 - Handling text data
 - Handling network data
- Scaling to massive data

Grading

- **15% Homework assignment**
 - Involving skills required from data scientists
 - Groups of 5 students
 - Homework of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)
- **50% Final exam**
 - First part: Multiple-choice questions
 - Second part: Mini data analysis project
 - Done on laptop, individually, on campus
 - Final exams of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)
- **35% Project (more details soon)**
 - Your own freestyle data analysis
 - Done in groups of 5 students (same as for homework)
 - Milestones spread throughout the semester
 - Projects of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)

13

14



Grading

- **15% Homework assignment**
 - Involving skills required from data scientists
 - Groups of 5 students
 - Homework of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)
- **50% Final exam**
 - First part: Multiple-choice questions
 - Second part: Mini data analysis project
 - Done on laptop, individually, on campus
 - Final exams of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)
- **35% Project (more details soon)**
 - Your own freestyle data analysis
 - Done in groups of 5 students (same as for homework)
 - Milestones spread throughout the semester
 - Projects of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)

This class will be hard work,
but it will get you a job.



15

16

Grading (cont'd)

- To obtain a meaningful grade distribution, **scaling/shifting** will be applied to each of {homework, project, exam} before taking weighted average (standard practice at EPFL)
- While intermediate grades are a good indication of where you stand, remember there might be some wiggle
 - **Don't rely on intermediate grades**

Deadlines

- **Homework**
 - **Homework**
 - Release Nov 5th 2025
 - Due Nov 26th 2025
- **Final exam**
 - Date TBD
- **Project deliverables**
 - **Project milestone P1**
 - Due Oct 1st 2025
 - **Project milestone P2**
 - Due Nov 5th 2025
 - **Project milestone P3**
 - Due Dec 17th 2025



All deadlines are 23:59 CET

17

18

Meeting logistics: Lectures

- **Wednesdays 8:15–10:00**
- If you want to see it live, come to class! (No live streaming)
- Lectures are also recorded and made available after class

Meeting logistics: Lab sessions

- Fridays 3:15–4:45
- In person only:
 - [GCC 330](#)
 - [CE 14](#)
 - [BCH 2201](#)
- Labs are complementary to lectures, not simply more detail on same
- You solve exercises that we make available the day before, can ask questions and get help from assistants
- In certain weeks: homework/project office hours (probably on Zoom, in parallel to exercises)

19

Weekly quizzes

- Available online on Moodle after every lecture
- 5 questions, to be answered within 10 minutes of starting
- Quiz 2: the first quiz with lecture material questions
- Quiz i is about lecture material of week i
- Goal:
 - Engage continuously with course material
 - Think (not just find right slide)
- Not graded, for you to recap lecture materials

20

Project

- We'll provide a number of datasets
- You need to form and pitch a crisp project idea
- Free to combine with other datasets (at your own risk)
- Goal: not a loose collection of results – tell a story with the data!
 - Data stories of [2017](#), [2018](#), [2019](#), [2020](#), [2021](#), [2022](#), [2023](#), [2024](#)
 - Nice [example](#) data story

21

Homework and projects: GitHub



- De-facto standard for managing and sharing code
- All students in this class need a GitHub account
- Homework and project submissions done via GitHub
- ADA GitHub repository:

<https://github.com/epfl-ada/2025>



22

Main communication channel: ed

- Class forum, available via Moodle
- Also accessible directly, outside of Moodle:
<https://edstem.org/eu/courses/2502/discussion>
(sign in using the same email address as for Moodle)
- Central place to ask all class-related questions
- Don't send us emails
- Mandatory! We'll send important announcements on Ed only
- Help each other (without cheating, of course)

23

Watch-at-home videos

- Throughout the semester, we'll release videos with supplemental information; e.g.,
 - Intro to lab sessions ([already available!](#))
 - Project instructions
 - Homework postmortem

24

General note on communication

- Multiple platforms used in ADA for various tasks (as in real life): Ed, GitHub, Google docs, ADA website
- To avoid confusion,
 - familiarize yourself with [communication guidelines](#)
 - all materials will be linked from the website as a central point of entry: <https://ada.epfl.ch>
 - all discussions will take place on Ed

25

Commercial break



26

Group registration

- Must form teams within 2 weeks, starting now
- Get started immediately to find 4 other teammates
- By **Fri Sep 26th 23:59**, complete the registration form (to be done by each team member individually):
<https://go.epfl.ch/ada2025-team-registration>

27

Prerequisites

Basics of

- **probabilities and stats**
- **databases**
- **programming**



- You won't survive if you can't program
- Homework, exam: Python required
- Project: up to you, but we support only Python
- Brush up your Python skills (many great online courses out there)

28

Python environments

- Homeworks and exams to be done as [Jupyter Notebooks](#)
- You will submit a pre-executed .ipynb file
 - We don't care how you produce it
 - Option 1: local Python installation (e.g., [Anaconda](#) + [JupyterLab](#))
 - Option 2: [Google Colab](#) = notebook hosted by Google
 - Option 3: [note](#) = notebook hosted by EPFL
- To get started: come to Friday's lab session ("[Exercise 0](#)")
- "[Homework 0](#)": do it yourself at home after lab session (optional, not graded)
- Doing Homework 0 is the best way of making sure you're set up correctly for later homework, project, exam

29

Python++



30



POLLING TIME

- “What is your prior experience with Python?”
- Scan QR code or go to <https://app.sli.do/event/f5usXPBvsT6GLWi5vLuAd6>



31

Instructor



Maria Brbić

Head TAs



Shuo Wen Siba Panigrahi

TAs: Teaching assistants



Vinko Sabolcic Mete Ismayilzada Tim Davidson Aoxiang Fan Shiqi Wang Sevda O gut Yulun Jiang Savaryaj Deshmukh Alba Carballo Castro



Artyom Gadetsky Yist Yu Schuyler James Stoller Pierre Beck Paula Sanchez Lopez

Vinko Sabolcic Mete Ismayilzada Tim Davidson Aoxiang Fan Shiqi Wang Sevda O gut Yulun Jiang Savaryaj Deshmukh Alba Carballo Castro

SAs: Student assistants (Master students)



Abdu Karim Mouakeh Zahra Taghizadeh Lysandre Costes Marija Zelic Alexander Proclewski Jean Silfert Alessandro Di Maria
Nastaran Hashemisanjani Sara Zatezalo Kyuhee Kim Amene Gafsi William Jallot Yassine Mustapha Wahidy

- Help each other on Ed
- Participate actively in classes and labs
- Give us **feedback**



34

Feedback

Give us feedback on this lecture here:

<https://go.epfl.ch/ada2025-lec1-feedback>

Feedback form available for each lecture and lab session

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- ...

35

Questions? 🤔

36

What is data science?

HARVARD
BUSINESS
REVIEW



Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

37

Why now?



38

“Data science”

- Most science is (or should be) based on data, per definitionem
- So how is “data science” different from plain old “science”?

Data volume explodes

“Between the dawn of civilization and 2003, we only created **five exabytes** of information; now [in 2010] we’re creating that amount **every two days.**”

Eric Schmidt, Google (2010)

39

40

Data variety explodes

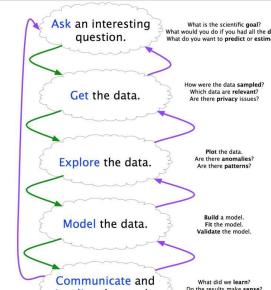


Text (indexed Web pages, email),
networks (Web graph, knowledge graph), **images, maps, logs** (search logs, server logs, GPS logs),
speech, ...



41

Needed: A method to the madness

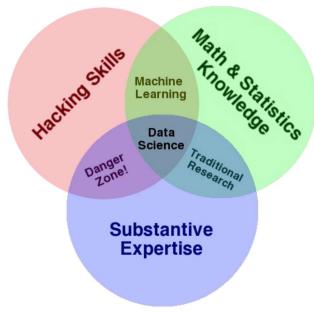


- **Scientific method 1.0:**
 - Focused on “Model the data”
 - Scientist has hypothesis prior to analyzing the data
- **Scientific method 2.0:**
 - Data-driven science
 - Systematic cycle (see diagram)
 - “Explore the data” becomes increasingly important

Data as a first-class citizen

42

Scientist 2.0



“A data scientist is someone who can obtain, scrub, explore, model, and interpret data, blending hacking, statistics, and machine learning. Data scientists not only are adept at working with data, but appreciate data itself as a first-class product.”

Hilary Mason, chief scientist at bit.ly

43

44



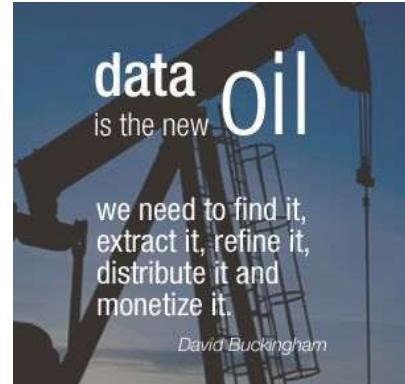
Following

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

Josh Wills, Data Scientist at Slack

45

46



Programming Life like a software: How Digital Biology Will Disrupt Everything?

X (formerly Twitter)
Seth Bannon (@sethbannon) on X
"Where do I think the next amazing revolution is going to come? And this is going to be flat out one of the biggest ones ever."

There's no question that digital biology is going to be it."

Jensen Huang, founder & CEO of NVIDIA. (72 kB) ▾



47

48

More data often beats better algorithms



Alan Halevy, Peter Norvig, and Fernando Pereira, Google

The Bitter Lesson

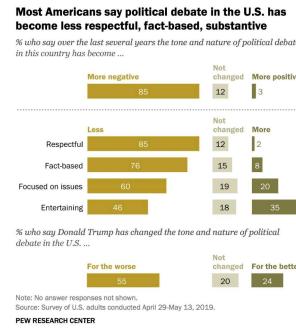
Rich Sutton

March 13, 2019

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are already the most effective, and by a large margin. The ultimate source of power in AI has been computation, not knowledge or performance. As computation has become exponentially more powerful, its cost has been exponentially falling cost per unit of computation. Most AI research has been conducted as if the computation available to the agent were constant (in which case leverage computation would be irrelevant). In fact, computation has been increasing exponentially over a slightly longer time than a typical research project, massively more computation has been available to agents than they could have used. This makes a difference in the shorter term, researchers seek to leverage their human knowledge of the domain, but the only thing that matters in the long run is the leveraged cost of computation. In the short term, researchers can ignore computation because they tend to. Time spent on one is time not spent on the other. There are psychological constraints on how much computation researchers can ignore, but these knowledge approach tends to complicate methods in ways that make them less suited to taking advantage of general methods leveraging computation. There were many examples of AI researches belated learning of this bitter lesson, and it is instructive to review some of the most prominent.

<http://www.incompleteclass.net/Inclass/BitterLesson.html>

21st-century politics



We ask: Do these subjective impressions reflect the true state of US political discourse?

ADA will teach you the tools to answer such questions using data (see next slides)

https://www.pewresearch.org/politics/wp-content/uploads/sites/4/2019/06/PP_2019.06.19_Political-Discourse_FINAL.pdf

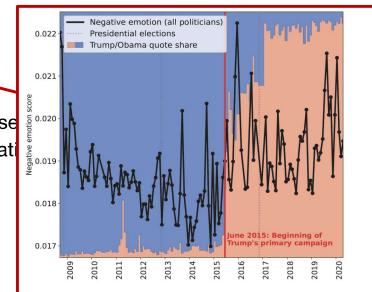
Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- Regression analysis for disentangling data
- Causal analysis of observational data
- Learning from data
- Handling text data
- Handling network data
- Scaling to massive data



Syllabus, revisited

- Handling data
- **Visualizing data**
- Describing data
- Regression analysis for disease
- Causal analysis of observations
- Learning from data
- Handling text data
- Handling network data
- Scaling to massive data



51

52

Syllabus, revisited

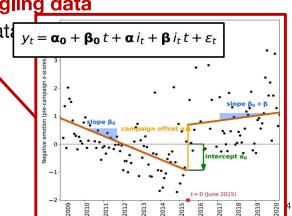
- Handling data
- Visualizing data
- **Describing data**
- Regression analysis for disentangling data
- Causal analysis of observational data
- Learning from data
- Handling text data
- Handling network data
- Scaling to massive data

"Is the effect real, or could it have been produced by chance?"

53

Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- **Regression analysis for disentangling data**
- Causal analysis of observational data
- Learning from data
- Handling text data
- Handling network data
- Scaling to massive data



Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- Regression analysis for disentangling data
- **Causal analysis of observational data**
- Learning from data
- Handling text data
- Handling network data
- Scaling to massive data

"What caused the observed increase in negativity?"

55

Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- Regression analysis for disentangling data
- Causal analysis of observational data
- **Learning from data**
- Handling text data
- Handling network data
- Scaling to massive data



56

Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- Regression analysis for disentangling data
- Causal analysis of observational data
- Learning from data
- **Handling text data**
- Handling network data
- Scaling to massive data

Research question ("Did political discourse become more negative?") is a question about language == text

57

Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- Regression analysis for disentangling data
- Causal analysis of observational data
- Learning from data
- Handling text data
- **Handling network data**
- Scaling to massive data

"Who speaks about whom in what way?" → Construct "who-mentions-whom" network

58

Syllabus, revisited

- Handling data
- Visualizing data
- Describing data
- Regression analysis for disentangling data
- Causal analysis of observational data
- Learning from data
- Handling text data
- Handling network data
- **Scaling to massive data**



59

Curious to learn more?

Full paper available at <https://www.nature.com/articles/s41598-023-36839-1>

The article discusses the shift in political discourse tone in the US during the 2016 presidential campaign. It uses a large dataset of 235 million quotes from 127 million news articles to analyze the change in language used by political candidates. The study found a significant increase in negative language, particularly in Donald Trump's speeches. The authors apply psycholinguistic tools to analyze the tone of political discourse in online media. The results show a clear shift towards more negative language, particularly in Donald Trump's speeches compared to Hillary Clinton's. The study also found that negative language increased during the 2016 primaries and continued through the general election. The authors conclude that the shift in language reflects a broader societal trend of increasing negativity in politics, with people using more negative language in their political discourse.

60

TODO before Friday's lab session

- Sign up for Ed [here](#) and familiarize yourself with it
- If you're not on GitHub yet, sign up for GitHub
- Start looking for 4 teammates
 - You may use “Group formation” category on Ed
- Check out [Google Colab](#) and [noto](#) (to see if you want to use either of them)
- Check out Exercise 0 [here](#) (in prep for Fri lab session)

61

Any feedback? -- Let us know!

Give us feedback on this lecture here:
<https://go.epfl.ch/ada2025-lec1-feedback>

Feedback form available for each lecture and lab session

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more details?
- What would you like the instructor to wear next time?
- ...

62