

Mechatronic Engineering program
Python for machine learning
and data science:
1: Introduction and data management

Ziemowit Dworakowski
AGH University of Krakow

1

Before we begin...

Lecture presentations (and all other course materials) will be available on my webpage at least 2 days before given lecture or laboratory:

<http://galaxy.agh.edu.pl/~zdw/students.html>

I recommend making notes only for stuff that is NOT on the slides. During lectures just focus on understanding relations and thought process. Knowledge will be easier to memorize after you see *the big picture*

This is the second time that this course is offered, some of it is still an experiment. Your feedback is important – if you feel something is not going as it should, notify me as soon as possible, don't wait till the end of semester.

2

What this course is all about?

We will be learning Python (a bit)

We will be learning Machine Learning (more than before)

We will be learning actual teamwork

We will be learning to think like data scientists

3

What should you know already?

„AI stuff” from *Signal Processing and Identification (...)* courses:

- Basics of optimization (gradient and 1+1 methods)
- Basics of classification (algorithms that we've used, basic vocabulary and rules)



↓

If you attended the polish version of the course before or just want a quick reminder, download my lecture slides on optimisation and classification, read through them and you should be good to go:

http://galaxy.agh.edu.pl/~zdzw/Materials/SP/2022_PSIWSUM_Notes_EN_01_Optimization.pdf
http://galaxy.agh.edu.pl/~zdzw/Materials/SP/2022_PSIWSUM_Notes_EN_03_Classification.pdf

4

Course team

 <p>Dr hab. inż. Ziemowit Dworakowski,</p> <p>Course coordinator, Lectures and projects</p> <p>D3, TBD zdzw@agh.edu.pl Office hours: Thursday, 8:00 – 10:00</p>	 <p>mgr inż. Adam Machynia</p> <p>Laboratories and projects</p> <p>D3, TBD machynia@agh.edu.pl Office hours: Thursday, 10:00 – 11:00</p>
---	--

Please notify respective teacher via email at least a day before you plan to use office hours.

5

What about it being English-only course?

- Lectures will be supported by notes
- Most of the tasks will be performed in teams, make sure that in each team we have at least one person with decent English.
- All of the written work (reports) should be grammatically correct, without typos, etc. But you have **ChatGPT** and **Grammarly** so that's not an issue
- Your English skills will not affect your grades in tests – you just need to be communicative, that is enough.

6

Syllabus for the course is available here:
<https://syllabusy.agh.edu.pl/pl/1/2/19/1/4/5/55>

7

So what exactly are contents of this course?

We have:

- 5 lectures
- 5 project classes
- 10 laboratories

L1: Introduction & data management
L2: Basics of scientific method and data interpretation
L3: How to do regression in practical cases?
L4: Classification revisited – from a learner perspective
L5: Challenges and problems in data science, data leakage, unbiased data interpretation, etc.

8

So what exactly are contents of this course?

We have:

- 5 lectures
- 5 project classes
- 10 laboratories

P1: Teamwork organization, agile methodology
P2: Waterfall methodology, project planning and scheduling
P3: Business presentation, selling your ideas
P4: Problem session – how to deal with issues in team
P3: Final presentations and assessment, feedback session

9

So what exactly are contents of this course?

We have:

5 lectures

5 project classes

8 laboratories

Lecture tests

L1: Python & Colab, loading and visualizing data

L2*: Drawing hypotheses and understanding datasets

L3: Feature selection

L4* & L5: Classification and regression

L6* & L7: Assessment and improvement of decision systems

L8: Conclusions and work evaluation

10

OK, apart from the tests, how will we get a grade?

- 1) Passing all laboratories is required. Some will be graded, some (project-related) are just for a pass
- 2) All the tests are required (with a positive grade)
- 3) At least 4 out of 5 project classes are required, including the final one
- 4) You will prepare project report and do a project presentation/defense

Final grade = (2* Project grade + Average test grade + lab grade) / 4

11

How will projects be graded?

- Each project will include research of one particular dataset. The actual task will be formulated mid-semester based on your conclusions from initial research – and will be different for different teams working on the same dataset.

- Several of the laboratories will be devoted strictly to project-related tasks. During these laboratories you will be working on required project parts.

- Final project grade will be affected by:

- Presentation quality (team)
- Report quality, both merit and form (team)
- Project defense (individual)

12

Laboratories are again using „reversed classroom” approach...

So again we have 3.0, 4.0 and 5.0 tasks, you are supposed to work at home before classes, you (hopefully) don't prepare any reports. Recommended scope includes doing at least tasks „for 3.0” on your own.

„Project laboratories” will not be graded (they contribute to the final project report). They will, however, contain „tasks for 3.0” – which are required for project positive grade and „tasks for 5.0” which allow you to gain additional points for higher project grade. You are free to pick and choose which of these you find interesting and worth pursuing.

13

What if something goes wrong?

If you don't pass a laboratory (either due to lack of preparation or due to absence) – you prepare a standard report with one additional task selected by the LA. You may do so twice in total (after that passing conditions will be set individually).

You can attempt any test up to three times.

Note, that right to 3 „attempts” extends also to project reports.

Note that my webpage contains a document on report writing (including laboratory and project reports)

14

What should you prepare in advance?

Project classes:

You should probably sign in to **Asana**
(We will use it to manage project work)

Laboratories:

On my webpage there is a **matlab-python handout** prepared by us to facilitate your initial steps in Python. Download it (possibly: print it as well?) and have it ready before the first laboratory.

15

What is python?

- High-level programming language
- Easy for beginners
- Significant indentation
- Easily readable
- Functional
- Dynamically typed
- Object-oriented
- Modular
- Free and open-source
- Garbage collected
- Go-to option for machine learning

16

Where will we code?

Standard (PC) approach:

- 1) Install Python
- 2) Configure environment
- 3) Install a preferred IDE
- 4) Add libraries as needed
- 5) Solve conflicts as needed

Colab (online) approach:

- 1) Have google account
- 2) Log in to colab, write your codes there
- 3) ... (that's it!)

Feel free to configure your own PC environment and use it, you will probably learn much more this way. However, you are on your own in case of conflicts, problems or stability.

17

How to share your codes in team?

a.k.a: **Why will encourage you to use GIT?**

- Version control
- Standard tool if you want to actually code in the future
- Possibility of easy cooperation
- Possibility of simultaneous off-line work

On my webpage there is a supplementary handout for colab-GIT configuration

18

How will we approach python?

Matlab is similar to Python

You all know Matlab



We'll learn via differences between Python and Matlab

You will be expected to work the differences on your own, using handout provided by us:



<http://galaxy.agh.edu.pl/~zdw/students.html>

19

Do we have an elephant in the living room?

a.k.a. what about ChatGPT and other LLMs?

What is mandatory:

- Acknowledgement of form of usage if it is used (e.g. „report was rewritten for clarity using ChatGPT“)
- Source files (Full prompts) available on-demand

What you can do:

- Use it to correct codes (bug fixes)
- Use it to prepare fragments of reports from prompts or drafts (note, you should store these prompts or drafts!)
- Use ChatGPT to translate reports from polish to english

What you should do:

- Use ChatGPT to correct your texts (maintain consistent and easy-to-read style, correct typos, grammatical errors, etc.)

What might get you into trouble:

- Use ChatGPT to prepare complete codes
- Use ChatGPT to gain new knowledge (merit questions regarding e.g. machine learning)
- Use ChatGPT to do interpretation tasks for you

20

Data management

21

Lets look at a typical dataset:

22

Data types

Numerical

Floating point

Temperature
34.31
24.51
23.23
31.19
33.10
32.54

Integer

Price
1400 \$
1200 \$
700 \$
800 \$
1100 \$
1000 \$

Categorical

Binary (Dichotomus)

Employment
Yes (1)
No (0)
No (0)
Yes (1)
No (0)
No (0)

Polytomous

Education level
Bachelor
Bachelor
Master
Doctorate
Bachelor
Master

23

Data types

There are two ways of approaching categorical variables:
 - Use them to divide and analyze dataset

24

Data types

There are two ways of approaching categorical variables:

- Use them to divide and analyze dataset
- Map them into numerical variables

a) Use **dummy variables**
 b) Use integers

Bachelor	Master	Doctorate	Education
1	0	0	1
1	0	0	1
0	1	0	2
0	0	1	3
1	0	0	1
0	1	0	2

Only if a variable has „natural progression“!

We can do this, because „Master“ is between „Bachelor“ and „Doctorate“ in education level

Education level
Bachelor
Bachelor
Master
Doctorate
Bachelor
Master

Categorical

Polytomous

25

Data types

There are two ways of approaching categorical variables:

- Use them to divide and analyze dataset
- Map them into numerical variables

a) Use **dummy variables**
 b) Use integers

Food	Clothing	Leisure	Industry
1	0	0	1
1	0	0	1
0	1	0	2
0	0	1	3
1	0	0	1
0	1	0	2

Only if a variable has „natural progression“!

We can't do this, because „Clothing“ is not higher in any sense than „Food“ and lower than „Leisure“

Industry
Food
Food
Clothing
Leisure
Food
Clothing

Categorical

Polytomous

26

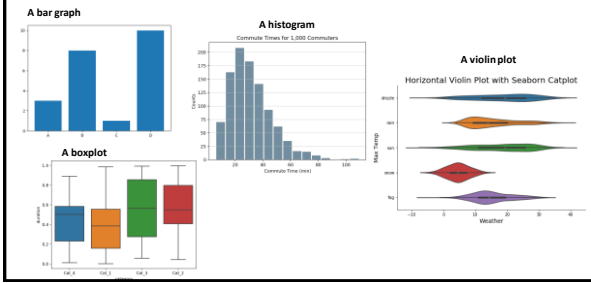
How to visualize data?

A time series plot

A scatterplot

27

How to visualize data?



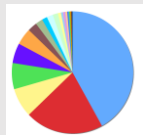
28

How **not** to visualize data?

29

Data visualization errors


1. Using pie charts




- You look at volumes but you should look at angles
- Very small values tend to be perceived similarly
- It is hard to compare categories

30


Percentage of pie charts that do not show the data clearly




Percentage of charts in this slide that are not necessary



Percentage of pie charts that could be replaced by bar charts



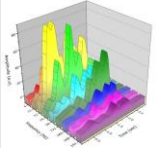


Percentage of 3D and exploded pie charts that are even worse than standard pie charts



31

Data visualization errors




1. Using pie charts
2. Using fancy 3D visualizations



32

Data visualization errors

1. Using pie charts
2. Using fancy 3D visualizations
3. Using unclear colormaps



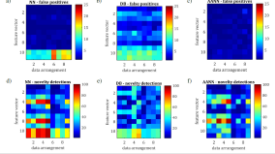
↑ ? ↑ ? ↑ ?

If you can name more than 3 colors of a colormap, it is probably not good. If you can see any sharp change in colors, it is bad as well.

33

Data visualization errors


1. Using pie charts
2. Using fancy 3D visualizations
3. Using unclear colormaps
4. **Using too complex visualizations**



34

Data visualization errors

1. Using pie charts
2. Using fancy 3D visualizations
3. Using unclear colormaps
4. Using too complex visualizations
5. **Overcomplication of simple relations**



35

Data visualization errors

„If you torture the data long enough, it will tell you anything“

John W. Tukey

36

Things to remember:

1. What are the features of Python language?
2. What types of data are there? How can we approach categorical data?
3. Enumerate and explain different visualization techniques (axes descriptions, usage)
4. What are the most common data visualization errors?
