# FEDS - Modelling and Analysing of Forestal and Environmental Data

An introduction to data science

## How to learn?

Do not expect the teachers to teach you.

They will present some information to you, but it is entirely 100% up to you to either make the most of it, or waste your time here, and go home and get a normal dumb job.

# How I will present you the informations!

- I will jump between websites, slides, coding ... will be intense
- I offer you to use some of your time and do challenges
- I get you to know stuff you have maybe never seen before
- I ask you to read and try some things just on your own
- I am here to answer questions popping up
- but you are also here to answer questions:)

# Topics we will touch in this course

Setting up a data science environment

Learn to manage data, analyse data, visualise data, and model data

Learn the Python language and data science libraries - today's most used data science tools

Learn to use Jupyter notebooks and Jupyter lab as interactive work tool

Learn to use machine learning methods

Learn to create computational essays

# Why Python?

# Languages

for data processing

- We can use all of them.
- Some are better fitting.
- We need sometimes to combine them.
- There are more.



Products ~ Qu

Quality Models ~

Markets ~

**Coding Standards** 

Schedule a demo

TIOBE Index

Jan 2023	Jan 2022	Change	Program	ming Language	Ratings	Change						
<u>*</u> ***********************************												
1 *	1 		•	Python	16.36%	+2.78%						
2	2		9	С	16.26%	+3.82%						
3	4	^	<b>©</b>	C++	12.91%	+4.62%						
4	3	•	<u>*</u>	Java	12.21%	+1.55%						
5	5		<u>@</u>	C#	5.73%	+0.05%						
6	6		VB	Visual Basic	4.64%	-0.10%						
7	7		JS	JavaScript	2.87%	+0.78%						
8	9	^	SQL	SQL	2.50%	+0.70%						
9	8	•	ASM	Assembly language	1.60%	-0.25%						
10	11	^	php	PHP	1.39%	-0.00%						
11	10	•	<u> </u>	Swift	1.20%	-0.21%						
12	13	^	- <b>GO</b>	Go	1.14%	+0.10%						
13	12			D	1.0.40/	-0.21%						
•	12	••••••	R	R	1.04%	-0.21%						
14	15	^	450	Classic Visual Basic	0.98%	+0.01%						
15	16	^	<b></b>	MATLAB	0.91%	-0.05%						
16	18	^		Ruby	0.80%	-0.08%						

# Is it just a hype?

Long-term history of programming language popularity

Programming Language	2023	2018	2013	2008	2003	1998	1993	1988
Python	1	5	8	7	13	28	17	-
С	2	2	1	2	2	1	1	1
Java	3	1	2	1	1	17	-	-
C++	4	3	4	3	3	2	2	6
C#	5	4	5	8	12	-	-	-
Visual Basic	6	15	-	-	-	-	-	-
JavaScript	7	7	10	9	8	21	-	-
Assembly language	8	12	-	-	-	-	-	-
SQL	9	-	-	-	7	-	-	-
PHP	10	8	6	5	6	-	-	-
Objective-C	16	18	3	45	47	-	-	-
Ada	29	27	17	18	15	7	8	2
Lisp	31	31	13	15	14	9	5	3
Pascal	242	128	15	20	99	11	3	7
(Visual) Basic	-	-	7	4	4	3	6	5

https://www.tiobe.com/tiobe-index/

# How to set up a data science environment for use with Python?

Local computer

#### Desktop/Laptop/Server



Install Docker or Docker Desktop



Choose Jupyter Docker image

https://jupyter-docker-stacks.readthedocs.io/en/latest/



Run Jupyter notebook or Jupyter lab in your browser

#### Pro's:

- Image/Container with up-todate libraries
- Multiple instances
- Fast testing possible

#### Con's:

Need to learn Docker basics

#### CONDA

Install Miniconda or Anaconda



- 1. Create virtual environment
- 2. install needed 3rd-party libraries
- 3. install Jupyter



Run Jupyter notebook or Jupyter lab in your browser

#### Pro's:

- Full control on the installed libraries
- Multiple venv's possible

#### Con's:

- Need to learn venv basics
- Need to keep libraries manually up-to-date

#### cloud service

Google Colab





Access Google Colab via your browser

https://colab.research.google.com/



- 1. Create a GitHub repository to exchange data files
- 2. Google drive etc is also possible

#### Pro's:

- · No installation of Python
- · instantly available
- Fast testing possible

#### Con's:

- Need some repository for data in/out
- Need paid version if you want to use it persistently (better choice for resources like
  CPU/GPU)

### What we will use?

• We will use the Google Colab (cloud service)



• We also use a Github repository to organise the course



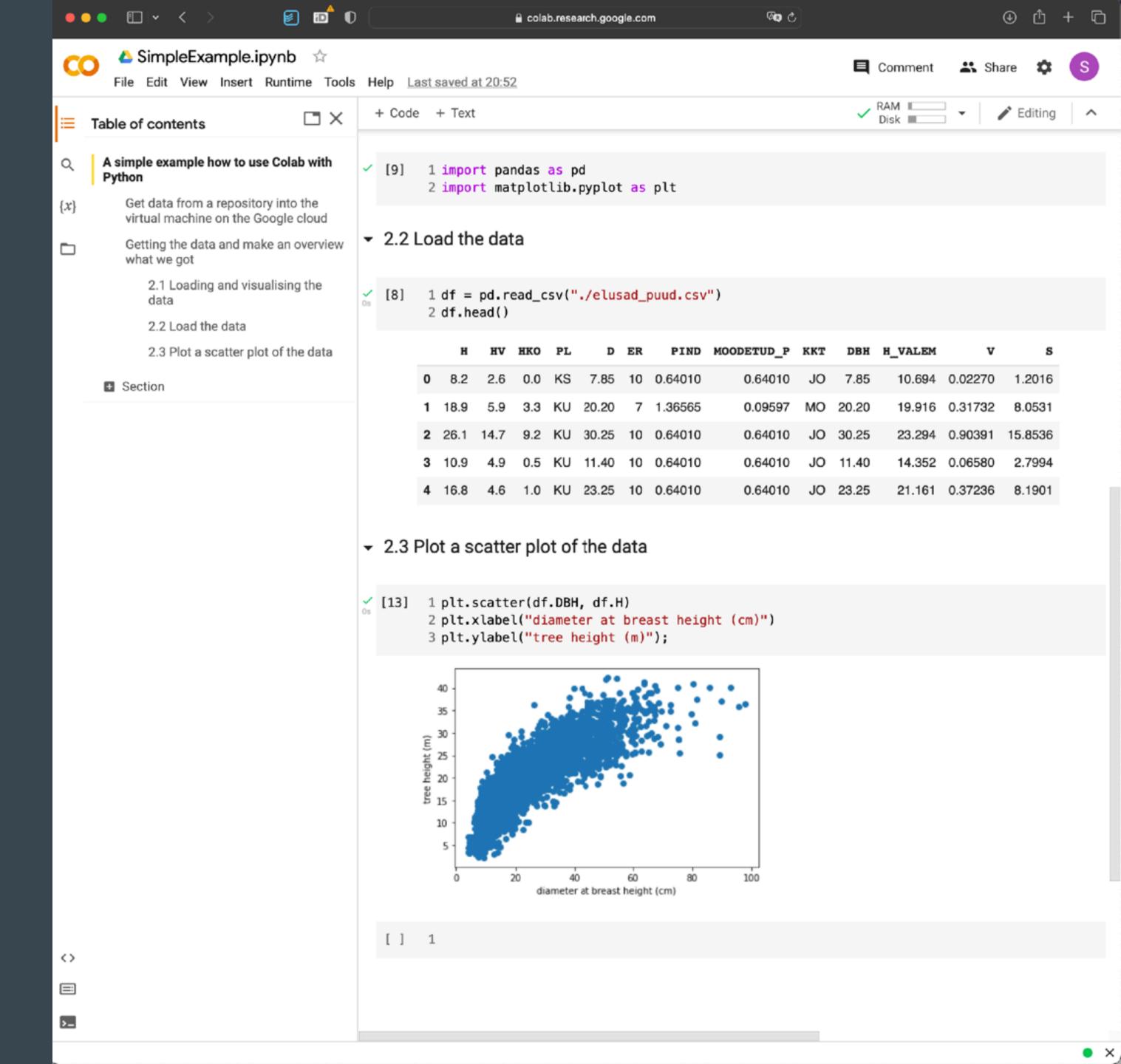
# Jupyter notebooks



The concept of computational documents

- Combines text and code
- Let you organise your data processing
- Let you visualise your results (fast)
- You have a good documentation of your work
- The paper is almost written on the fly  $\stackrel{ ext{@}}{=}$

Let's see an example



# Make it happen...

Instead of too much theory



we will cook



and experiment