

# efl Data Science Course

## Introduction

# About us: Who we are



## Original Mission:

- Investigate and co-shape Digital Finance 2.0
  - Web-based selfservices of customers
- Research was performed in three different Layers:
  - Customers in E-Finance
  - E-Financial Markets and Market Infrastructures
  - IT Infrastructures: Service Systems in E-Finance



## New Mission (Since 2019):

- Use expertise in Data Science to deliver cutting edge research in the fields of
  - Financial Services
  - Retail & Marketing
  - Health
  - Law
  - General, cross-sectional research

# About us: The efl



**Industry - Academic Partnership**

Universities



Sponsors



# About us: Lecturers



**M.Sc. Tino Cestonaro**

*Research Assistant*

Tino Cestonaro joined the team in April 2020. His research focuses on Market Microstructure, Financial Machine Learning,...



**M.Sc. Micha Bender**

*Research Assistant*

Micha Bender focuses on empirical financial market research. Furthermore, he is interested in the application of machine...



**M.Sc. Johannes Chen**

*Research Assistant*

Johannes Chen studied Business Informatics at the Technical University of Darmstadt. During his studies, he focused on...



**Dr. Benjamin M.  
Abdel-Karim**

*Alumnus*

Since March 2018, Benjamin M. Abdel-Karim is a research assistant of Prof. Dr. Oliver Hinz at the Chair of Information...



# Agenda

## Day 1 (Python Course) (30.10.2023)

**9:00 - 10:30 Uhr**

Python Basics  
Introduction and Primitive Data Types

**10:40 - 12:10 Uhr**

Data Structures  
Lists, Sets, Dictionaries

**13:30 - 15:00 Uhr**

Control Structures  
Loops (for, while), case distinction  
(if, else)

**15:15 - 16:45 Uhr**

Functions  
Structure of Functions and Application

## Day 2 (Python Course) (01.11.2023)

**9:00 - 10:30 Uhr**

Helpful functions for data processing  
Libraries: os, re, csv

**10:40 - 12:10 Uhr**

Data types and data structures  
Libraries: numpy, pandas

**13:30 - 15:00 Uhr**

Data import and visualization  
Libraries: csv (cont'd), matplotlib

**15:15 - 16:45 Uhr**

Outlook: Data Science  
Exemplary implementation of a KDD process

## Day 3 (Data Science) (03.11.2023)

**9:00 - 10:30 Uhr**

Introduction to Data Science  
Terminology and basic concepts

**10:40 - 12:10 Uhr**

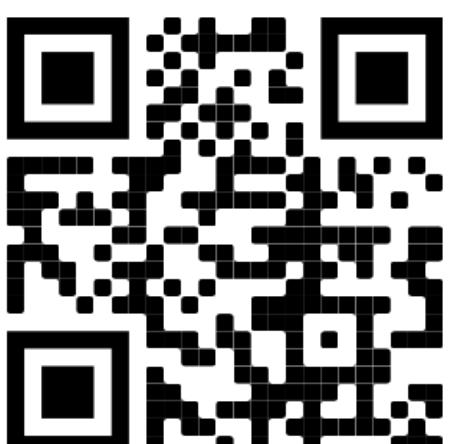
Working with Data  
Preprocessing, explorative data analysis

**13:30 - 15:00 Uhr**

Data Analysis I  
Classification

**15:15 - 16:45 Uhr**

Data Analysis II  
Neural Networks



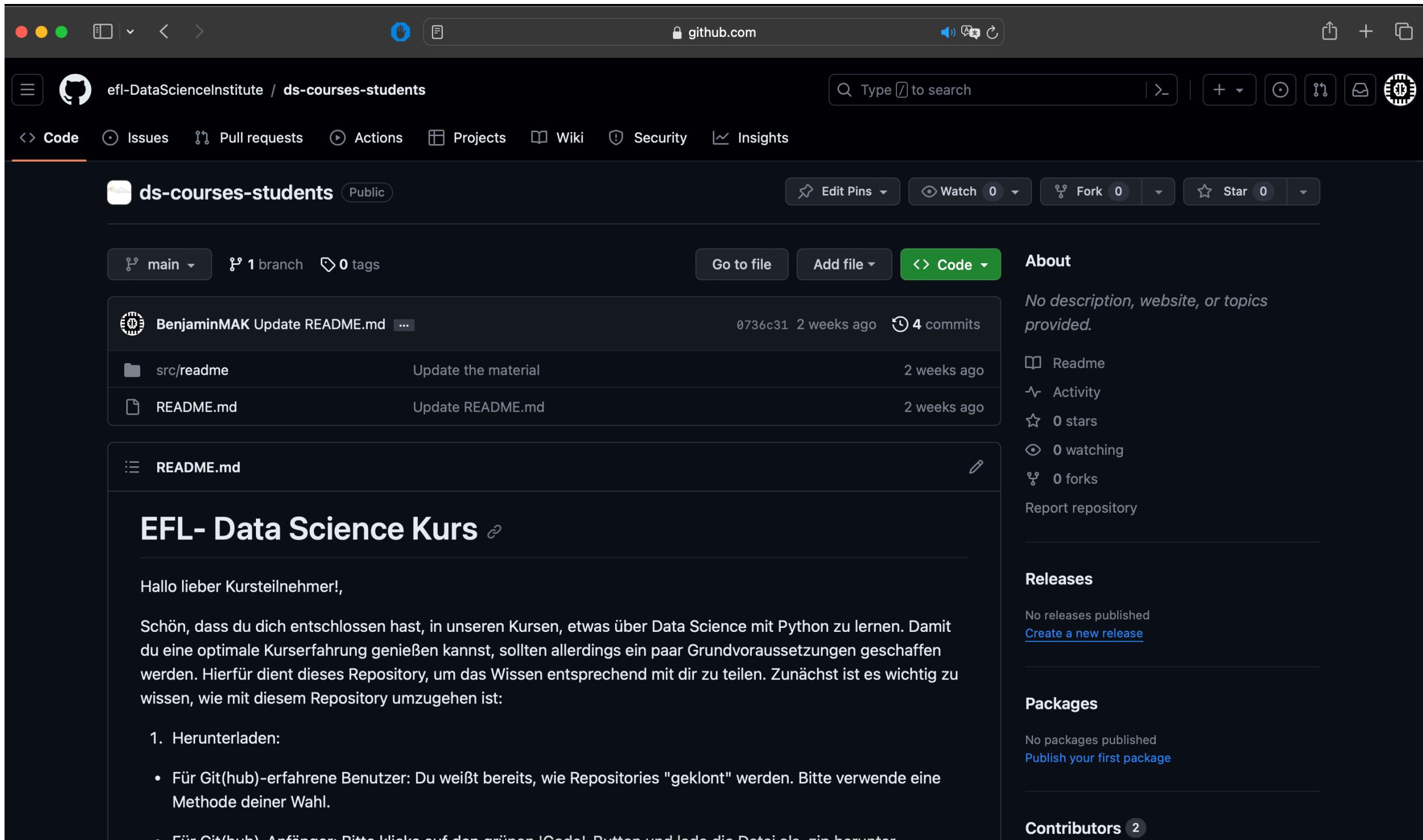
<https://www.eflab.de/teaching>

# Certificates



- Completion of tasks from the last class
- Submission: Description of how the tasks were solved
- Formalities:
  - Min. 2 pages (Arial 11, 1.5 Line spacing, 3 CM Correction margin )
  - Code must be delivered separately (code folder)

# Course Material?



The screenshot shows a GitHub repository page for 'efl-DataScienceInstitute / ds-courses-students'. The repository is public and contains one branch ('main') and one tag ('0.1'). The README.md file is displayed, containing the following content:

```

EFL- Data Science Kurs

Hallo lieber Kursteilnehmer!,  

Schön, dass du dich entschlossen hast, in unseren Kursen, etwas über Data Science mit Python zu lernen. Damit du eine optimale Kurserfahrung genießen kannst, sollten allerdings ein paar Grundvoraussetzungen geschaffen werden. Hierfür dient dieses Repository, um das Wissen entsprechend mit dir zu teilen. Zunächst ist es wichtig zu wissen, wie mit diesem Repository umzugehen ist:  

1. Herunterladen:  

  • Für Git(hub)-erfahrene Benutzer: Du weißt bereits, wie Repositories "geklont" werden. Bitte verwende eine Methode deiner Wahl.  

  • Für Git(hub)-Anfänger: Bitte klicke auf den grünen 'Code'-Button und lade die Datei als .zip herunter.

```

The repository has no description, website, or topics provided. It has 4 commits from BenjaminMAK, 0 stars, 0 forks, and 0 watching.



<https://github.com/efl-DataScienceInstitute/ds-courses-students>

# Why you are here

Kenntnis in Python und anderen Softwareprogrammen werden auf dem heutigen Arbeits- und Praktikumsmarkt oft vorausgesetzt. Um auch über die in den Vorlesung vermittelten Inhalte hinaus etwas über statistische Softwareprogramme zulernen, möchte ich an dem Kurs teilnehmen.

This semester I am taking "Business Informatics" with Professor Dr. Peter Gomber, which brought you to my attention.

Since the importance of Data Science as a tool for business decision making has already been discussed in class, an interest has already been sparked here. Now I hope to gain a deeper understanding through your course.

Auffrischen von Python, neue Einblicke in Software

*Bei meiner aktuellen Werksstudentenstelle merke ich immer wieder, wie Entscheidungen sehr oft datengetrieben werden und welche Rolle gut aufgearbeitete Daten spielen. Daher würde ich gerne, um mich beruflich und persönlich weiterzubilden, gerne an dem Kurs teilnehmen.*

# Why Coding?



NACH ABSCHLUSS IN ...  
**INFORMATIK, MATHEMATIK,  
 WIRTSCHAFTSINFORMATIK**

## TOP BERUFSTITEL\*

IT Manager  
**92.275 €**

IT-Projektmanager/in  
**86.695 €**

Senior Developer  
**78.635 €**

Analyst  
**74.469 €**

Software-Architekt/in  
**74.152 €**

## TOP BRANCHE

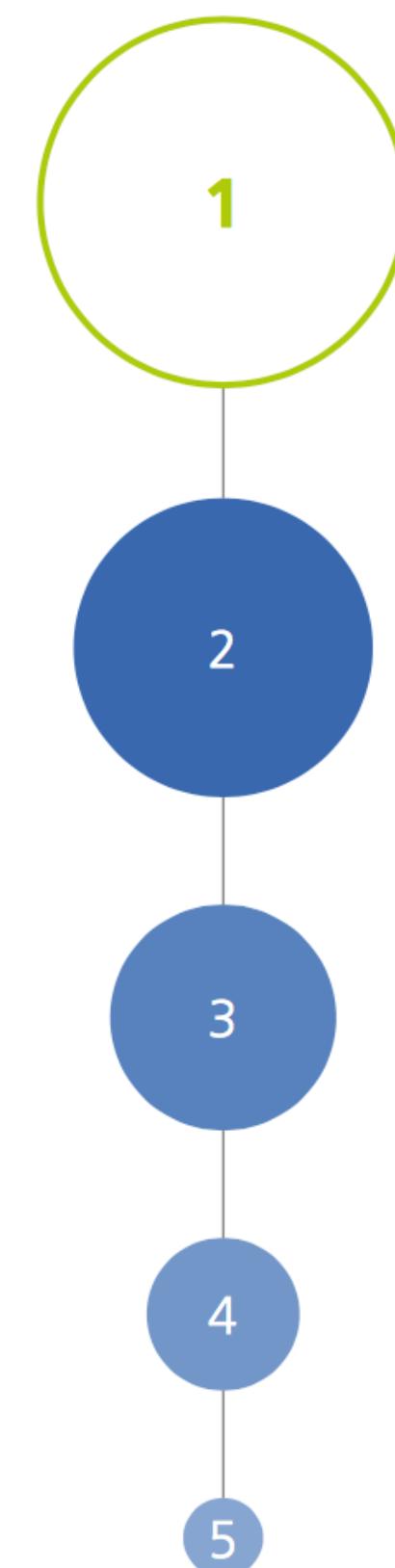
Banken  
**85.067 €**

Chemie- und Erdöl-  
 verarbeitende Industrie  
**84.033 €**

Konsum- und  
 Gebrauchsgüter  
**82.586 €**

Unternehmensberatung,  
 Wirtschaftsprüfung und Recht  
**82.336 €**

Finanzdienstleister  
**81.331 €**


\*Berufstitel ohne Management- & Personalverantwortung  
 Bruttodurchschnittsgehalt (inkl. variabler Bezüge), Mittelwert

# What is coding?

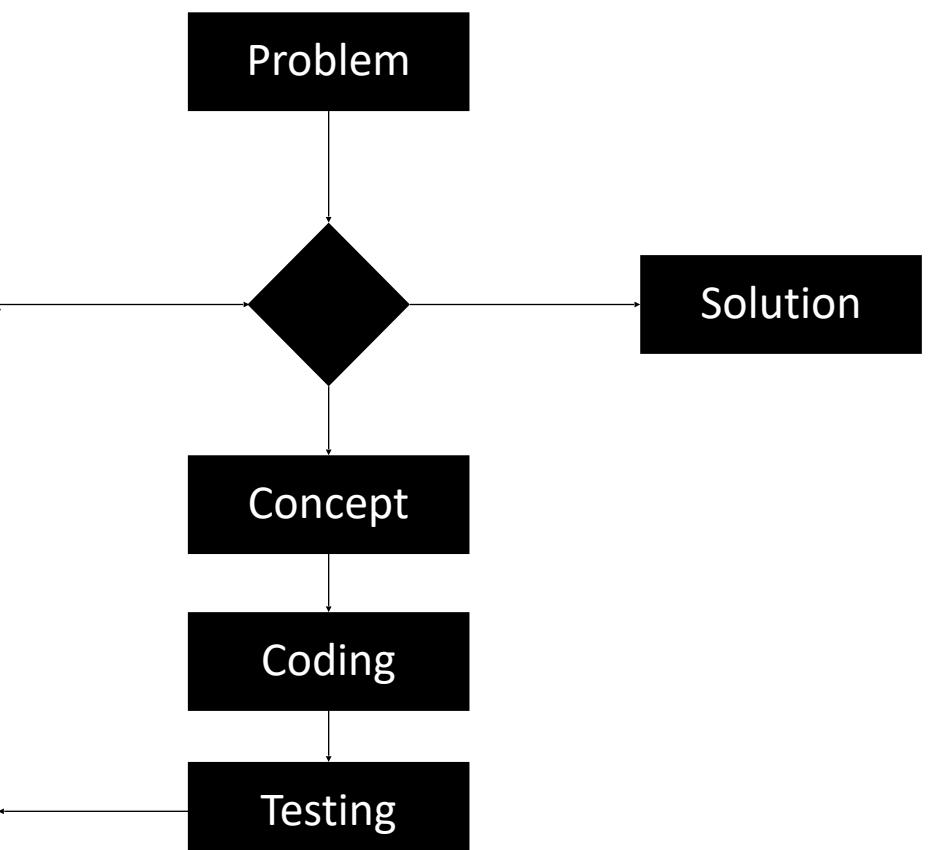


## Computer programming

Article [Talk](#)

From Wikipedia, the free encyclopedia

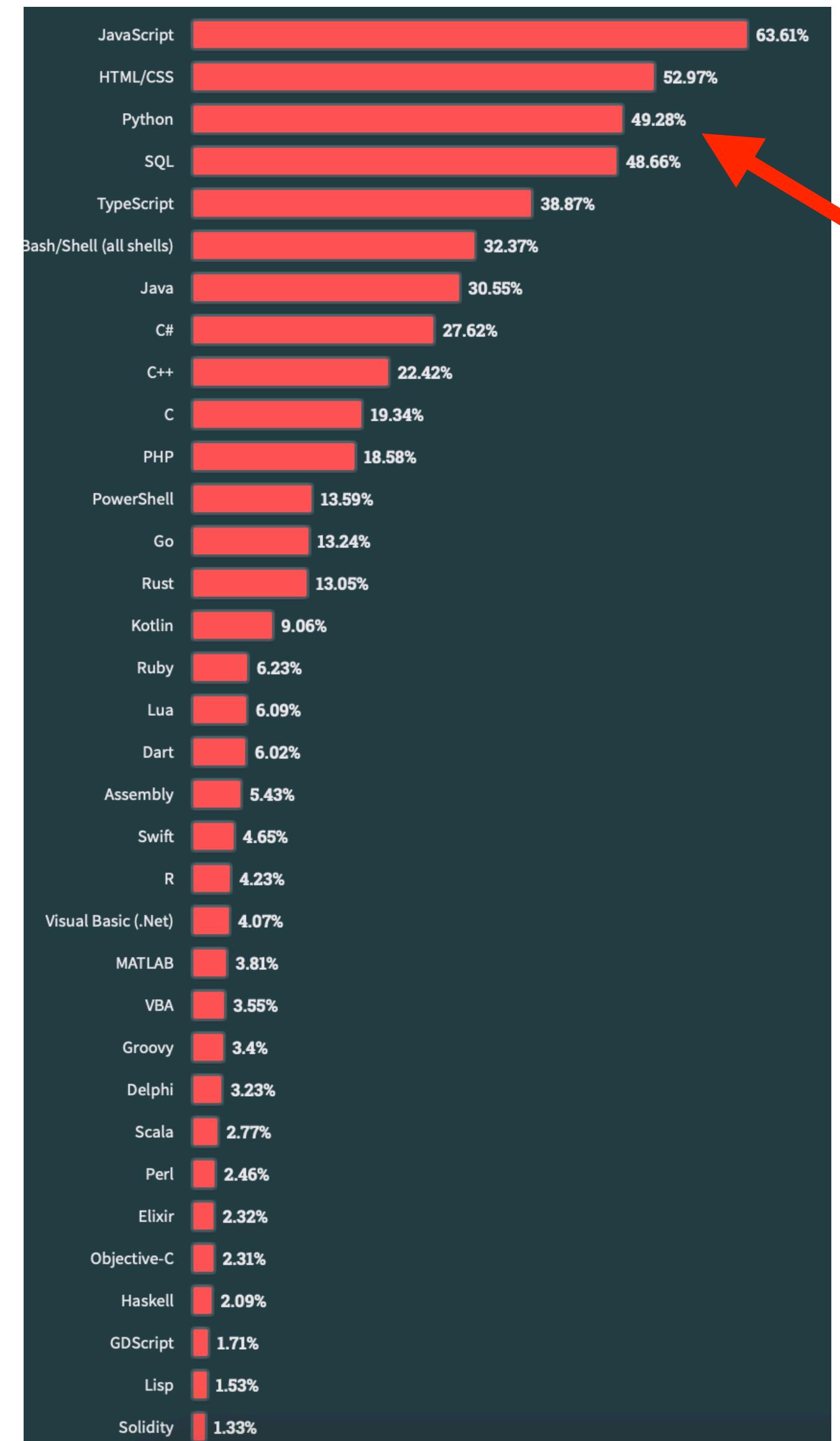
**Computer programming** or **coding** is the composition of sequences of instructions, called **programs**, that computers can follow to perform tasks.<sup>[1][2]</sup> It involves designing and implementing **algorithms**, step-by-step specifications of procedures, by writing **code** in one or more **programming languages**. Programmers typically use **high-level programming languages** that are more easily intelligible to humans than **machine code**, which is directly executed by the **central processing unit**. Proficient programming usually requires expertise in several different subjects, including knowledge of the **application domain**, details of programming languages and generic code **libraries**, specialized algorithms, and formal **logic**.



[https://en.wikipedia.org/wiki/Computer\\_programming](https://en.wikipedia.org/wiki/Computer_programming)

# Motivation: Why Python?

- Released in 1991 by Guido van Rossum
- With the explosive growth of ‘big data’ in disciplines such as bioinformatics, neuroscience and astronomy, programming know-how is becoming ever more crucial (Perkel 2015, p. 125).



<https://survey.stackoverflow.co/2023/>

# Motivation: The Elementary Basics in Python

Pros Python:

- General purpose language
- Higher programming language
- It's simple
- Fast to read
- Structuring by indenting
- No {} or ; => Faster to Code
- Data types are managed dynamically. There is no static type check like in java
- Widespread in science
- Extensive Support Libraries (important data science, math and many more)
- Integration Feature
- Productivity (Many Frameworks such as unit testing)



Ozgur, C., Colliau, T., Rogers, G., & Hughes, Z. (2017). MatLab vs. Python vs. R. Journal of data Science, 15(3), 355-371.  
Abdel-Karim, B. M. (2022). Data Science. In Data Science: Best Practices mit Python (pp. 57-62). Wiesbaden: Springer Fachmedien Wiesbaden.

# Cons of Python

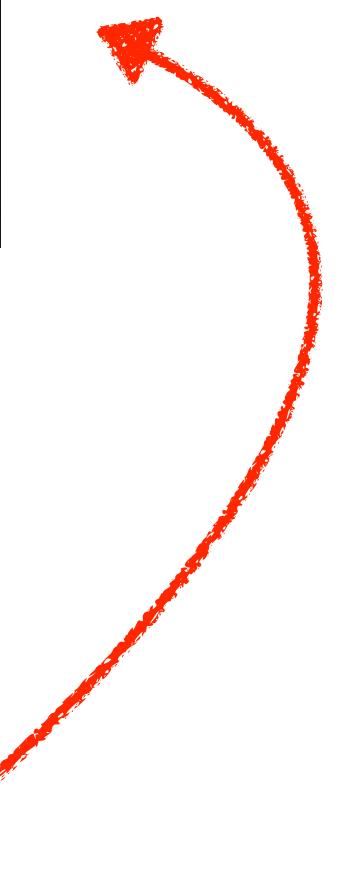
Some Points:

- Python is an interpreted language, this is slow as compared to C/C++
- See Interpreted vs. Compiled programming languages
- Python can have runtime errors (dynamical typing feature)
- Consumes a lot of memory space
- Not easy to test



Small Satellite\*

\*) <https://www.nbcnews.com/id/wbna43313086>



Not in Python ...

Prechelt, L. (2000). An empirical comparison of c, c++, java, perl, python, rex and tcl. *IEEE Computer*, 33(10), 23-29.

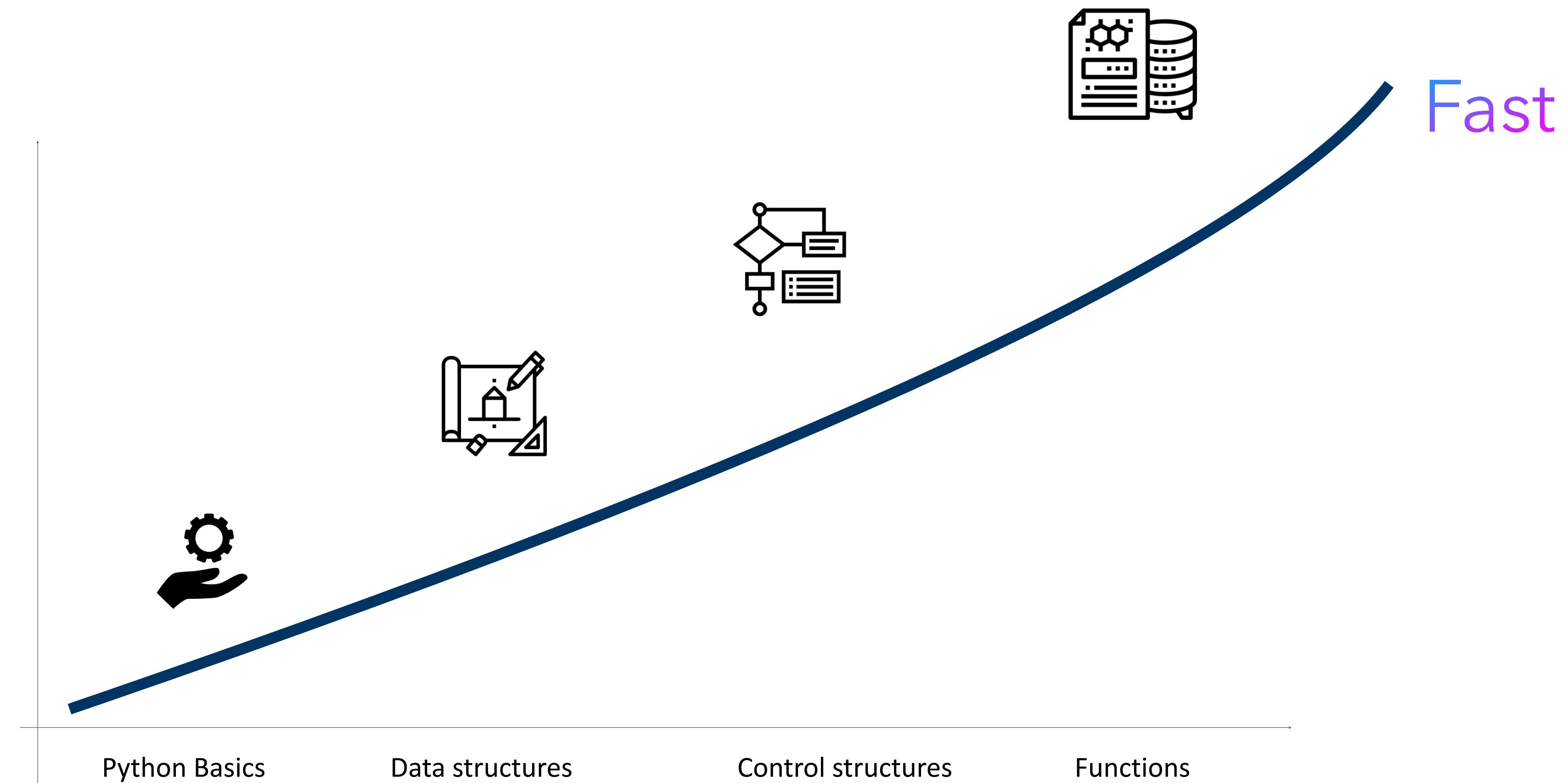
Newhall, T. K. (1999). Performance Measurement of Interpreted, Just-in-Time compiled, and Dynamically Compiled Executions. The University of Wisconsin-Madison.

# Agreement

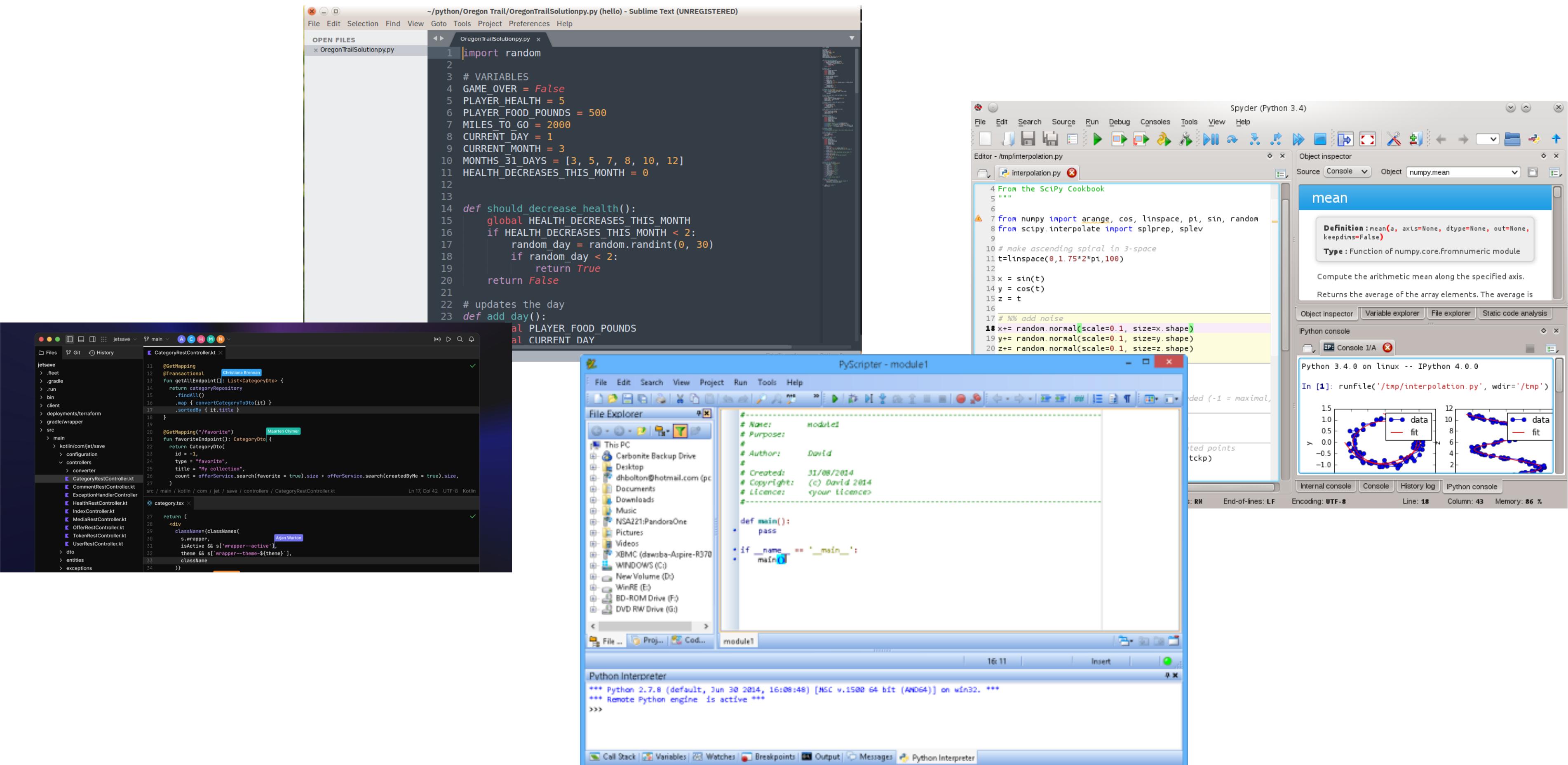


- Introduction to programming => All Questions are allowed!!!
- For beginners
- The module is interactive => Use your computer
- **We develop the solutions together!**
- Please be on time
- There are no dumb questions
- Nobody knows everything
- Copying solutions is plagiarism

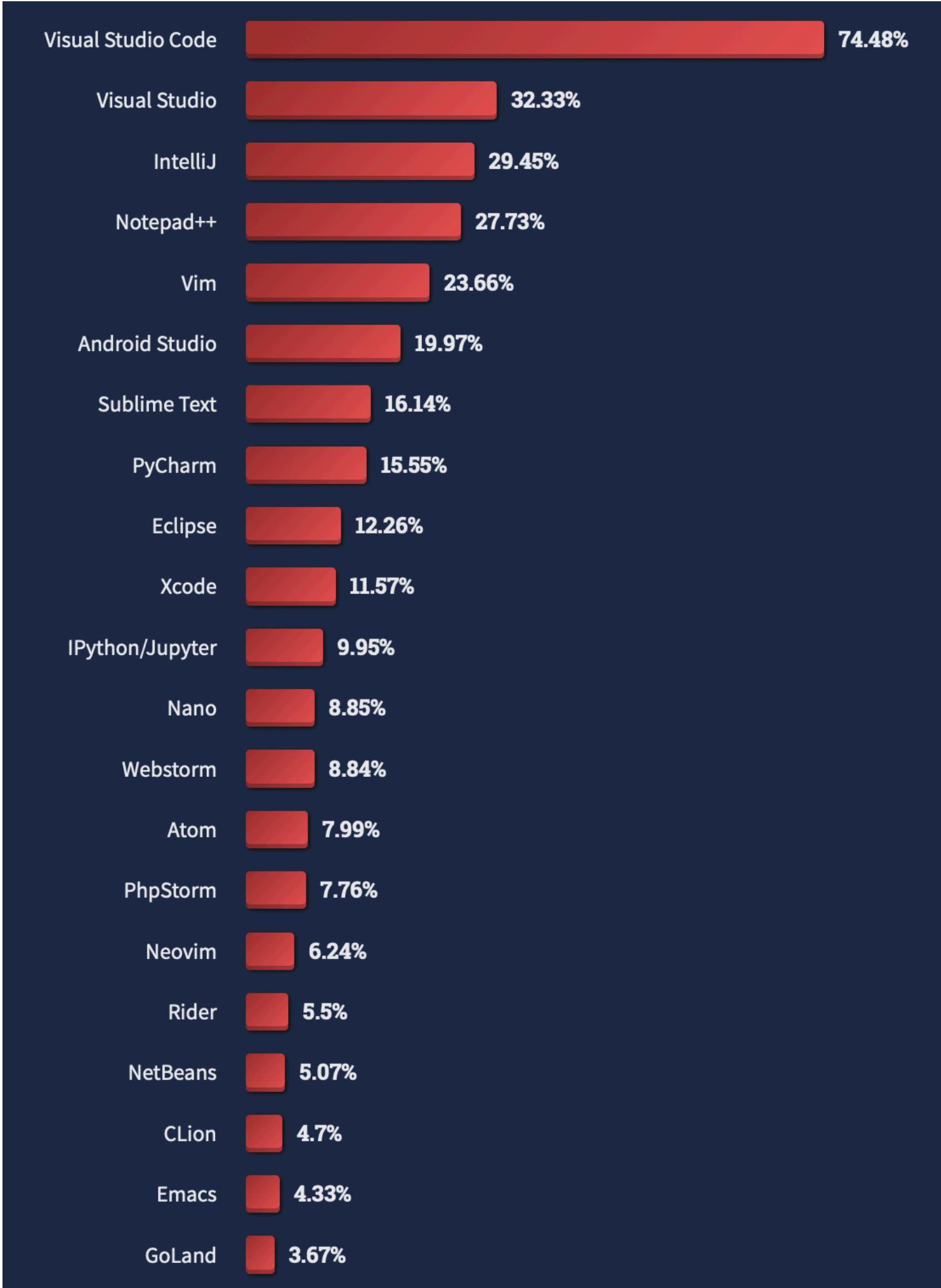
# Note: Learning Curve for Today



# Integrated Development Environment IDE a Deeper look? The Problem of Choice



# The Problem of Choice

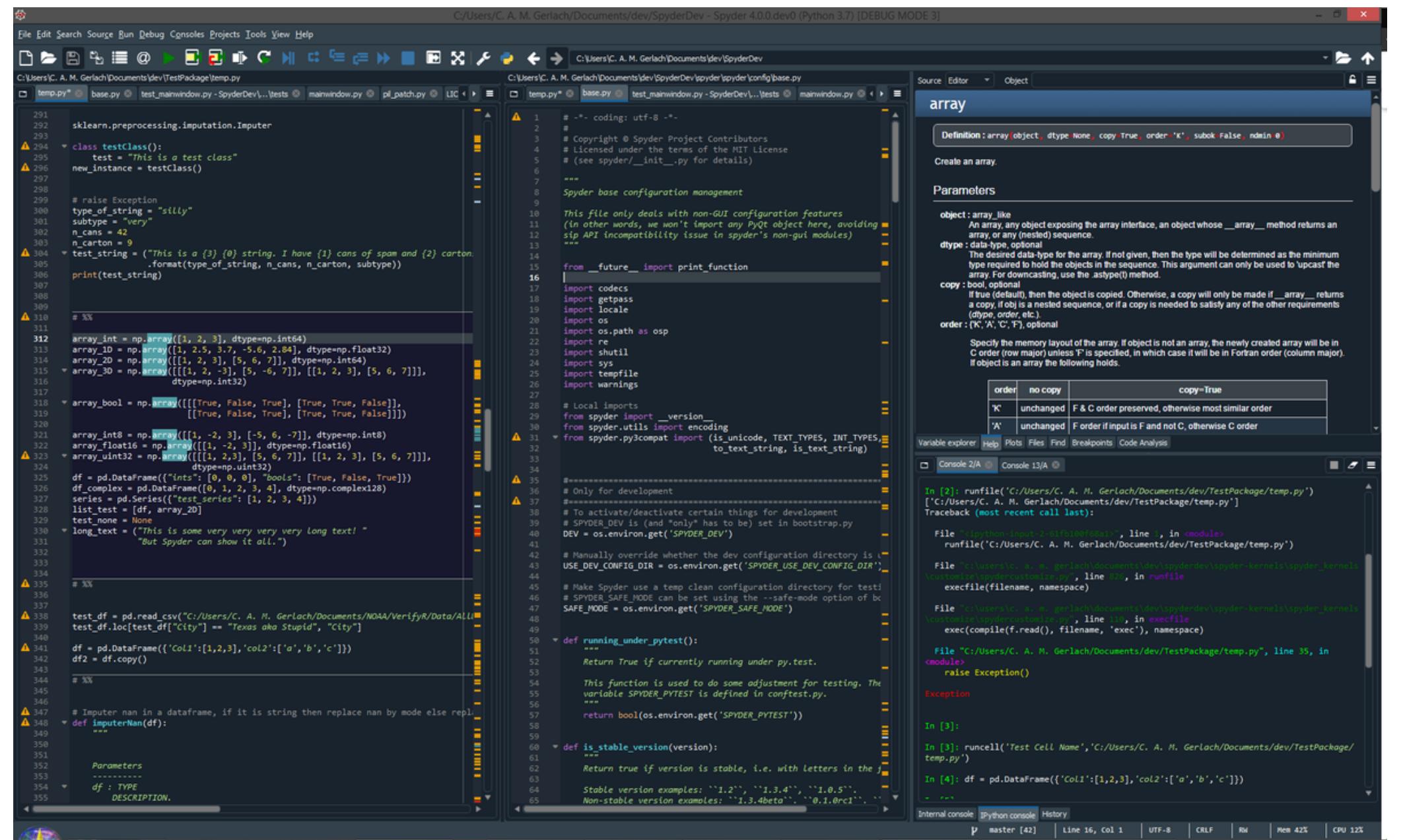


<https://survey.stackoverflow.co/2022/#most-popular-technologies-new-collab-tools-prof>

# What We Use for This Course



# SPYDER



The screenshot shows the Spyder 4.0.0.dev0 Python 3.7 [DEBUG MODE 3] interface. The code editor displays a script named `temp.py` containing Python code related to array operations and configuration management. The documentation viewer on the right provides detailed information about the `array` class, including its definition, parameters, and usage examples. The console window at the bottom shows the output of running the script, including tracebacks and execution results.

```

File Edit Search Source Run Debug Consoles Projects Tools View Help
C:/Users/C. A. M. Gerlach/Documents/dev/spyderDev - Spyder 4.0.0.dev0 (Python 3.7) [DEBUG MODE 3]
C:/Users/C. A. M. Gerlach/Documents/dev/spyderDev/spyder/config/base.py
temp.py base.py test_mainwindow.py-SpyderDev-.Tests mainwindow.py pil_patch.py LIC
temp.py base.py test_mainwindow.py-SpyderDev-.Tests mainwindow.py

array
Definition: array(object, dtype=None, copy=True, order='K', subok=False, ndmin=0)
Create an array.

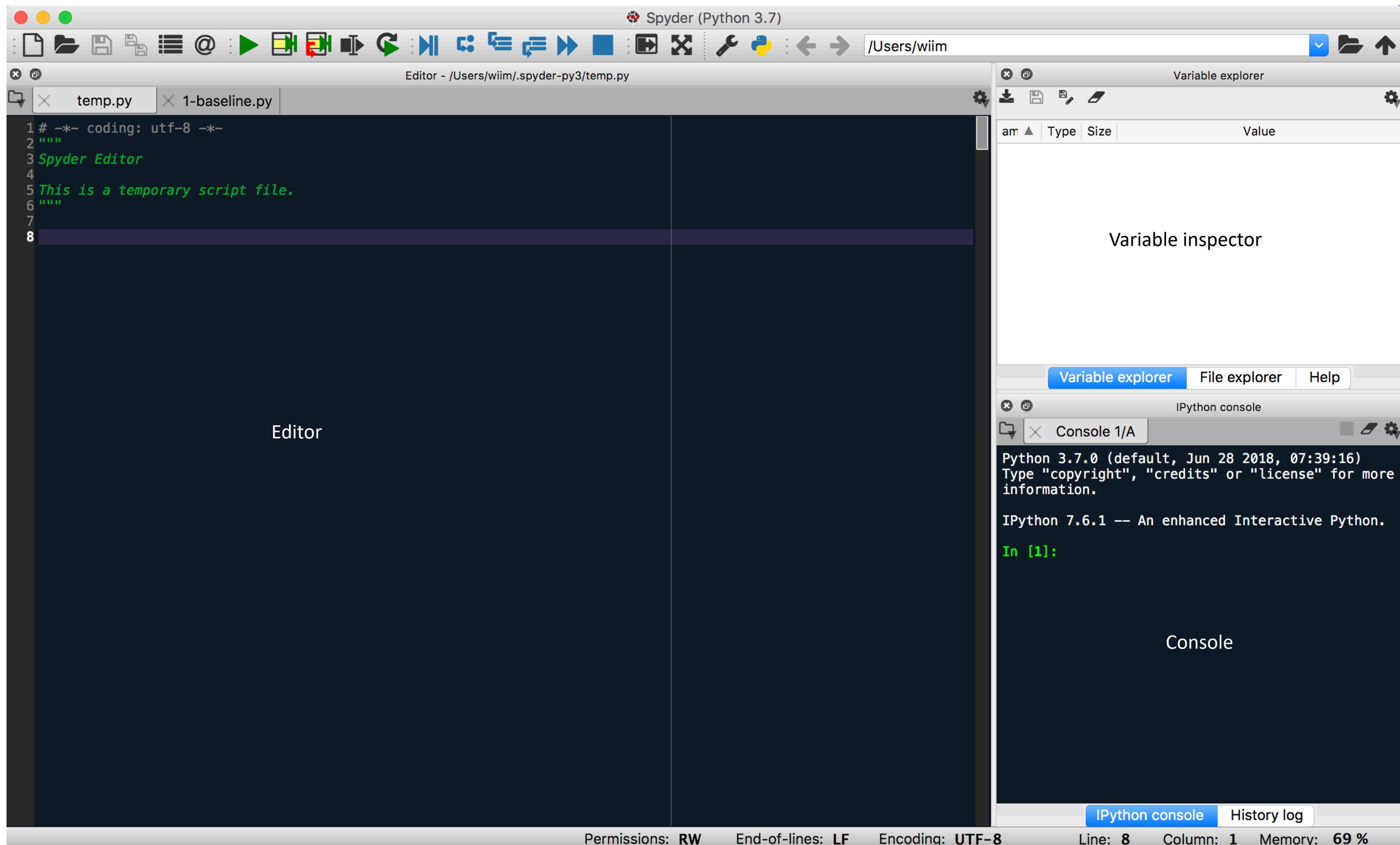
Parameters
object : array_like
An array, any object exposing the array interface, an object whose __array__ method returns an array, or any (nested) sequence.
dtype : data-type, optional
The desired data type for the array. If not given, then the type will be determined as the minimum type required to hold the objects in the sequence. This argument can only be used to 'upcast' the array. For 'downcasting', use the astype() method.
copy : bool, optional
If False, then the object is copied. Otherwise, a copy will only be made if __array__ returns a copy, if it is a nested sequence, or if a copy is needed to satisfy any of the other requirements (datatype, order, etc.).
order : {'C', 'F', 'A'}, optional
Specify the memory layout of the array. If object is not an array, the newly created array will be in C order (row major) unless 'F' is specified, in which case it will be in Fortran order (column major). If object is an array the following holds:
order      no copy          copy=True
          'C'      unchanged    F & C order preserved, otherwise most similar order
          'A'      unchanged    Order if input is F and not C, otherwise C order
Variable explorer Help Plots Files Find Breakpoints Code Analysis
Console 2/A Console 13/A
In [2]: runfile('C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py')
[C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py]
Traceback (most recent call last):
File "C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py", line 1, in <module>
  runfile('C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py')
File "C:/Program Files/Anaconda3/lib/python3.7/site-packages/spyder/utils/widgets/runfile/runfile.py", line 105, in runfile
  execfile(filename, namespace)
File "C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py", line 1, in <module>
  execfile(filename, namespace)
File "C:/Program Files/Anaconda3/lib/python3.7/site-packages/spyder/utils/widgets/runfile/runfile.py", line 105, in execfile
  exec(compile(f.read(), filename, 'exec'), namespace)
File "C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py", line 35, in <module>
  raise Exception()
Exception
In [3]:
In [3]: runcell('Test Cell Name','C:/Users/C. A. M. Gerlach/Documents/dev/TestPackage/temp.py')
In [4]: df = pd.DataFrame(['Col1':[1,2,3], 'Col2':['a','b','c']])
...
Internal console Python console History
P master [4] Line 16, Col 1 UTF-8 CR LF RD Mem 42% CPU 12%

```

# Practice and Questions

- Take your computer and let's get started!
- Create a Folder 'Day 1' on our Desktop
- Launch Spyder and select the folder

# In a Nutshell: Integrated Development Environment

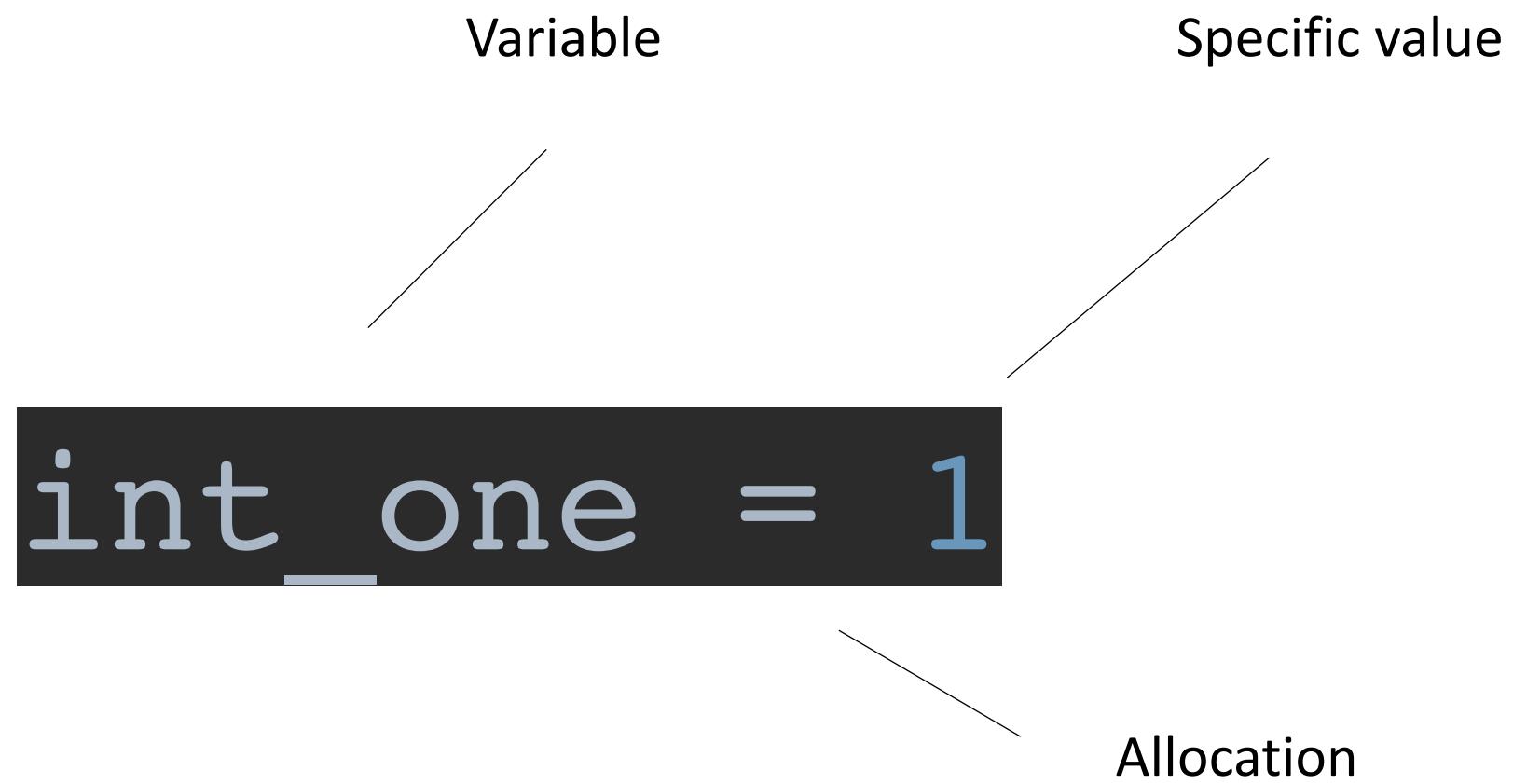


# In a Nutshell: Hello World

```
print('Hello World')
```

# Abstraction (Computer Science)

## The first step - The concept of variable



We assign the value 1 to the variable 'one'.

Now we can continue working with the variable 'one'.

Advantage: We are independent of concrete value

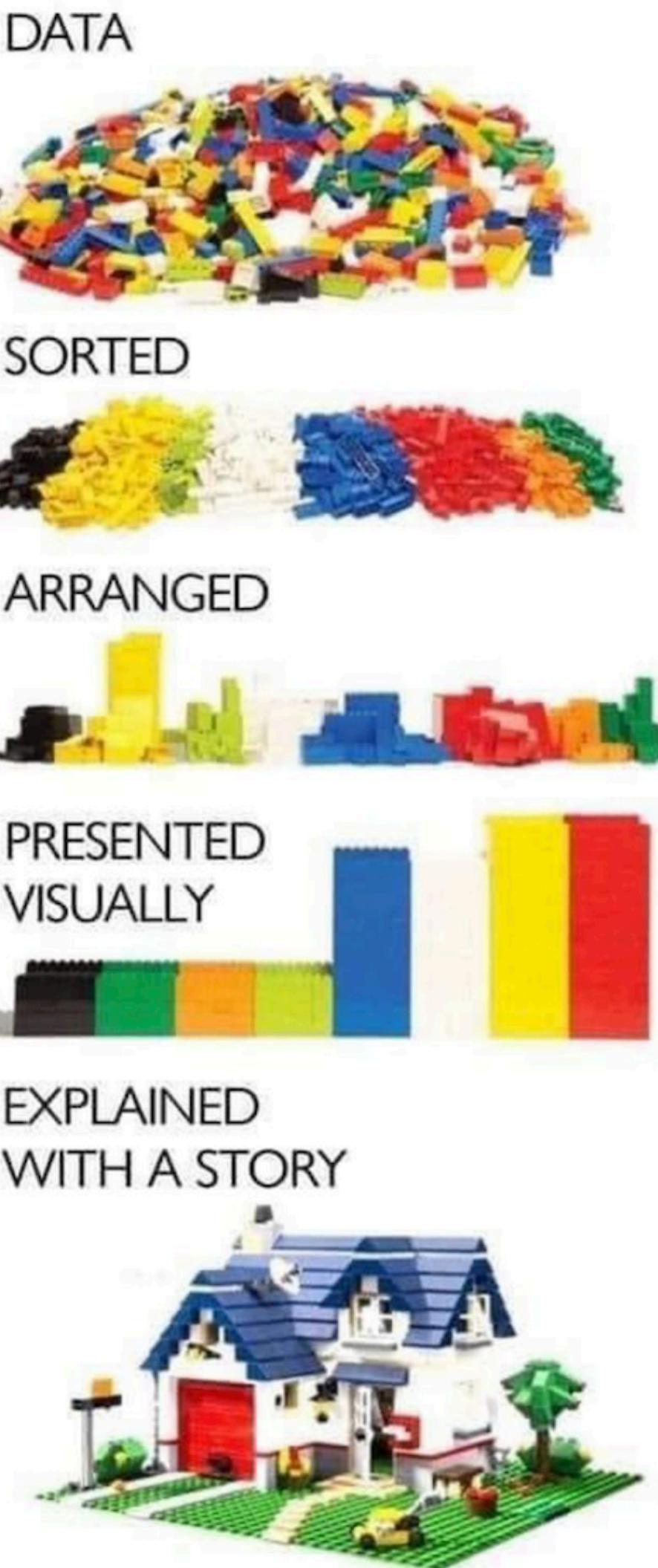
*The essence of abstractions is preserving information that is relevant in a given context, and forgetting information that is irrelevant in that context.*

– John V. Guttag



# Different Building Blocks

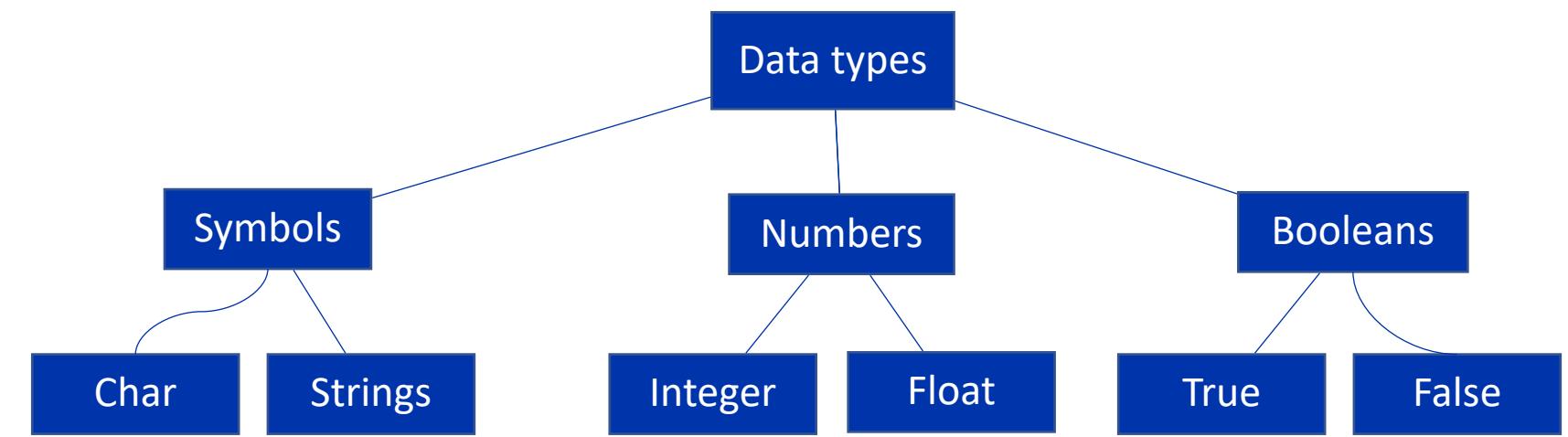
- Integer: int (signed integers) – They are often just called integers or ints, are positive or negative whole numbers with no decimal point.
- Float: Float (floating point real values, double) – Also called floats, they represent real numbers and are written with a decimal point dividing the integer and fractional parts. Floats may also be in scientific notation, with E or e indicating the power of 10.
- Strings are a sequence of chars. We can create them simply by enclosing characters in quotes. “Hello World” is a String! Therefore strings in Python are bytes representing Unicode characters. In Detail: Python does not have a character data type, a single character is simply a string with a length of 1.
- Boolean: We can perform logical operations with True and False in combination with AND and OR.



# Primitive Types

- Integer: int (signed integers) – They are often just called integers or ints, are positive or negative whole numbers with no decimal point.
- Float: Float (floating point real values, double) – Also called floats, they represent real numbers and are written with a decimal point dividing the integer and fractional parts. Floats may also be in scientific notation, with E or e indicating the power of 10.
- Strings are a sequence of chars. We can create them simply by enclosing characters in quotes. “Hello World” is a String! Therefore strings in Python are bytes representing Unicode characters. In Detail: Python does not have a character data type, a single character is simply a string with a length of 1.
- Boolean: We can perform logical operations with True and False in combination with AND and OR.

# Primitive Data Types: A first overview



\* Simplified Illustration

Python Documentation: <https://docs.python.org/3/library/stdtypes.html>

Other classification: Chun, W. (2001). Core python programming (Vol. 1). Prentice Hall Professional.

# Documentation

## Source Code Documentation

@author: name

@since: first implementation date

@version: date of last update

@source: if you using links etc.

@code: special code note

@param: if special parameter is used or you have to describe.

```
# I am a comment
@author: My Name
@since: 2022-10-03
@update: 2023-10-26
@version: v.0.0.2
print('Hello World')
```

More Information: PEP 8 – Style Guide for Python Code: <https://peps.python.org/pep-0008/>

# Naming convention

- Names of attributes, variables, methods start with a small letter
  - may use letters without ß or similar
  - which points to the data type like i, s or l
- This is standard in professional software development.
- Camel Case: Compound words are written in programming language. Every new word is capitalized.

(More PeP8 Style)

- name = Is the name of...
- bscript\_main = Simple code file that does something
- CName = Class (later more)
- float\_name = Variable that saves a floating point (double)
- int\_name = Variable that saves an integer value
- str\_name = Variable that saves a string value
- b\_name = Boolean for true or false values
- list\_name = Object from type list
- method\_name = Self-written function

(Camel Case Style)

- Name = Is the name of...
- bScriptName = Simple code file that does something
- CName = Class (later more)
- dName = Variable that saves a floating point (double)
- iName = Variable that saves an integer value
- sName = Variable that saves a string value
- bName = Boolean for true or false values
- LName = Object from type list
- fName = Self-written function

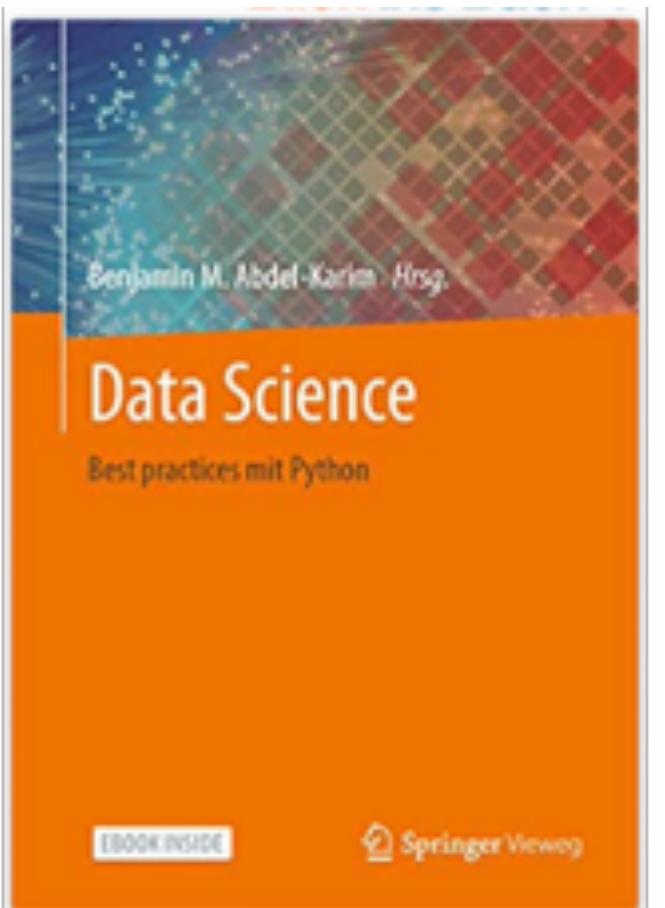
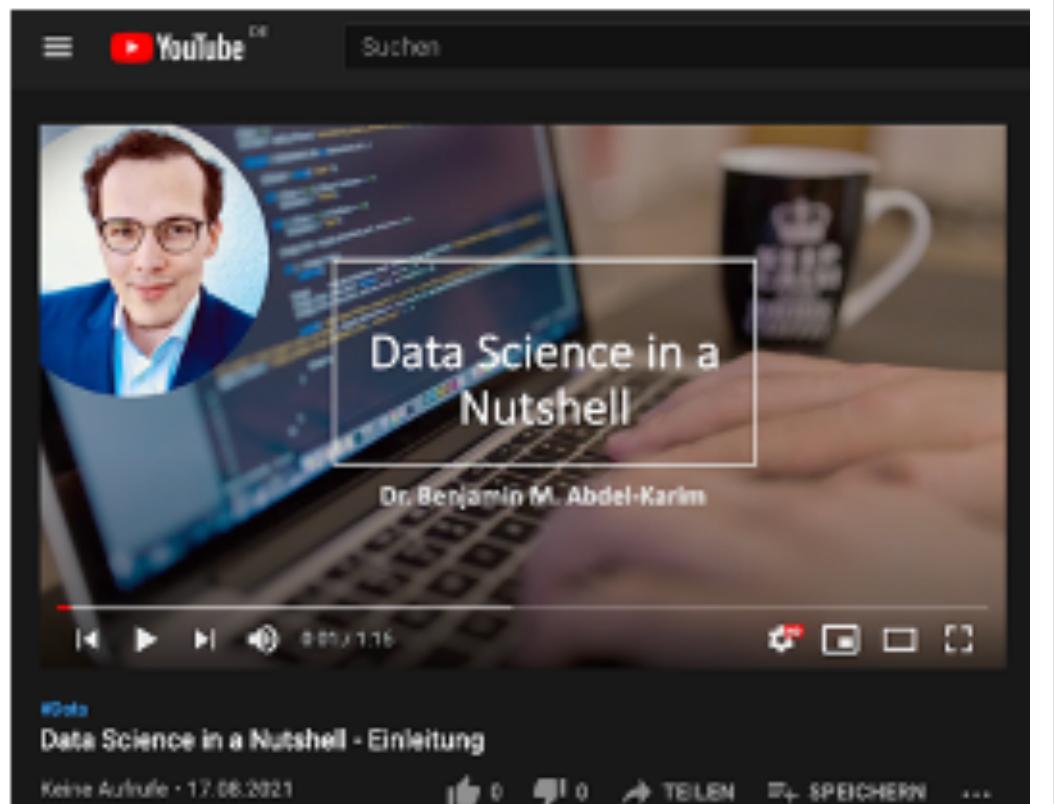
# Code is some kind of Art. Therefore...



Jackson Pollock Art

These are all approaches. Find your own style!

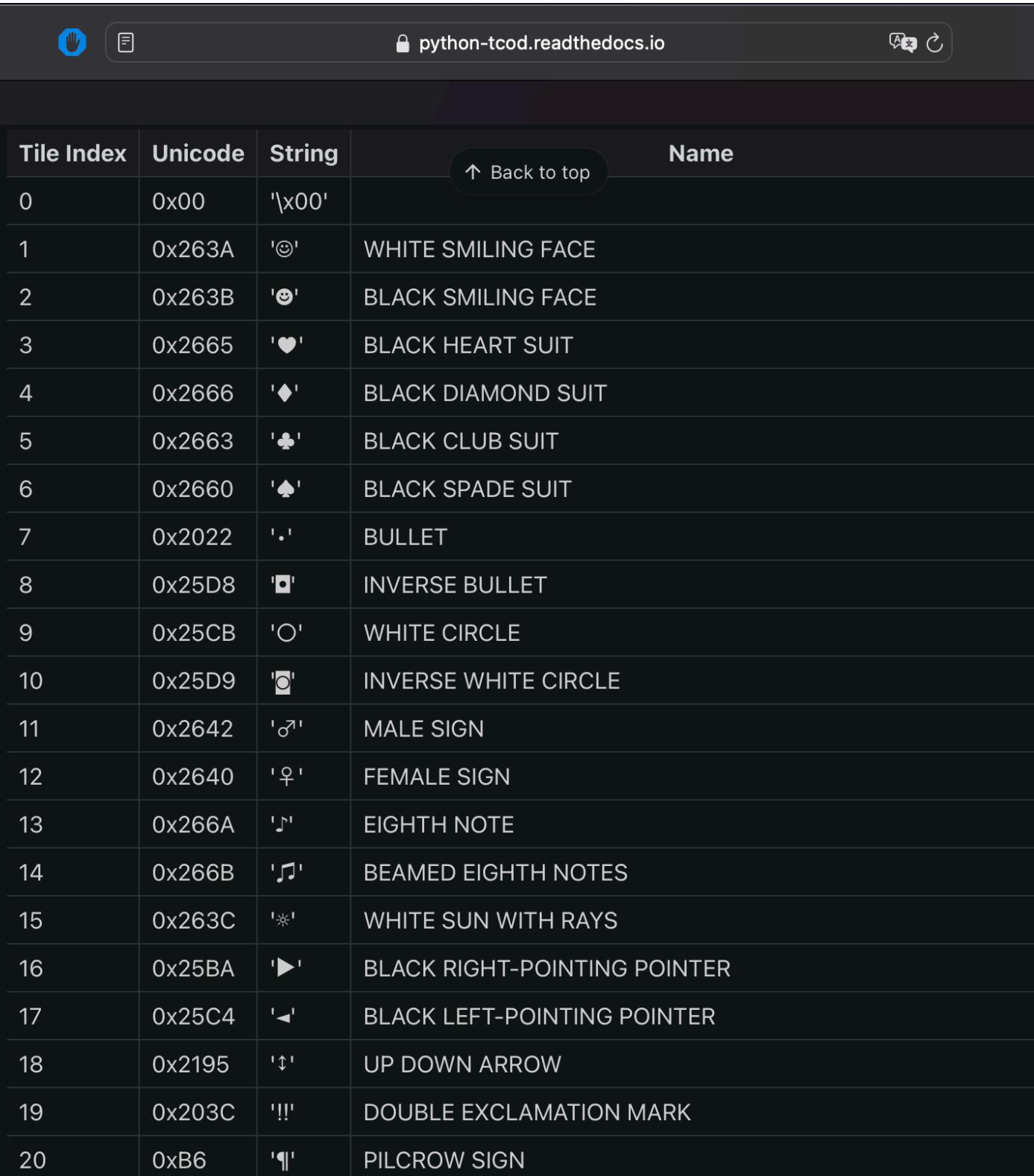
# Thank You



# Appendix

# Unicode Table

`ord(c) inverse chr()`



Tile Index	Unicode	String	Name
0	0x00	'\x00'	
1	0x263A	'☺'	WHITE SMILING FACE
2	0x263B	'☻'	BLACK SMILING FACE
3	0x2665	'♥'	BLACK HEART SUIT
4	0x2666	'♦'	BLACK DIAMOND SUIT
5	0x2663	'♣'	BLACK CLUB SUIT
6	0x2660	'♠'	BLACK SPADE SUIT
7	0x2022	'•'	BULLET
8	0x25D8	'▣'	INVERSE BULLET
9	0x25CB	'○'	WHITE CIRCLE
10	0x25D9	'◐'	INVERSE WHITE CIRCLE
11	0x2642	'♂'	MALE SIGN
12	0x2640	'♀'	FEMALE SIGN
13	0x266A	'♪'	EIGHTH NOTE
14	0x266B	'♫'	BEAMED EIGHTH NOTES
15	0x263C	'☀'	WHITE SUN WITH RAYS
16	0x25BA	'▶'	BLACK RIGHT-POINTING POINTER
17	0x25C4	'◀'	BLACK LEFT-POINTING POINTER
18	0x2195	'↕'	UP DOWN ARROW
19	0x203C	'❗'	DOUBLE EXCLAMATION MARK
20	0xB6	'¶'	PILCROW SIGN

<https://python-tcod.readthedocs.io/en/latest/tcod/charmap-reference.html>