



Artificial Intelligence | Basics of Algorithms & Application

Exercise 2: “*Introduction to Python*”

Timo Sturm & Dr. Dominik Jung

ki@is.tu-darmstadt.de

Prof. Dr. Peter Buxmann | Information Systems | Software & Digital Business

School of Business, Economics & Law

TU Darmstadt

1

Introduction to Python

- 1.1 Overview
- 1.2 Data Import & Management
- 1.3 Data Visualization & Exploration
- 1.4 Development Environments
- 1.5 Resources & Hands-On Exercises

1

Introduction to Python

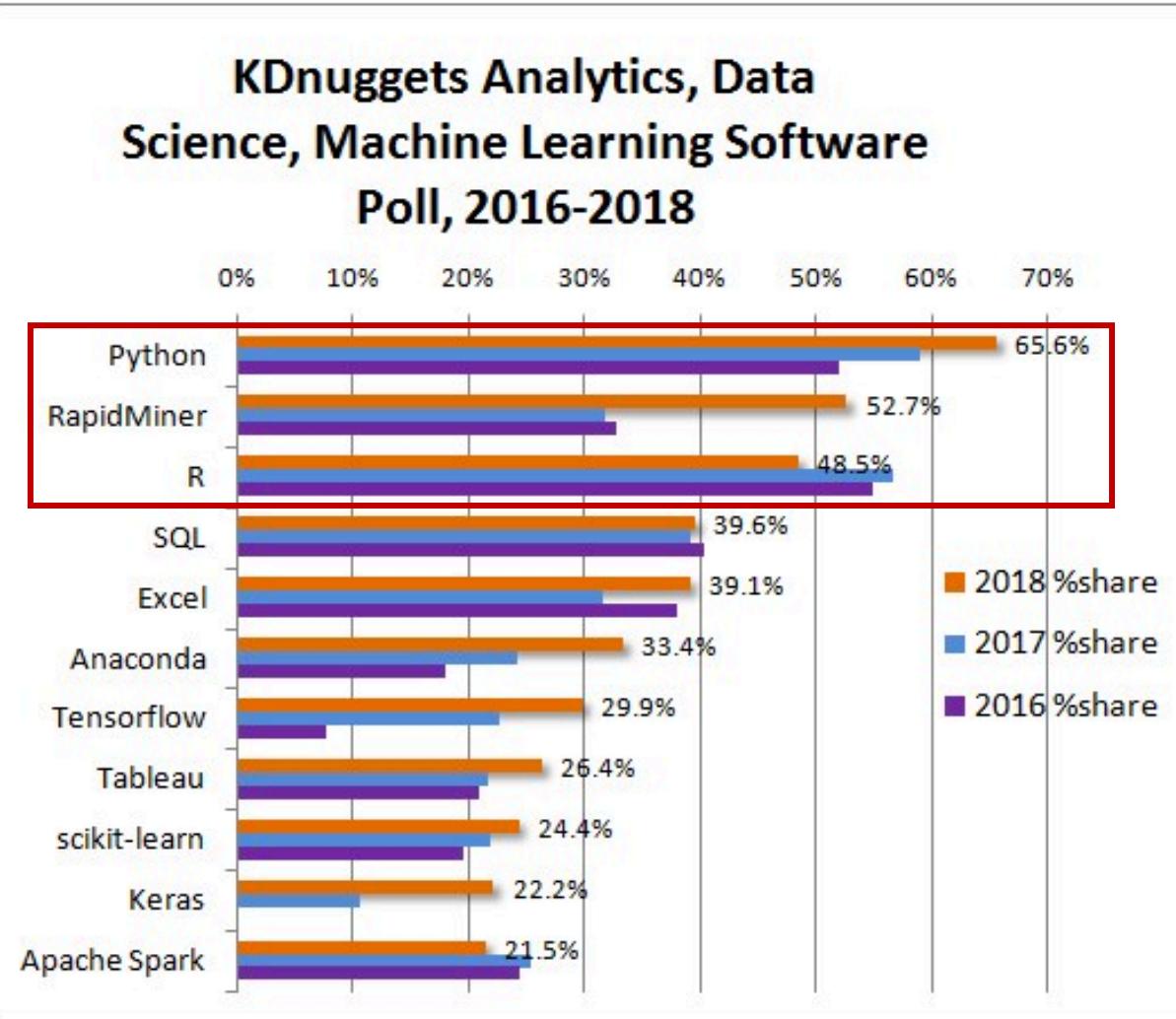
1.1 Overview

1.2 Data Import & Management

1.3 Data Visualization & Exploration

1.4 Development Environments

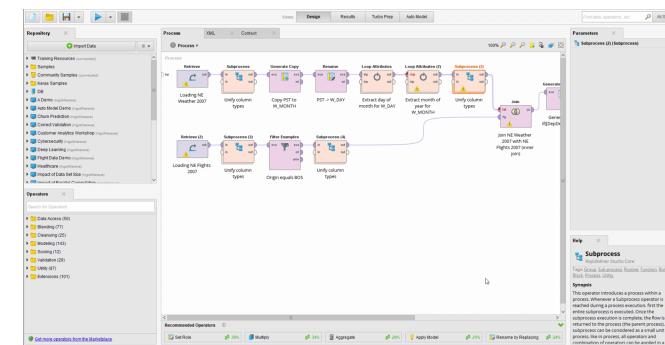
1.5 Resources & Hands-On Exercises



- Python & R currently represent the **most popular programming languages** in the data science area



- RapidMiner has become the **leading GUI-based tool** for conducting data science projects and its popularity is increasing even further



- Interpreted, high-level, general-purpose programming language
- **Origin:** Developed by *Guido van Rossum* in Amsterdam at the Dutch National Research Institute for Mathematics and Computer Science
 - **Goal:** Promote simplicity = a programming language should be *easy-to-use & easy-to-read*
 - **Multi-paradigm language:** Supports *object-oriented, structured, functional, aspect-oriented*, and more programming paradigms
 - **Integration of other programming languages:** If speed is important, time-critical functions can be moved to extension modules written in other programming languages, e.g., C
- Currently, **two Python versions** are used side by side by developers
 - **Latest Version:** Python 3
 - **Still widely used:** Python 2 (*official support will end after 2020!*)
 - With Python 3, many fundamental changes were introduced
 - Many libraries were not transferred to Python 3, so developers have to keep using Python 2 when they want to use such libraries in their Python code
 - As a result, developers have to switch between both versions depending on their use of libraries



- **Goal:** Promote simplicity = a programming language should be *easy-to-use & easy-to-read*

Zen of Python

```
>>> import this

The Zen of Python, by Tim Peters

Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.
Complex is better than complicated.
Flat is better than nested.
Sparse is better than dense.
Readability counts.
Special cases aren't special enough to break the rules.
Although practicality beats purity.
Errors should never pass silently.
Unless explicitly silenced.
In the face of ambiguity, refuse the temptation to guess.
There should be one-- and preferably only one --obvious way to do it.
Although that way may not be obvious at first unless you're Dutch.
Now is better than never.
Although never is often better than *right* now.
If the implementation is hard to explain, it's a bad idea.
If the implementation is easy to explain, it may be a good idea.
Namespaces are one honking great idea -- let's do more of those!
```

I. Indentations

- In most programming languages, the indentation of code is only used to improve readability without affecting any semantics
 - In Python, the indentation matters as it is used to define code blocks!
 - **Example:** Only the left code snippet is correct:
 - **Why?**
 - In the **right code snippet**: The print function is on the same level as the if statement, so the consequent, which is aimed to be executed if the if-condition is true, is missing
 - In the **left code snippet**: The consequent is defined through the indentation of the print function

```
if(x < y):  
    print(y) ✓
```

```
if(x < y):  
print(y) ✗
```

II. Variables

- When declaring variables, no data type needs to be indicated explicitly
 - The data type is automatically determined based on the declared value
 - ```
a = 5
b = -0.537
c = "hello!"
```

## III. Comments

- Single-line comments are indicated with a hash sign: #
- Multi-line comments are indicated by surrounding the text with three apostrophes: """

```
this is a single-line comment
print("hello!")
'''
this multi-line comment
can cover more than just
one line
'''
```

| Value Type                                                           | Description                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    | Examples                                                                                                                                                                                                                                                                                                                                                                                                  |
|----------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>boolean</b>                                                       | <ul style="list-style-type: none"> <li>Represent either true or false</li> </ul>                                                                                                                                                                                                                                                                                                                                                                                                                                               | <i>True, False</i>                                                                                                                                                                                                                                                                                                                                                                                        |
| <b>Numeric</b><br><i>(int, long, float, complex)</i>                 | <ul style="list-style-type: none"> <li><b>int</b>: a whole number</li> <li><b>long</b>: an integer of unlimited size</li> <li><b>float</b>: a real number written with a decimal point</li> <li><b>complex</b>: a number of the form <math>a + bJ</math>, with:           <ul style="list-style-type: none"> <li><math>a</math> &amp; <math>b</math>: floats</li> <li><math>J</math>: square root of -1 (= imaginary number)</li> </ul> </li> </ul>                                                                            | <ul style="list-style-type: none"> <li><b>int</b>: 1, -5, -200, 5432, ...</li> <li><b>long</b>: 1, -5, -200, 5432, 43243944949494994, ...</li> <li><b>float</b>: 0.54, -0.63678, 0.0005, ...</li> <li><b>complex</b>: 2+3j, 5j, 100+6j, ...</li> </ul>                                                                                                                                                    |
| <b>String</b>                                                        | <ul style="list-style-type: none"> <li>Representation of text</li> </ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                       | <i>"hello", "a", "xyz43[4]", ...</i>                                                                                                                                                                                                                                                                                                                                                                      |
| <b>Collections (Arrays)</b><br><i>(list, tuple, set, dictionary)</i> | <ul style="list-style-type: none"> <li><b>list</b>: a collection which is <i>ordered</i> and <i>changeable</i>; <i>allows</i> duplicates</li> <li><b>tuple</b>: a collection which is <i>ordered</i> and <i>unchangeable</i>; <i>allows</i> duplicates</li> <li><b>set</b>: a collection which is <i>ordered</i> and <i>unindexed</i>; <i>does not allow</i> duplicates</li> <li><b>dictionary</b>: a collection which is <i>unordered</i>, <i>changeable</i>, and <i>indexed</i>; <i>does not allow</i> duplicates</li> </ul> | <ul style="list-style-type: none"> <li><b>list/set</b>: ["banana", "apple", "peach"], [1.1, 0.583, 0.001, 54.43], ...</li> <li><b>tuple</b>: ("banana", "apple", "peach"), (1.1, 0.583, 0.001, 54.43), ...</li> <li><b>dictionary</b>:<br/>           {<br/>             "firstName": "Ronald",<br/>             "lastName": "McDonald",<br/>             "yearOfBirth": 1981<br/>           }</li> </ul> |

- **def:** *Define a function*  
→ functions can return multiple values!

```
def myAwesomeFunction():
 x = 1
 y = 2
 return x,y

def main():
 a, b = myAwesomeFunction()
 print(b)
```

- **class:** *Define a class*

```
class student():
 def __init__(self, firstName, lastName, age):
 self.firstName = firstName
 self.lastName = lastName
 self.age = age
```

- **if/else:** *Define a conditional control flow*

```
if x == 42:
 print("I also like the Hitchhikers guide to the galaxy")
elif x < 35:
 x = x + 1
else:
 x = x * 2
```

- **for & while loop:** *Define a loop with some ending condition*

```
z = 0

for x in range(4,10):
 z = z + x
 print(z)

while z < 1000:
 print(z)
 z = z + 1
```

## I. Data Structures & Statistics

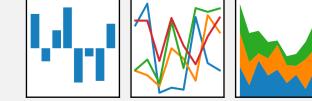
- *NumPy*
- *Pandas*
- ...



NumPy

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



## II. Visualization

- *Matplotlib*
- *Plotly*
- ...



plotly

## III. Machine Learning

- *Scikit-learn*
- *Tensorflow*
- *Keras*
- *XGBoost*
- *Pytorch*
- ...



Keras



→ Many more useful libraries  
for data science exist!

→ See, e.g., these (non-scientific) overview articles:

- Medium (2018): "[Top 20 Python libraries for data science in 2018](#)"
- Kdnuggets (2018): "[Top 10 Python Data Science Libraries](#)"
- Big Data Made Simple (2019): "[Top 20 Python libraries for Data Science](#)"
- ...

## 1

### Introduction to Python

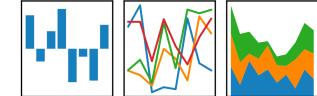
1.1 Overview

1.2 Data Import & Management

1.3 Data Visualization & Exploration

1.4 Development Environments

1.5 Resources & Hands-On Exercises



- The Python library *pandas* provides a rich set of data **import** and **export** capabilities
- Both **flat-file-based** (e.g., CSV, JSON, Excel, ...) & **server-based** (e.g., SQL, Google Big Query, ...) imports and exports are supported:

| Format Type | Data Description     | Reader         | Writer       |
|-------------|----------------------|----------------|--------------|
| text        | CSV                  | read_csv       | to_csv       |
| text        | JSON                 | read_json      | to_json      |
| text        | HTML                 | read_html      | to_html      |
| text        | Local clipboard      | read_clipboard | to_clipboard |
| binary      | MS Excel             | read_excel     | to_excel     |
| binary      | OpenDocument         | read_excel     |              |
| binary      | HDF5 Format          | read_hdf       | to_hdf       |
| binary      | Feather Format       | read_feather   | to_feather   |
| binary      | Parquet Format       | read_parquet   | to_parquet   |
| binary      | Msgpack              | read_msgpack   | to_msgpack   |
| binary      | Stata                | read_stata     | to_stata     |
| binary      | SAS                  | read_sas       |              |
| binary      | Python Pickle Format | read_pickle    | to_pickle    |
| SQL         | SQL                  | read_sql       | to_sql       |
| SQL         | Google Big Query     | read_gbq       | to_gbq       |

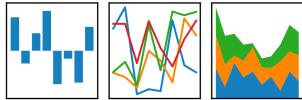
Source & more information: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/io.html](https://pandas.pydata.org/pandas-docs/stable/user_guide/io.html)

- Each import operator returns a pandas **dataFrame** (see next slide)
- **Example:** `read_csv()`

```
import pandas as pd

data = pd.read_csv('URL/To/Some/CSV/file.csv', delimiter=',')
print(data)
```

pandas  
 $y_{it} = \beta'x_{it} + \mu_i + \epsilon_{it}$



- **Result output example:**

|   | last_name | first_name | age | student? |
|---|-----------|------------|-----|----------|
| 0 | Skywalker | Luke       | 22  | True     |
| 1 | Targaryen | Daenerys   | 22  | False    |
| 2 | Weasley   | Ronald     | 13  | True     |
| 3 | White     | Walter     | 50  | False    |
- The import functions automatically identify fitting variable types for the values contained in the indicated file, e.g., in the example above:
  - all values of `last_name` & `first_name`: `str`
  - all values of `age`: `integer`
  - all values of `student?`: `boolean`

# Pandas: DataFrame – The primary data structure

- “**Two-dimensional size-mutable, potentially heterogeneous tabular data structure with labeled axes (rows and columns). [...]. The primary pandas data structure.**“
  - [Pandas DataFrame API Reference](#)

→ i.e.: a modifiable table-like data structure with namable columns and rows that can hold different data types (e.g., one column may hold integers while another one holds booleans)

- DataFrames are heavily used in data science projects as they allow for an **easy** and **efficient data management**
  - Enable efficient column-wise data manipulations which are especially useful for data preparation
  - Other popular libraries (e.g., scikit-learn) are generally used in combination with pandas dataFrames



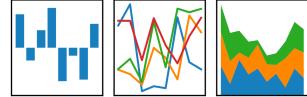
DataFrame Example:

|   | last_name | first_name | age | student? |
|---|-----------|------------|-----|----------|
| 0 | Skywalker | Luke       | 22  | True     |
| 1 | Targaryen | Daenerys   | 22  | False    |
| 2 | Weasley   | Ronald     | 13  | True     |
| 3 | White     | Walter     | 50  | False    |

# Pandas: Data Selection (with dataFrames)

- Based on some pandas *dataFrame* df, multiple possibilities exist to select a required data subset. Some popular approaches for selecting data subsets are:

pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



| Approach            | Focus                                          | Example Code                                                                                        | Example Code Explanation                                                                                           |
|---------------------|------------------------------------------------|-----------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------|
| (1) By column name  | Single column                                  | • <code>df["last_name"]</code>                                                                      | • Select column named "last_name"                                                                                  |
|                     | Multiple columns                               | • <code>df[["last_name", "age"]]</code>                                                             | • Select the columns named "last_name" and "age" (in this order)                                                   |
| (2) By position     | Single row                                     | • <code>df.iloc[1]</code>                                                                           | • Select second row                                                                                                |
|                     | Range of rows                                  | • <code>df.iloc[:3]</code><br>• <code>df.iloc[1:4]</code><br>• <code>df.iloc[[1,3,2]]</code>        | • Select row number 0 to 2<br>• Select row number 1 to 3<br>• Select row number 1, 3, and 2 (in this order)        |
|                     | Single column                                  | • <code>df.iloc[:,2]</code>                                                                         | • Select third column                                                                                              |
|                     | Range of columns                               | • <code>df.iloc[:,0:4]</code><br>• <code>df.iloc[:,1:3]</code><br>• <code>df.iloc[:,[1,0,2]]</code> | • Select column number 0 to 3<br>• Select column number 1 to 2<br>• Select column number 1, 0, and 2 in this order |
|                     | Single cell                                    | • <code>df.iloc[1,2]</code>                                                                         | • Select cell in row 1 and column 2                                                                                |
|                     | Range of cells                                 | • <code>df.iloc[1:3,[1,0,2]]</code>                                                                 | • Select row 1 to 2 and column 1, 0, and 2 (in this order)                                                         |
| (3) Condition-based | Condition: minimum age                         | • <code>df[df["age"] &gt; 15]</code>                                                                | • Select every row where age > 15                                                                                  |
|                     | Condition: specific first name                 | • <code>df[df["first_name"] == "Luke"]</code>                                                       | • Select every row where first_name equals "Luke"                                                                  |
| (4) Query method    | Query focus: age range                         | • <code>df.query('(age &lt; 30) &amp; (age &gt; 15)')</code>                                        | • Select every row where 30 > age > 15                                                                             |
|                     | Query focus: specific first name & minimum age | • <code>df.query('(first_name != "Luke") &amp; (age &gt; 15)')</code>                               | • Select every row where first_name is not "Luke" and age > 15                                                     |

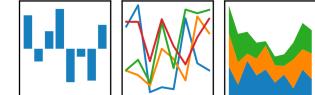
For the example code snippets, imagine you have the following DataFrame stored in a variable called df:

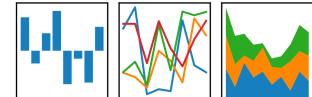
|   | last_name | first_name | age | student? |
|---|-----------|------------|-----|----------|
| 0 | Skywalker | Luke       | 22  | True     |
| 1 | Targaryen | Daenerys   | 22  | False    |
| 2 | Weasley   | Ronald     | 13  | True     |
| 3 | White     | Walter     | 50  | False    |

- Pandas and other Python libraries offer a variety of possibilities for data manipulation that can be used to perform an efficient data preparation
- As the required functionality varies from project to project, we recommend that you explore possible capabilities by your own based on some specific use case
  - See, e.g., [pandas documentation](#)
- However, to provide you with some basic insights, capabilities for handling missing data are discussed in the following slides

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$





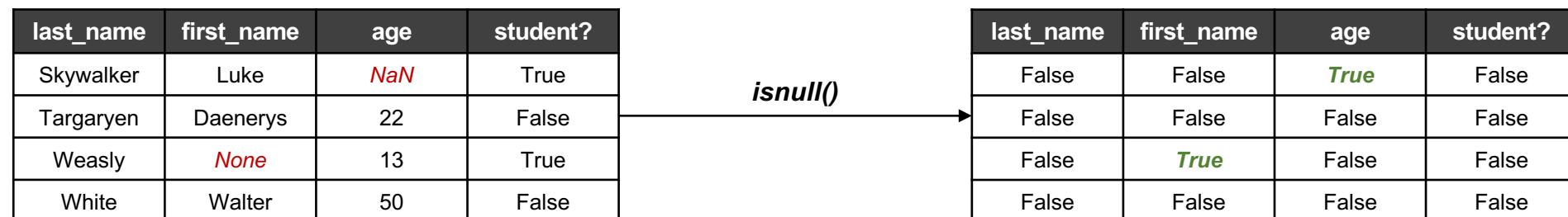
## 1) Identify missing data

- a. In Python, two value types may indicate missing values:

- **None**: official built-in Python object used to indicate a missing value
  - clearly indicates a missing value 😊
  - produces unnecessary processing overhead & must be treated as a special value 😞
- **NaN**: value type contained in the *numpy* library and is used to indicate that a variable is “Not a Number” (= NaN)
  - very efficient & can be easily treated like a number (i.e., can be used in arithmetic) 😊
  - does only represent missing numerical data & always requires the *numpy* library 😞

- b. Luckily, pandas treats both value types interchangeable & provides two helpful functions to identify them:

- **isnull()**: generates a table filled with booleans indicating missing values (= values which are either *None* or *NaN*)
- **notnull()**: generates the opposite of *isnull()*



## 2) Handle missing data

- Pandas provides two helpful functions to handle missing values (= *None* & *NaN*):
  - dropna()**: Returns a filtered version in which rows containing any missing values are discarded
  - fillna()**: Returns a copy of the provided data set in which missing values are replaced with an indicated value or based on a defined replacement method\* (i.e.: ‘backfill’, ‘bfill’, ‘pad’, ‘ffill’)

| last_name | first_name  | age        | student? |
|-----------|-------------|------------|----------|
| Skywalker | Luke        | <i>NaN</i> | True     |
| Targaryen | Daenerys    | 22         | False    |
| Weasly    | <i>None</i> | 13         | True     |
| White     | Walter      | 50         | False    |

*dropna()*

| last_name | first_name | age | student? |
|-----------|------------|-----|----------|
| Targaryen | Daenerys   | 22  | False    |
| White     | Walter     | 50  | False    |

| last_name | first_name  | age        | student? |
|-----------|-------------|------------|----------|
| Skywalker | Luke        | <i>NaN</i> | True     |
| Targaryen | Daenerys    | 22         | False    |
| Weasly    | <i>None</i> | 13         | True     |
| White     | Walter      | 50         | False    |

*fillna(value="missing")*

| last_name | first_name     | age            | student? |
|-----------|----------------|----------------|----------|
| Skywalker | Luke           | <i>missing</i> | True     |
| Targaryen | Daenerys       | 22             | False    |
| Weasly    | <i>missing</i> | 13             | True     |
| White     | Walter         | 50             | False    |

*fillna(method="ffill")*

| last_name | first_name      | age       | student? |
|-----------|-----------------|-----------|----------|
| Skywalker | Luke            | 22        | True     |
| Targaryen | Daenerys        | 22        | False    |
| Weasly    | <i>Daenerys</i> | 13        | True     |
| White     | Walter          | <i>13</i> | False    |

## 1

### Introduction to Python

1.1 Overview

1.2 Data Import & Management

1.3 Data Visualization & Exploration

1.4 Development Environments

1.5 Resources & Hands-On Exercises

- Python offers a big variety of visualization libraries
- Popular libraries are (inter alia): [matplotlib](#) & [plotly](#)
- But many more exist that offer varying foci!
  - Please feel free to explore any other libraries to identify the ones that fit your needs ☺
- In the following slides, some basic visualizations are discussed to provide you some insights on how to use the libraries for visualizing your data
  - **However:** This is just the tiny tip of the iceberg of visualizations you can generate! (Again: It's worth to explore the libraries!)

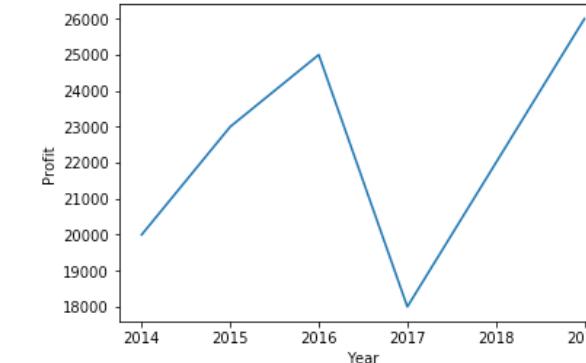
## a) Use matplotlib *explicitly* to plot data

```
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt

import data from a csv file into a pandas DataFrame:
data = pd.read_csv('profitPerYear.csv', delimiter=',')

plot the two dimensions year and profit and name the x and y axis accordingly:
plt.plot(data['year'], data['profit'])
plt.xlabel('Year')
plt.ylabel('Profit')
```

output



Example  
dataFrame:

| year | profit |
|------|--------|
| 2014 | 20000  |
| 2015 | 23000  |
| 2016 | 25000  |
| 2017 | 18000  |
| 2018 | 22000  |
| 2019 | 26000  |

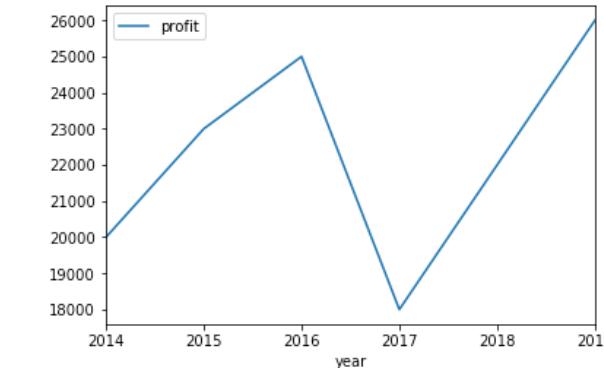
## a) Use matplotlib *implicitly* to plot data (with a pandas DataFrame)

```
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt

import data from a csv file into a pandas DataFrame:
data = pd.read_csv('profitPerYear.csv', delimiter=',')

plot the dataframe based on the indicated values for x and y axis:
data.plot(x='year', y='profit')
```

output



# Create a Bar Chart with *matplotlib*

## a) Use matplotlib *explicitly* to create a bar chart

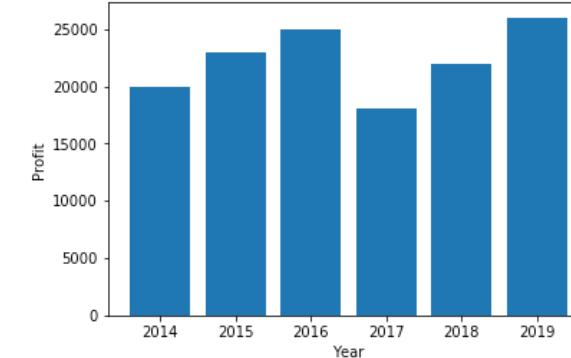
- Simply replace the function `plot()` with `bar()`

```
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt

import data from csv into a pandas DataFrame:
data = pd.read_csv('profitPerYear.csv', delimiter=',')

plot the two dimensions year and profit and name the x and y axis accordingly:
plt.bar(data['year'], data['profit'])
plt.xlabel('Year')
plt.ylabel('Profit')
```

output



Example  
dataFrame:

| year | profit |
|------|--------|
| 2014 | 20000  |
| 2015 | 23000  |
| 2016 | 25000  |
| 2017 | 18000  |
| 2018 | 22000  |
| 2019 | 26000  |

## b) Use matplotlib *implicitly* to create a bar chart (with a pandas DataFrame)

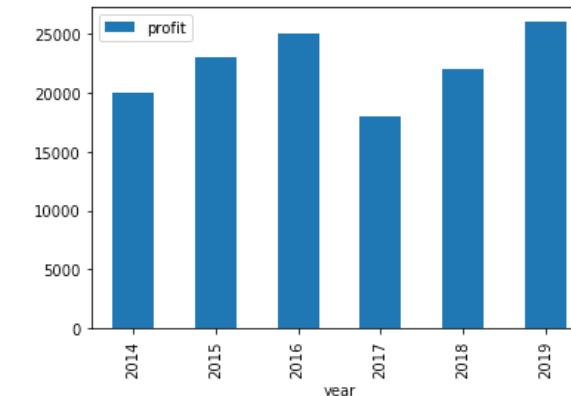
- Simply indicate 'bar' as chart type with the additionally introduced `kind` parameter

```
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt

import data from csv into a pandas DataFrame:
data = pd.read_csv('profitPerYear.csv', delimiter=',')

plot the dataframe based on the indicated values for x and y axis:
data.plot(x='year', y='profit', kind='bar')
```

output



## 1

### Introduction to Python

1.1 Overview

1.2 Data Import & Management

1.3 Data Visualization & Exploration

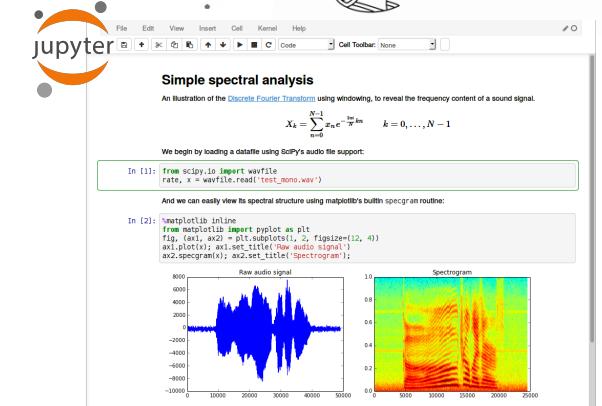
1.4 Development Environments

1.5 Resources & Hands-On Exercises

# Overview of Python Development Environments

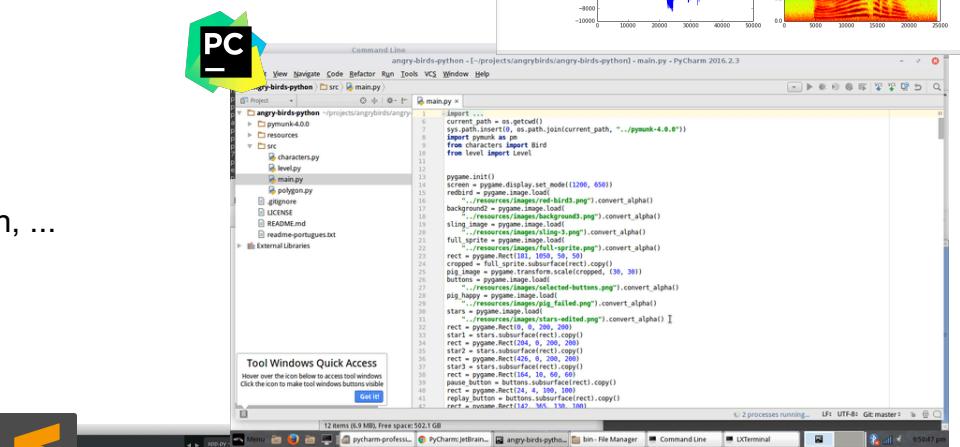
## a) Interactive Notebooks: Jupyter Notebooks

- Document-like coding environment: allows to structure & execute your code along text snippets
- Allows a straightforward documented script-like programming
- Currently very popular in the data science community as it allows an easy documented exploration of data & AI algorithms! (but not really suited to implement huge applications)



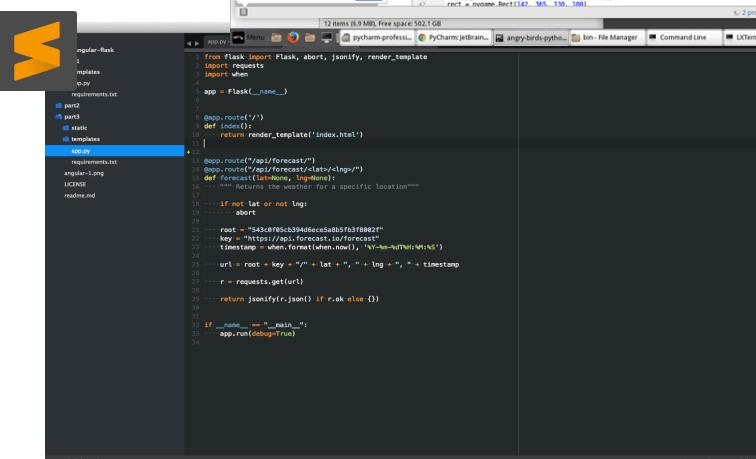
## b) Python-specific IDE: PyCharm

- Integrated development environment (IDE) for Python developed by JetBrains
- Includes: code analysis, a graphical debugger, an unit tester, version control integration, ...
- A commercial and a (free) community version exists
- Especially useful for managing big complex programming projects
- **Alternatives:** Spyder, Thonny, ...



## c) General Code/Text Editor: Sublime Text

- Lightweight editors that provide helpful capabilities for text editing, such as, e.g., text highlighting based on the chosen programming language
- Especially useful for quickly creating scripts
- **Alternatives:** Atom, Visual Studio Code, ...





- Open-source-distribution of Python (commercial version exists as well)
- Very popular among data scientists!
  - As it comes with many tools and libraries used in data science, including:
    - A package manager (called “conda”) for installing and maintaining Python libraries
    - Environment manager for managing different installations of Python (e.g., to switch between a Python 2 and 3 installation)
    - Development tools like Jupyter Notebook
    - ...and much more!
- See: <https://www.anaconda.com/>

## 1

### Introduction to Python

- 1.1 Overview
- 1.2 Data Import & Management
- 1.3 Data Visualization & Exploration
- 1.4 Development Environments
- 1.5 Resources & Hands-On Exercises

- **Official Documentation:**

- [Python 3](#)
- [Python 2](#)

- **Free eBook:**

- [Swaroop C H \(2013\): „A Byte of Python“](#)  
*(PDF version: scroll „introduction“ page down for a download link)*

- **Free Online Courses:**

- [Codecademy's Python 2 & 3 courses](#)
- [Google's Python Class](#)
- [Tutsplus Python Tutorial](#)  
*(combines many further resources)*

- Install the most recent version of the Anaconda Python distribution (Python 2 or 3 version)
  - URL: <https://www.anaconda.com/distribution/#download-section>
- Explore the data management and exploration capabilities of Python by completing the tasks of the exercise 2 sheet