

Big Data analysis with Python



State of the art and overview of the applications in the area of the diploma thesis

Supervisor: Dr hab. inż. Dariusz Król
Presented by: Oleksii Kyrylchuk 223224
13.10.2017



HR EXCELLENCE IN RESEARCH

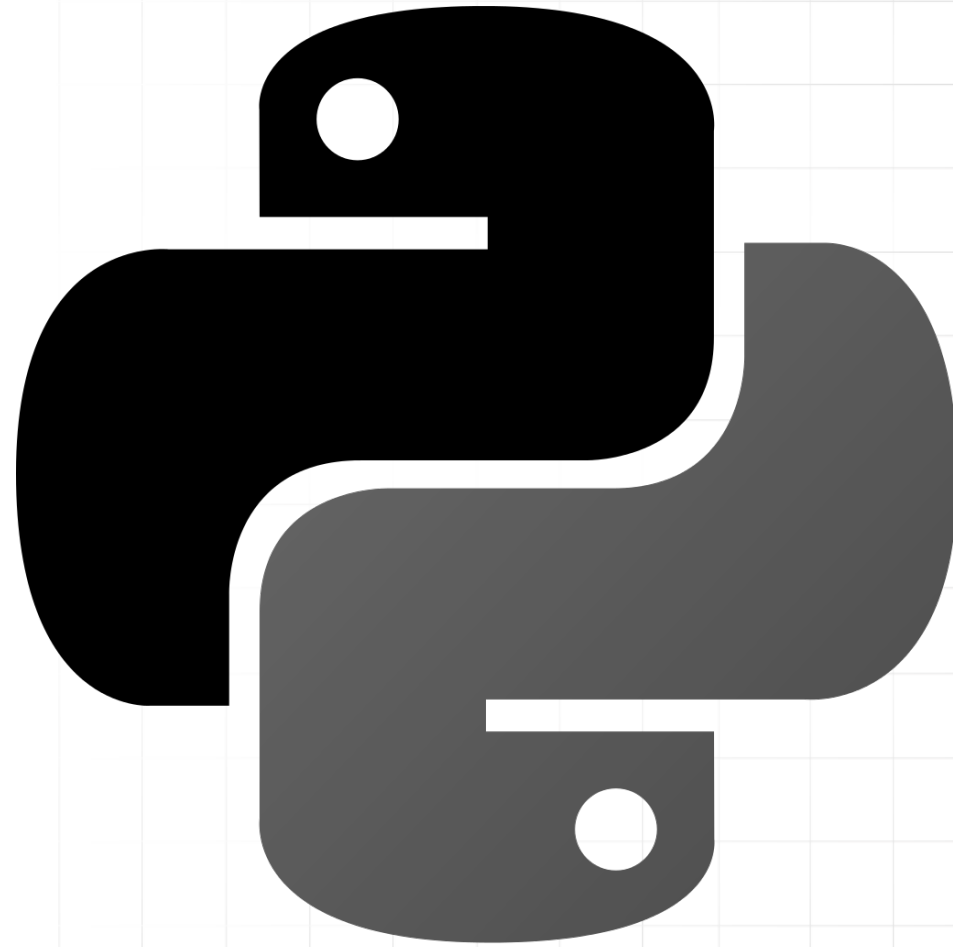


Wrocław University
of Science and Technology

My topic

“Big Data analysis with
Python”

“Przetwarzanie dużych
zbiorów danych w
środowisku Python”



Why did I choose this topic?



HR EXCELLENCE IN RESEARCH



Wrocław University
of Science and Technology

Why big data?

- Analyze stock data
 - Predict the future stock price
 - Plot the stock data over time
 - ...
- Analyze facebook data
 - Plot the event location data
 - Analyze the popularity of posts
 - ...
- Analyze YouTube data
 - Make a computer generated song/video

Why Python?

- Quickly achieve good results
- Wide variety of tools available
- Very popular language in data science

Brief overview of the DS stack

Extract
Gather

- Web scraping, API scraping, Data source packages, Flat files, ...

Transform
Clean

- NumPy, pandas, built-in lists, ...

Load
Store

- APIs - sqllite3, pycopg2, hadoop, ORMs - SQLAlchemy, peewee ...

Analyze

- Tensorflow, matplotlib, scikit-learn, bokeh, pandas, ...

Step 1

Data extraction/gathering



HR EXCELLENCE IN RESEARCH



Wrocław University
of Science and Technology

Processing files

- CSV reader¹ – A built-in package
- Numpy² – A more advanced package, oriented on scientific computing
- Pandas³ – A package, that implements R's DataFrame, and has a lot of extra features, including plotting

¹ <https://docs.python.org/3/library/csv.html>

² <http://www.numpy.org/>

³ <http://pandas.pydata.org/>

Data source packages

- There are packages, that allow you to download stock data automatically without having to scrape websites/APIs
- yahoo-finance¹
- googlefinance²
- pandas-datareader³ - fetches data from multiple sources

- Example:

```
import pandas_datareader as pdr  
pdr.get_data_yahoo('AAPL') # Returns a DataFrame
```

¹ <https://pypi.python.org/pypi/yahoo-finance>

² <https://pypi.python.org/pypi/googlefinance>

³ <https://pandas-datareader.readthedocs.io>

API scraping

- `json`¹ – a built-in package that allows us to convert json into a python dictionary
- `urllib3`² – a built-in http request package (synchronous)
- `requests`³ – an external request package, that simplifies the code
- `aiohttp`⁴ – a package, that allows you to asynchronously send http requests, which greatly improves the performance (asyncio/await syntax accessible in python 3.5+)

¹ <https://docs.python.org/3/library/json.html>

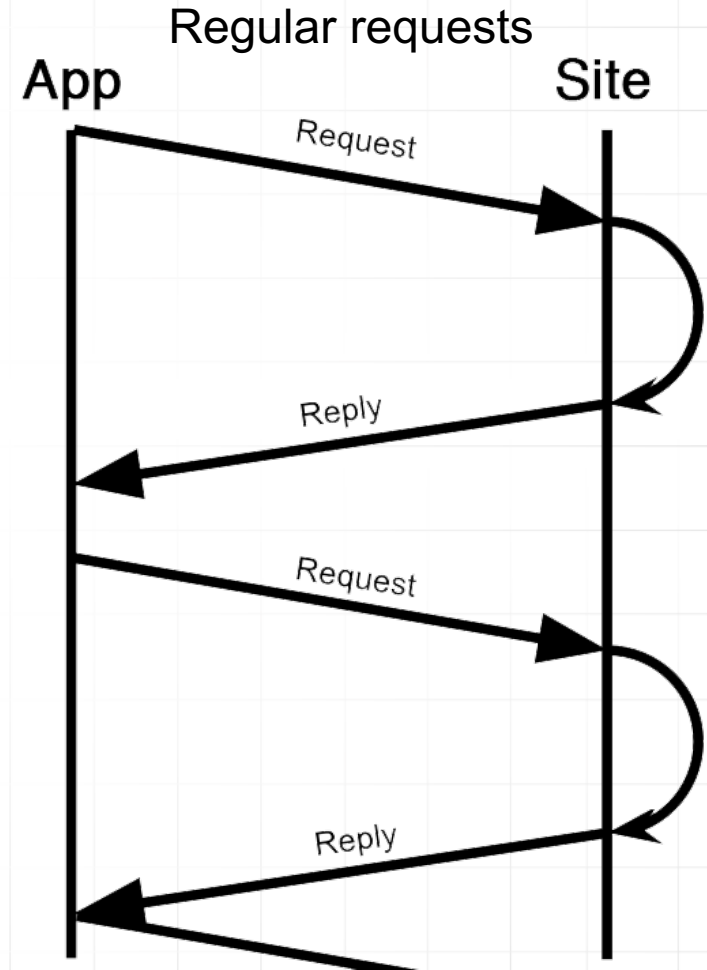
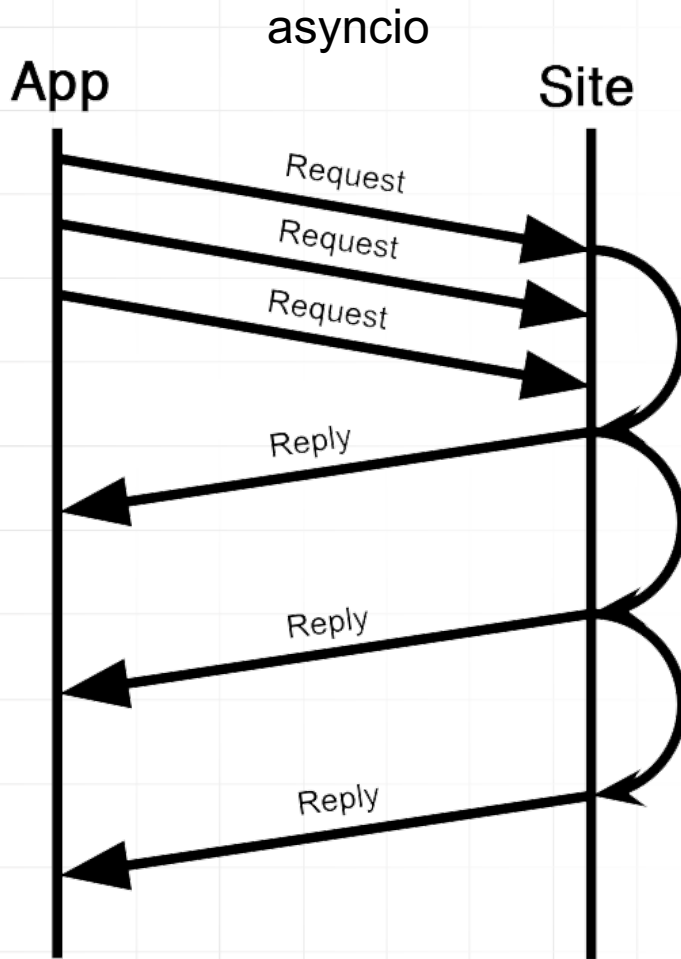
² <https://urllib3.readthedocs.io/en/latest/>

³ <http://docs.python-requests.org/en/master/>

⁴ <http://aiohttp.readthedocs.io/en/stable/>



API/web scraping - asyncio



Web scraping

- We use the same packages from the API scraping to download the data
- `re`¹ – a built-in package that implements regular expressions
- `BeautifulSoup`² – a package that allows us to navigate HTML much more easily
- `scrapy`³ – a package for making web spiders, that download webpages recursively

¹ <https://docs.python.org/3/library/re.html>

² <https://www.crummy.com/software/BeautifulSoup/>

³ <https://scrapy.org/>

Step 2

Data transformation, cleaning and storage



HR EXCELLENCE IN RESEARCH



Wrocław University
of Science and Technology

Transformation & cleaning

- Built-in lists¹
- numpy² – a lower level package that gives you access to array and matrix classes with a lot of useful methods
- pandas³ – one of the best packages for dealing with table data. Has an interface to import/export data to/from numpy, csv, zip files, etc.

¹ <https://docs.python.org/3.6/tutorial/datastructures.html>

² <http://www.numpy.org/>

³ <http://pandas.pydata.org/>



Data Storage – APIs vs ORMs

- APIs (drivers)
 - sqlite3¹
 - psycopg2²
- Use the database API to do everything
- Use SQL for querying
- Harder to code & maintain
- Usually faster
- ORMs
 - SQLAlchemy³
 - peewee⁴
- Easier to code & maintain
- They use the API packages to send the SQL to the database
- Sometimes slower

¹ <https://docs.python.org/3/library/sqlite3.html>

² <http://initd.org/psycopg/>

³ <https://www.sqlalchemy.org/>

⁴ <http://docs.peewee-orm.com/en/latest/>



Step 3

Data analysis

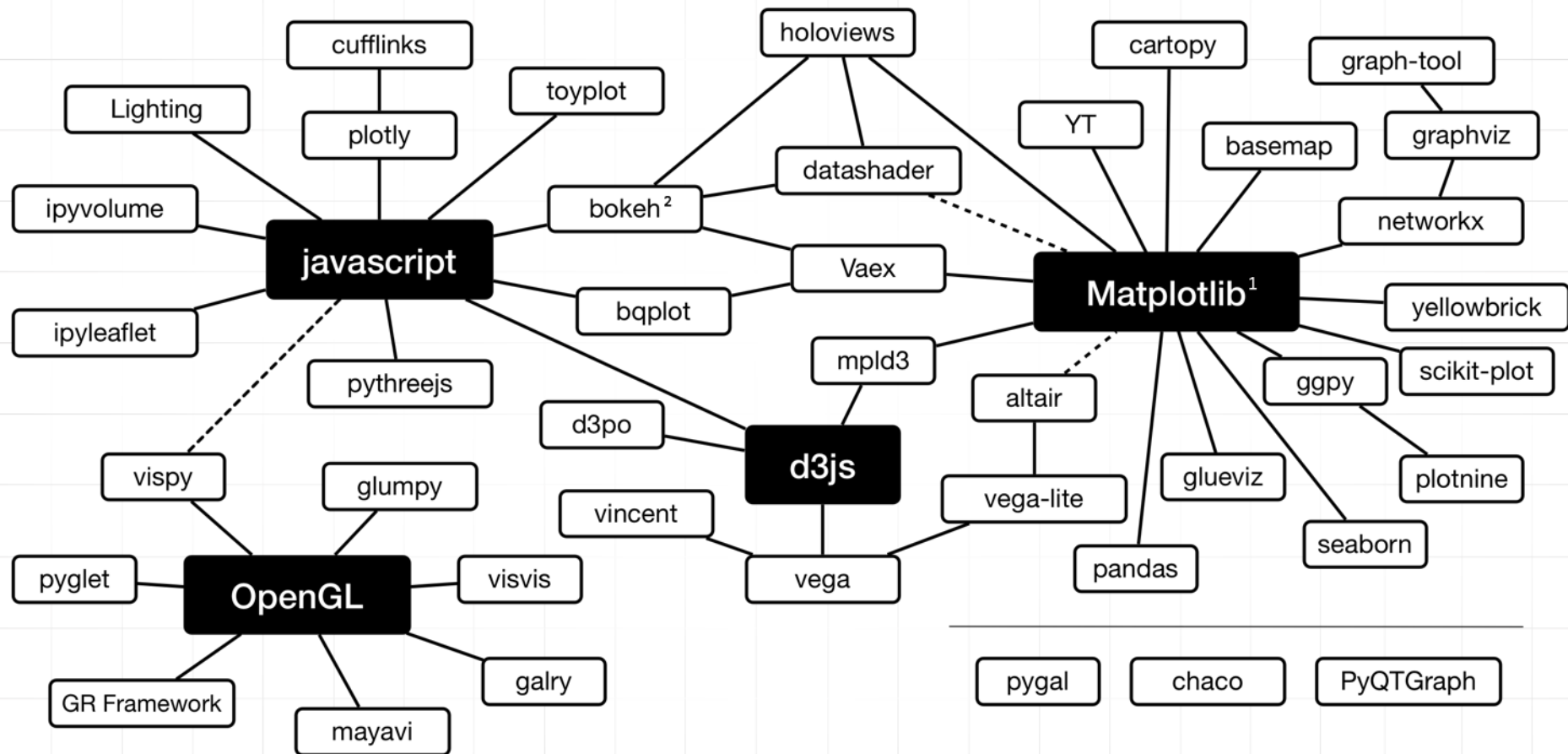


HR EXCELLENCE IN RESEARCH



Wrocław University
of Science and Technology

Data analysis – Visualization



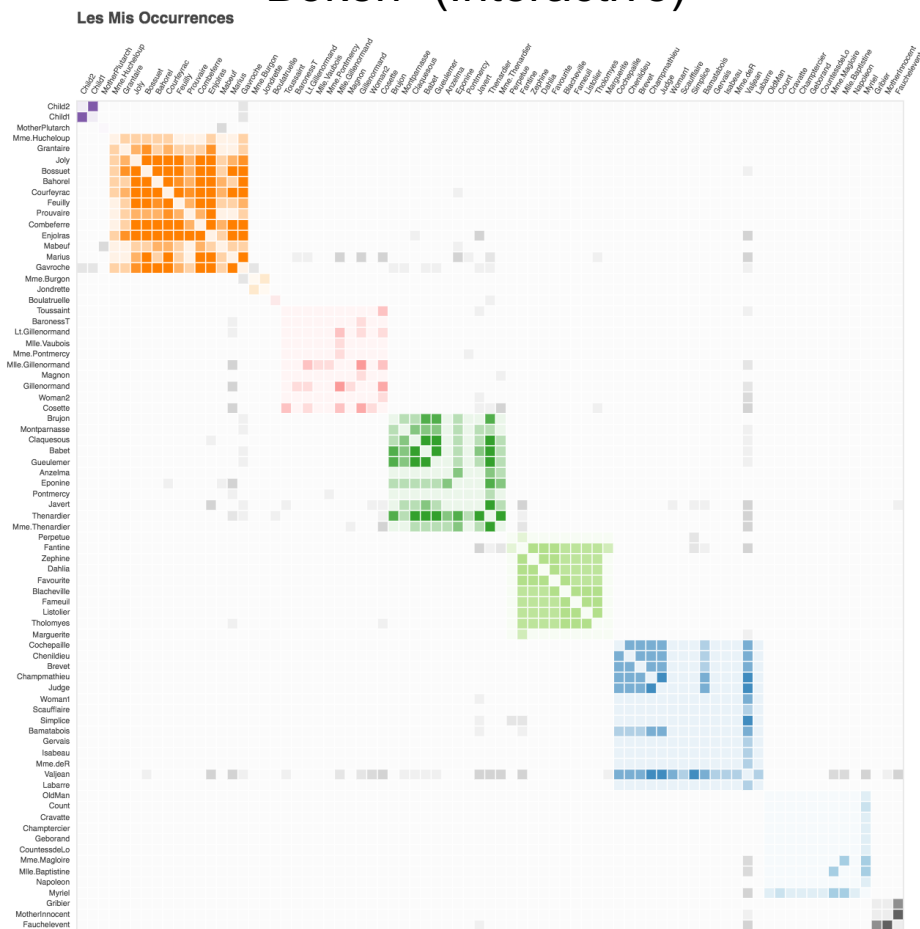
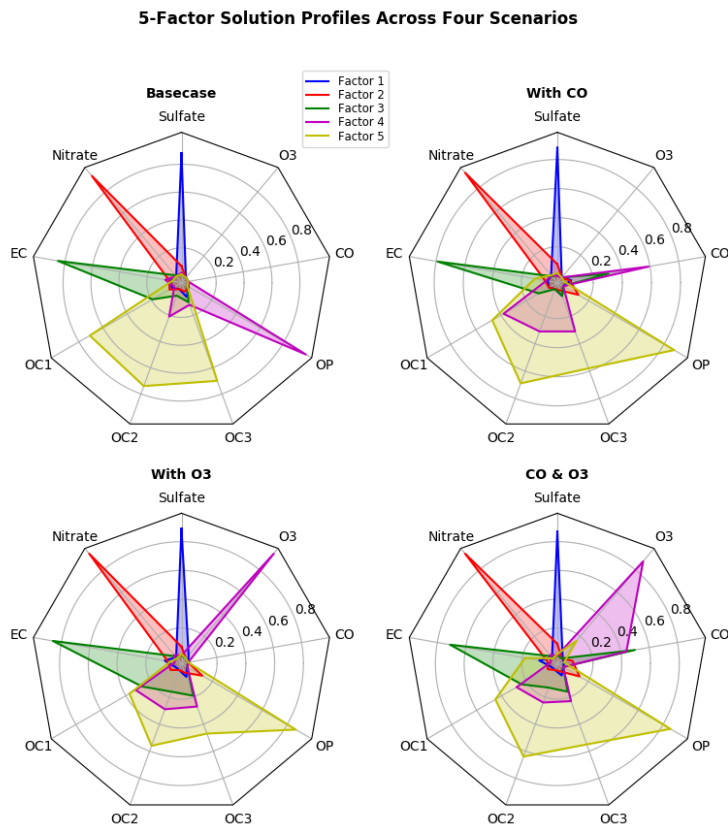
¹ <https://matplotlib.org/>

² <https://bokeh.pydata.org/en/latest/>

Data analysis – Example visualizations

Matplotlib¹

Bokeh² (Interactive)



Data analysis – Machine learning

- Scikit-learn¹ – the easiest way to get the results
- Numpy², Scipy³ – if we want to do everything from scratch
- Tensorflow⁴ – a more efficient way to do everything from scratch
- Theano⁵ – a competitor to Tensorflow
- Keras⁶ – A high level package that can use either Tensorflow or Theano as its backend

¹ <http://scikit-learn.org>

² <http://www.numpy.org/>

³ <https://www.scipy.org>

⁴ <https://www.tensorflow.org>

⁵ <http://deeplearning.net/software/theano>

⁶ <https://keras.io/>



Similar applications

- There are a lot of apps that implement a similar structure, but most of them are corporate and operate on corporate datasets
- Example projects include
 - Analyzing stock data¹
 - Analyzing medical data
 - Computer-generated art/art style^{2 3}
 - Computer-generated music
 - ...

¹ <http://sentdex.com/>

² <https://deepart.io/>

³ <http://nightmare.mit.edu/>



My thesis plan – A Roadmap

Task	Status	Month	March		April				May					June		J - S	October					November				Dec	
		Week	13	14	15	16	17	18	19	20	21	22	23	24	25-27	28-40	41	42	43	44	45	46	47	48	49	50	
Organization																											
Paperwork (Registering, etc.)	Completed														Session	Internship											
Specification	Completed																										
Studying available tools	In progress																										
Development																											
Implementing the facebook analyzer	In progress														Session	Internship											
Implementing the stock analyzer	In progress																										
Implementing the youtube analyzer	Not started																										
Writing																											
Introduction	In progress														Session	Internship											
Available technologies	Not started																										
Part 1 - Stock analysis	Not started																										
Part 2 - Facebook analysis	Not started																										
Part 3 - Youtube analysis	Not started																										

References

- Wes McKinney (2012) Python for Data Analysis – O'Reilly Media, Inc.
- Joel Grus (2015) Data Science from Scratch – O'Reilly Media, Inc.
- Jake VanderPlas (2017) The Python Visualization Landscape - PyCon 2017
<https://www.youtube.com/watch?v=FytuB8nFHPQ>
- Nikola Milosevic (2016) Equity forecast: Predicting long term stock price movement using machine learning