AMIS

# Data Analytics on Conference Session Catalog

## using Jupyter Notebooks

## handson workshop
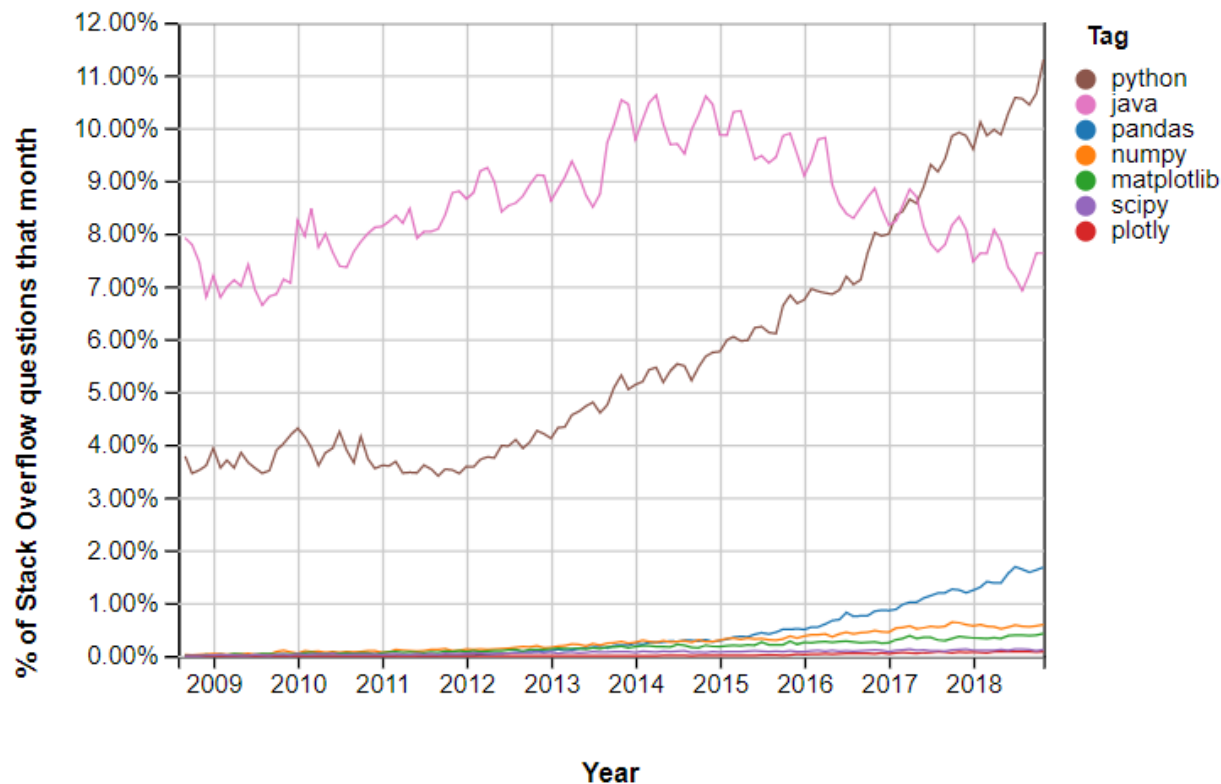
Lucas Jellema
Conclusion Machine Learning Gilde
21 februari 2019

# Hands On met Jupyter Notebooks in vier stappen

- Een werkende Jupyte Notebook Server omgeving
  - Lokaal – op basis van Docker container
  - Cloud – in een KataCoda omgeving
- Hello World Notebook
  - Eerste stappen met Notebook, Markdown, Python & Pandas
- Casus Oracle OpenWorld 2018 Session Catalog
  - Gather
  - Wrangle
  - Analyze
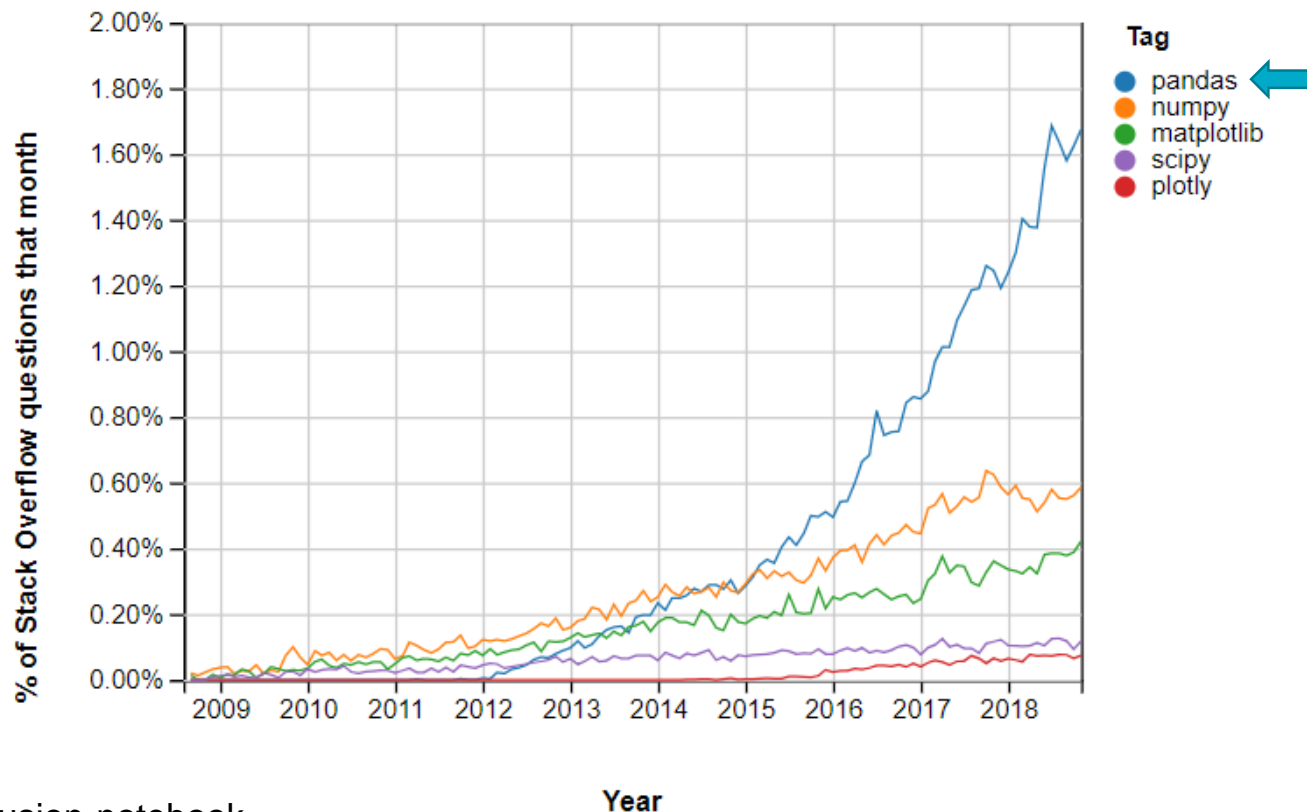  - Visualize
- Doe Het Zelf Notebook met Titanic Data Set

Resources: bit.ly/conclusion-notebook

# Groeiende belangstelling in Python…



Resources: bit.ly/conclusion-notebook

# .. en in een specifieke Python Library in het bijzonder
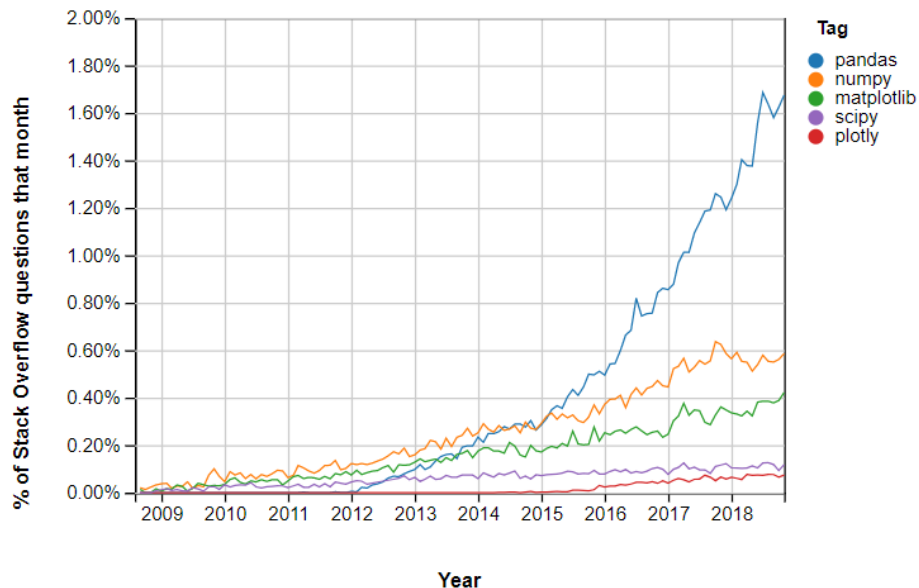


Resources: bit.ly/conclusion-notebook

# pandas (software)

From Wikipedia, the free encyclopedia

In computer programming, **pandas** is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series. It is free software released under the three-clause BSD license.[2] The name is derived from the term "panel data", an econometrics term for data sets that include observations over multiple time periods for the same individuals.[3]

**pandas**

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

| | |
|---|---|
| **Original author(s)** | Wes McKinney |
| **Developer(s)** | Community |
| **Initial release** | 11 January 2008; 11 years ago |
| **Stable release** | 0.23.4[1] / 3 August 2018; 5 months ago |
| **Repository** | github.com/pandas-dev/pandas |
| **Written in** | Python, Cython, C |
| **Operating system** | Cross-platform |
| **Type** | Technical computing |
| **License** | New BSD License |
| **Website** | pandas.pydata.org |

% of Stack Overflow questions that month vs Year

Tag: pandas, numpy, matplotlib, scipy, plotly

**QuantDare**

ARTIFICIAL INTELLIGENCE    ASSET MANAGEMENT    RISK MANAGEMENT    PYTHON    R    ALL     ABOUT US    TERMS OF USE & PRIVACY POLICY

**Daring** to quantify the markets |

The scientific blog of **ETS Asset Management Factory**

PYTHON

# Calculate monthly returns...with Pandas

*mgreco*    *27/09/2017*

💬 2

Calculating returns on a price series is one of the most basic calculations in finance, but it can become a headache when we want to do aggregations for weeks, months, years, etc. In Python, the Pandas library makes this aggregation very easy to do, but if we don't pay attention we could still make mistakes. Assuming that we want the

https://pandas.pydata.org

pandas

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

Fork me on GitHub

home // about // get pandas // documentation // community // talks // donate

# Python Data Analysis Library

*pandas* is an open source, BSD-licensed library providing high-performance, easy-to-use data structures and data analysis tools for the Python programming language.

*pandas* is a NumFOCUS sponsored project. This will help ensure the success of development of *pandas* as a world-class open-source project, and makes it possible to donate to the project.

A Fiscally Sponsored Project of

NUMF⊙CUS

OPEN CODE = BETTER SCIENCE

## v0.23.4 Final (August 3, 2018)

## VERSIONS

**Release**
0.24.1 - February 2019
download // docs // pdf

**Development**
0.25.0 - April 2019
github // docs

**Previous Releases**
0.24.0 - download // docs // pdf
0.23.4 - download // docs // pdf
0.23.3 - download // docs // pdf
0.23.2 - download // docs // pdf
0.23.1 - download // docs // pdf
0.23.0 - download // docs // pdf
0.22.0 - download // docs // pdf
0.21.1 - download // docs // pdf

# Pandas = Panel Data

Anatomy of a DataFrame

[ Personal_Data , Sales_Data , Region ]

There are 3 sheets in the workbook

In [21]: `#Display the records in the first sheet`
`sheet_3`

Out[21]:

| | SALES_ID | Sales_BY_Region | January | February | March | April | May | June | July | August | September | October | November | December |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | LA | 2000 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 1 | 2 | NY | 2200 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 2 | 3 | HS | 2400 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 3 | 4 | HS | 2100 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 4 | 5 | HS | 2300 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 5 | 6 | LA | 3200 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 6 | 7 | LA | 2210 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 7 | 8 | LA | 2320 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 8 | 9 | HS | 1945 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 9 | 10 | HS | 900 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 10 | 11 | LA | 1920 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 11 | 12 | NY | 1800 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 12 | 13 | HS | 1820 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |
| 13 | 14 | LA | 1450 | 2200 | 2420 | 2662 | 2928 | 3221 | 3543 | 3897 | 4287 | 4716 | 5187 | 5706 |

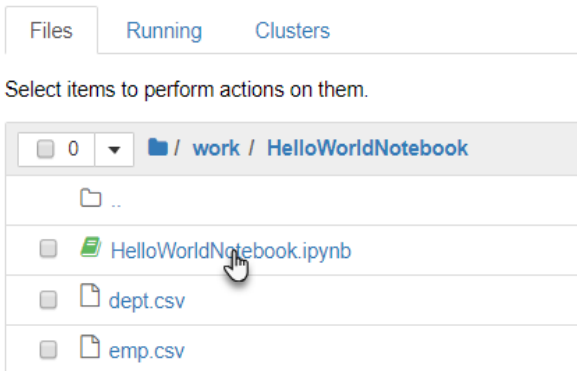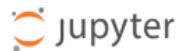# Een paar operaties op een Pandas Data Frame

- *df.dtypes* – alle data types van alle kolommen

- *df['column_name'].value_counts* – telling van aantal verschillende waarden

- *df.head(5)* – toon de eerste vijf rijen van het data frame
  - df['column_name'].head(5) (of .tail(10))
  - df[['column_name', 'column_name_2']].head(5)

- *pd.crosstab(df['column_name'], df['column_name_2'])* – kruistabel

- *df['column_name3'] = 2 \* df['column_name']* – voeg een kolom toe
  - df['column_name3'] = df['column_name'].apply(len) – bepaal waarde van nieuwe kolom door toepassen van een functie op een bestaande kolom
  - df['column_name3'] = df.apply(lambda row: 'Y' if row['column_name'] == 'HOT' else 'N') – bepaal waarde van nieuwe kolom op basis van conditie

- *df.drop('column_name3')* – verwijder kolom

# Hello World Notebook

- Als je omgeving draait
  - Lokaal of Katacoda
- Open dan het Notebook
  HelloWorldNotebook.ipynb
  in folder /work/HelloWorldNotebook
- Stap door de cellen
  lees de instructie
  voer de opdrachten uit



```python
import pandas as pd
# read csv file into Pandas Data Frame, using a semi colon is separator
hr= pd.read_csv("emp.csv",sep=';')
#show first five rows in the dataframe
hr.head(5)
```

Out[1]:

|   | empno | ename | job | mgr | hiredate | sal | comm | deptno |
|---|-------|-------|-----|-----|----------|-----|------|--------|
| 0 | 7369 | SMITH | CLERK | 7902.0 | 13/06/1993 | 800.0 | 0.0 | 20 |
| 1 | 7499 | ALLEN | SALESMAN | 7698.0 | 15/08/1998 | 1600.0 | 300.0 | 30 |
| 2 | 7521 | WARD | SALESMAN | 7698.0 | 26/03/1996 | 1250.0 | 500.0 | 30 |
| 3 | 7566 | JONES | MANAGER | 7839.0 | 31/10/1995 | 2975.0 | NaN | 20 |
| 4 | 7698 | BLAKE | MANAGER | 7839.0 | 11/06/1992 | 2850.0 | NaN | 30 |

# Press shift + tab + tab to get Help in a Jupyter Notebook



```
In [ ]: df.stack

In [ ]:      Signature: df.stack(level=-1, dropna=True)
             Docstring:
In [ ]:      Pivot a level of the (possibly hierarchical) column labels, returning a
             DataFrame (or Series in the case of an object with a single level of
             column labels) having a hierarchical
In [ ]:      of row labels.
             The level involved will automatically

In [ ]:      Parameters
```

Pressing shift + tab + tab to reveal the docum

You can also press `tab` directly following a dot to have a dropdown menu

```
In [ ]: df.
        df.abs
In [ ]: df.add
        df.add_prefix
        df.add_suffix
In [ ]: df.agg
        df.aggregate
In [ ]: df.align
        df.all
In [ ]: df.any
        df.append
```

Pressing tab following a DataFrame lists the 200+ available objects

# Data Analytics on Oracle OpenWorld 2018 Session Details

# High level overview

- Two events: Oracle OpenWorld 2018 and Oracle CodeOne 2018
- Over 2000 sessions
- Over 2500 speakers
- Various dimensions:

Filters                    CLEAR

▶ Intelligent Cloud Applications

▶ Oracle Cloud Platform

▶ Oracle Cloud Infrastructure

▶ Your Cloud Transformation
   Roadmap

▶ Your Cloud Success: Training
   and End to End Support

▶ The "Suite Spot": Connected
   and Intelligent Business

▶ Accelerate Growth: Solutions
   for Small to Medium (SMB)
Business

▶ Real Stories, Real Customers

▶ Sessions By Topic

▶ Sessions By Role

▶ Sessions By Industry

▶ Session Type

▶ Day

# High level overview

- Details per session:

# High level overview

- Details per speaker:

# Flow for Oracle OpenWorld 2018 Session Data Analytics

# Data flow for Oracle OpenWorld 2018 Session Data

Speakers:
Read JSON file, extract speaker details, visualize and analyze speaker data

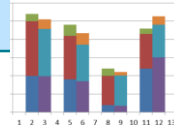Fetch all raw session data per Event and per Session Type from API and write to 44 local JSON files

Read all JSON files, discard unneeded attributes, derive new attributes, deduplicate records and write to a single local JSON file

Gender guesser

API for session details in JSON

{ api }

Jupyter Notebook & Python

JSON files

Jupyter Notebook & Python

JSON file

Jupyter Notebook & Python

Jupyter Notebook & Python

Report & Visualization on Speakers

Report & Visualization on Sessions

oow2018-sessions-wrangled.json

/datalake

/datawarehouse

Sessions:
Read JSON file, visualize and analuyze session data