

KÖNIGSWEG



Analytics with Pandas and Jupyterlab
ALEXANDER CS HENDORF FOR PYLADIES HAMBURG



KÖNIGSWEG

We do digital excellence.

STRATEGY & INNOVATION

DATA & ARTIFICIAL INTELLIGENCE

BUSINESS TRANSFORMATION &
OPERATIONS

Get in touch with our specialists.

KÖNIGSWEG



LIKE THE PICTURES IN THE SLIDES?
COME IN THE BREAK AND MAKE ONE OF YOUR OWN!



Foto ändern



Mitglied von **PyData** - 142 Gruppen ?

PyData Südwest

📍 Mannheim, Deutschland

👤 737 Mitglieder · Öffentliche Gruppe ?

👤 Organisiert von **Alexander C. S. Hendorf** und 5 andere



Teilen:

[Über uns](#)

[Events](#)

[Mitglieder](#)

[Fotos](#)

[Diskussionen](#)

[Mehr](#)

Gruppe verwalten ▾

Event erstellen ▾

Worum es bei uns geht

Welcome to PyData Südwest!

...

[Mehr lesen](#)

Bevorstehende Events (1)

[Alle anzeigen](#)

Organisatoren



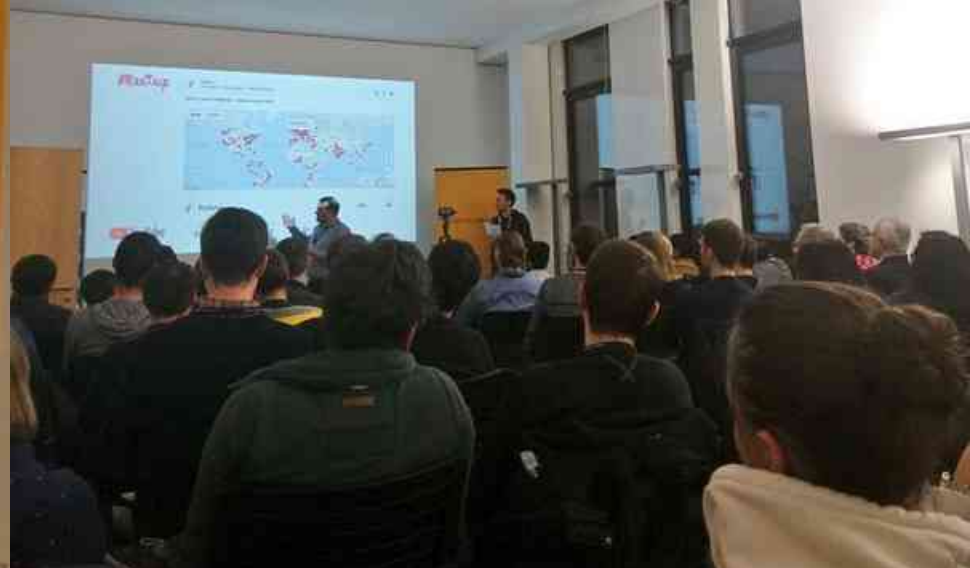
Alexander C. S. Hendorf und 5 andere

[Nachricht](#)

Mitglieder (737)

[Alle anzeigen](#)







Joint conference
October 14 to 16
at bcc, Berlin

Ticket sale opens April 6
CfP is open until May 5
Fin Aid is open until June 1

pycon.de. @pyconDE
berlin.pydata.org @pydataBER

Alexander C. S. Hendorf

- Managing Partner & Principal Consultant Information Technology
Consulting on AI & Data Science
- Python Software Foundation Fellow,
Program Chair EuroPython, EuroSciPy, PyConDE & PyData, MongoDB Master
Speaker Europa & USA MongoDB World New York / San José, PyCon Italy, CEBIT
Developer World, BI Forum, IT-Tage FFM, PyData, PyParis, PyData Südwest &
Frankfurt



ah@koenigsweg.com

 @hendorf



Prepared?

- Follow instructions: git <http://bit.ly/pandas-base1>
- Make sure you have the latest version: *git pull*



Introduction to Data Analytics with Pandas and Jupyterlab

Jupyter

- Ecosystem
- Benefits of Jupyter
- Jupyter Notebooks
- Jupyterlab

Pandas

- Benefits of Pandas
- How to work with Pandas
- Visualisation

iPython

IPython 5.3.0 -- An enhanced Interactive Python.

? -> Introduction and overview of IPython's features.

%quickref -> Quick reference.

help -> Python's own help system.

object? -> Details about 'object', use 'object??' for extra details.

```
[In [1]: n = 100000
```

```
[In [2]: import numpy as np
```

```
[In [3]: %timeit np.sum(1. / np.arange(1., n) ** 2)
```

The slowest run took 8.35 times longer than the fastest. This could mean that an intermediate result is being cached.

1000 loops, best of 3: 186 µs per loop

```
[In [4]: np.arange
```

np.arange	np.arcsin	np.arctan2	np.argmin	np.argwhere	np.array2string	np.array_repr
np.arccos	np.arcsinh	np.arctanh	np.argmax	np.argsort	np.array_equal	np.array_split
np.arccosh	np.arctan	np.argmax	np.argsort	np.array	np.array_equiv	np.array_str

Jupyter Notebooks

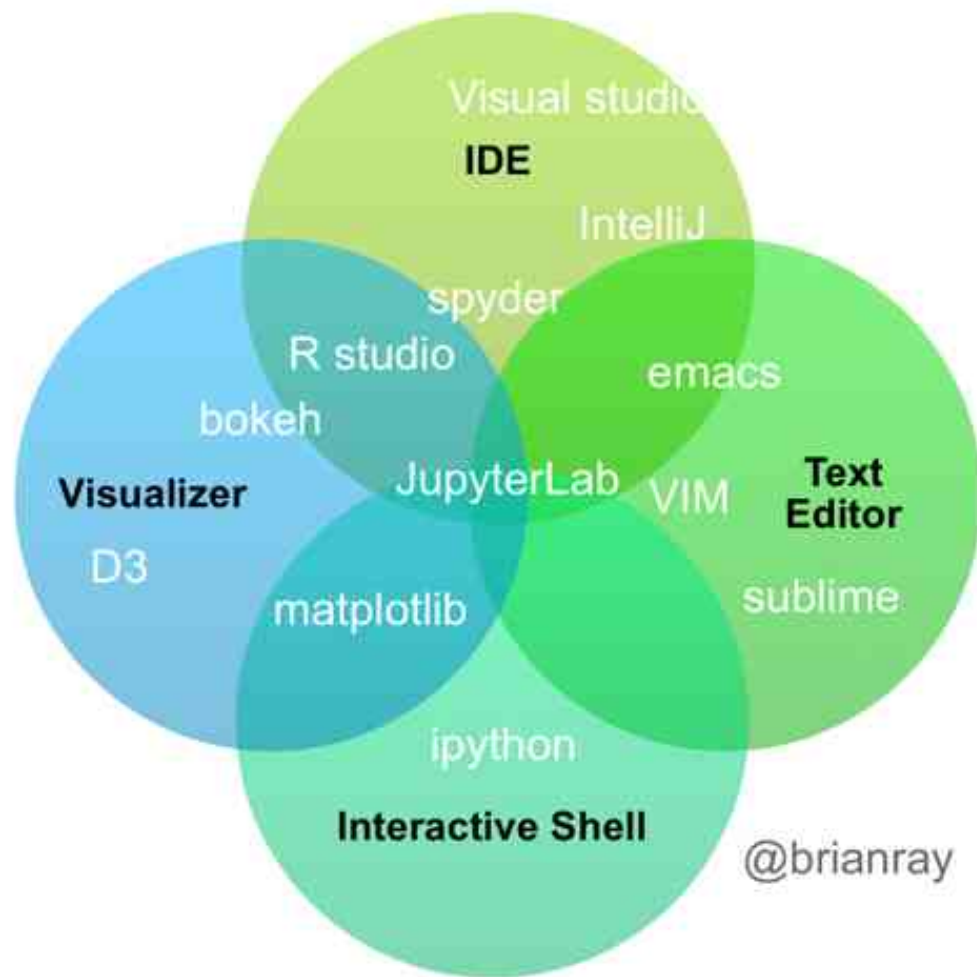
- Code, text (docs, background, research, references,...) und Visualisation
- Full programme / script
- IDE
- Explore iteratively
- Reproducible and customizable
- Export to other formats(HTML, PDF, Python module,...)



Jupyterlab

- Was in the works since 2016
- Released in January 2018
- The *future* Jupyter Notebooks
- Compact D1E
- Extendable





Hands on Jupyterlab



File Edit View Run Kernel Tabs Settings Help

Files

- notebooks
- Data.ipynb
- Fasta.ipynb
- Julia.ipynb

Running

Lorenz.ipynb

Code

Python 3

In this Notebook we explore the Lorenz system of differential equations:

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= \rho x - y - xz \\ \dot{z} &= -\beta z + xy\end{aligned}$$

let's call the function once to view the solutions. For this set of parameters, we see the trajectories swirling around two points, called attractors.

```
from lorenz import solve_lorenz
x_t = solve_lorenz(N=10)
```

lorezn.py

```
def solve_lorenz(N=10, max_time=4.0, sigma=10.0, beta=8./3, rho=28.0):
    """Plot a solution to the Lorenz differential equations."""
    fig = plt.figure()
    ax = fig.add_axes([0, 0, 1, 1], projection='3d')
    ax.axis('off')

    # prepare the axes limits
    ax.set_xlim([-25, 25])
    ax.set_ylim([-35, 35])
    ax.set_zlim([5, 55])

    def lorenz_deriv(x_y_z, t0, sigma=sigma, beta=beta, rho=rho):
        """Compute the time-derivative of a Lorenz system."""
        x, y, z = x_y_z
        return [sigma * (y - x), x * (rho - z) - y, x * y - beta * z]

    # Choose random starting points, uniformly distributed from -15 to 15
    np.random.seed(1)
    x0 = -15 + 30 * np.random.random((N, 3))
```

10.00

2.67

28.00

File Edit View Run Kernel Tabs Settings Help

Files

- notebooks
- transit.ipynb
- passenger.csv
- routes.json
- stops.json

Running

transit.ipynb

Code

Python 3

We plot the number of passengers at the Rosengartenstrasse stop.

```
load = df[df.stopNameShort=='ROSE'].passengerLoadStop
sns.distplot(load, kde=False)
plt.axvline(load.mean())
plt.title('Passenger Load at Rosengartenstrasse stop')
plt.xlabel('Number of passengers')
plt.ylabel('frequency')
```

Passenger Load at Rosengartenstrasse stop

Compare the median load at this stop with the medians of all stops.

```
sns.distplot(df.groupby('stopNameShort').passengerLoadStop.median(), kde=False)
plt.axvline(load.mean())
plt.title('Passenger load medians across all stops')
plt.xlabel('Median passenger load')
plt.ylabel('frequency')
```

Passenger load medians across all stops

stops.json

- features
- properties
- stopId
- stopNumber
- stopNameShort
- stopName
- geometry

passenger.csv

stopId	stopNameShort	stopName
2104	ROSE	Zürich, Rosengartenstr.
664	BUCH	Zürich, Bucheggplatz
2017	RADL	Zürich, Radbühl
498	BIRD	Zürich, Brühl
1708	NEJA	Zürich, Neufeldstr.
1800	OSAU	Zürich, Ochsenschw.
787	ENF	Zürich, Enge

Fun Fact

—

***„Anyone can learn Python,
at least for Data Analytics.“***

Business Fact

—

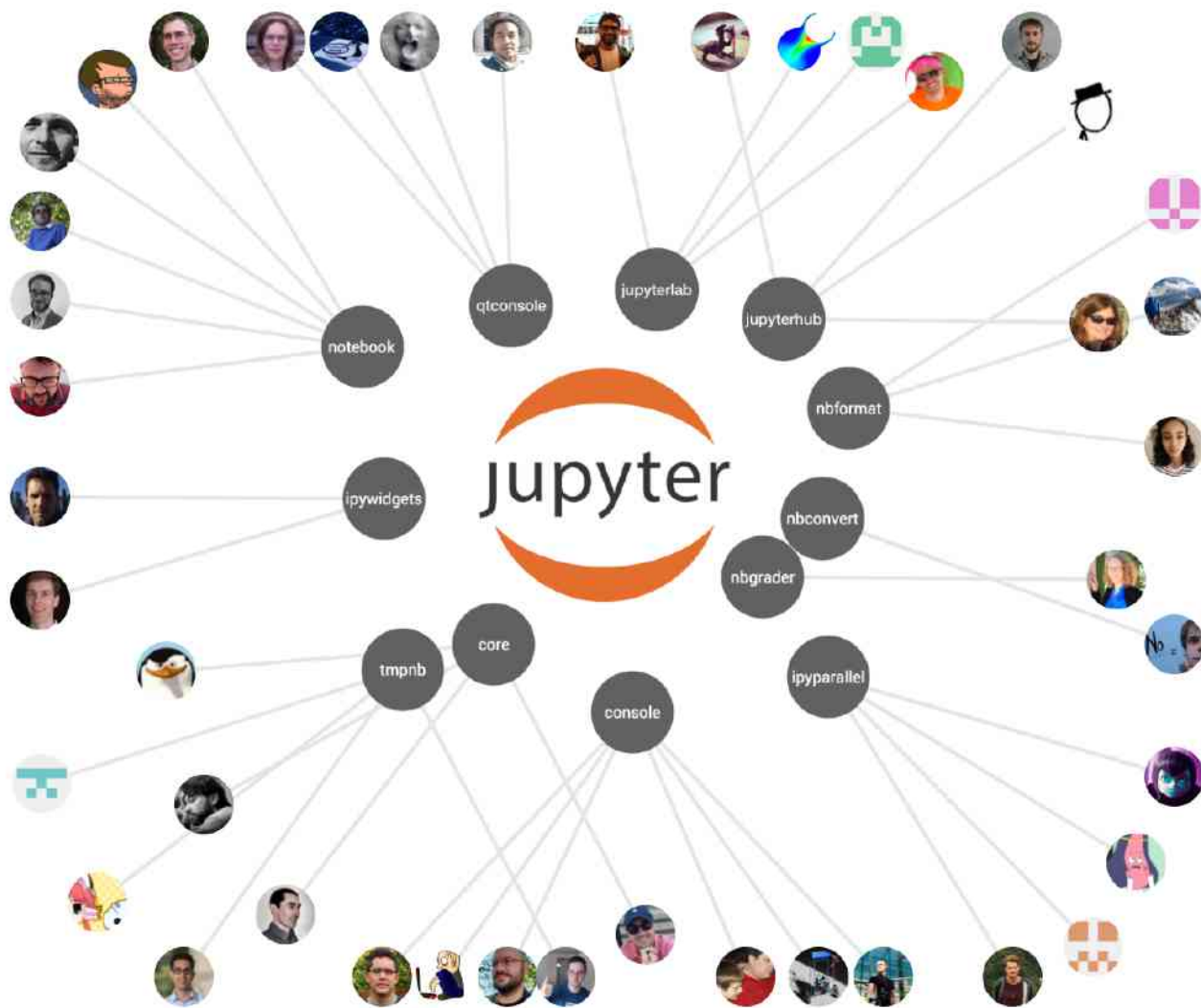
„Python is a perfect common language for a heterogeneous group.“

Jupyter



- *Jupyter supports Python, R, Julia...*
- *Language independent features:*
 - Notebook
 - Message queue
 - Qt-console
- *Open Source, [modified BSD license](#)*







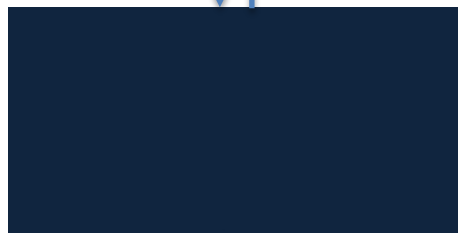
Architecture



Notebook

.ipynb file

HTTP/Websockets



Webserver

OMQ



Kernel

Jupyter Notebooks

- **Document:**
 - Executable code
 - Rich text elements: markdown, LaTeX
 - Visualisations
- **Notebook App:**
 - server-client application allowing editing and running notebook documents via a web browser
- **Kernels:**
 - computational engine
- **Dashboard:**
 - manager



Anaconda Distribution

- Anaconda CPython distribution (covers 2.7 + 3.6)
- Package management *conda*
- 1000+ data-science libraries
- Ensures **packages are compatible** with each other
(newer version of a package may have API changes)
- Provided by Continuum Analytics



Jupyterlab Extensions

- JupyterLab is designed as an extensible environment
- Extensions can customize or enhance any part of JupyterLab
- Extensions are npm packages – JavaScript ;)
- Installation via Extension Manager or command line



KÖNIGSWEG

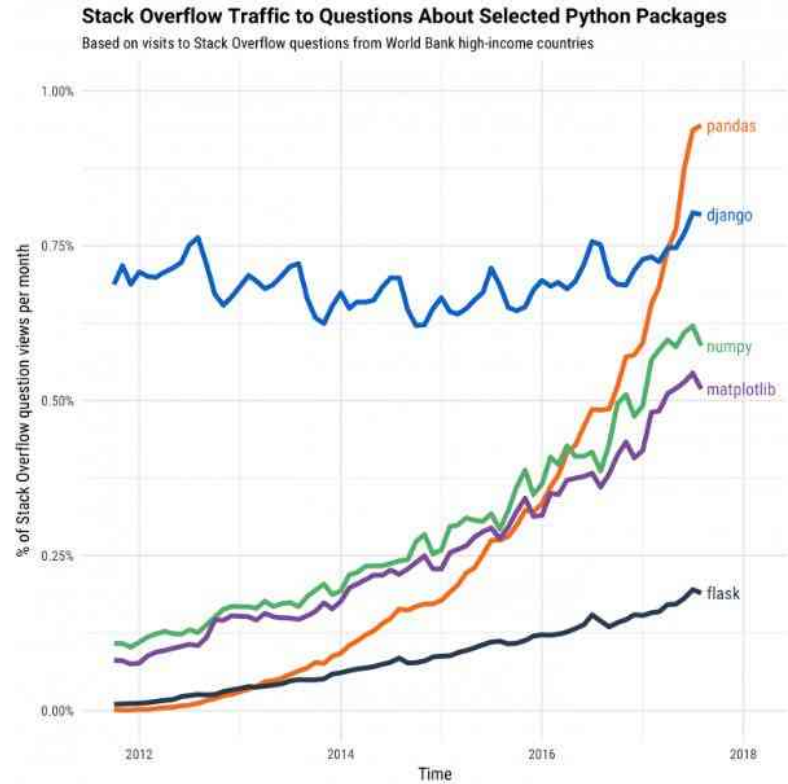


Pandas

- Open Source Python Library
- Praktische *'real-world'*-Datenanalyse - schnell, effizient & einfach
- Lückenloser Datenanalyse Workflow (ohne Wechsel in z.B R)
- 2008 begonnen von Wes McKinney,
nun PyData Stack bei Continuum Analytics ("Anaconda")
- Sehr Stabiles Projekt mit regelmäßigen Updates
- <https://github.com/pydata/pandas>



Development



Pandas Main Features

- Support for CSV, Excel, JSON, SQL, SAS, clipboard, HDF5,...
- Data cleansing
- Re-shape & merge data (joins & merge) & pivoting
- Data Visualisation
- Well integrated in Jupyter (iPython) notebooks
- Database-like operations
- Performant

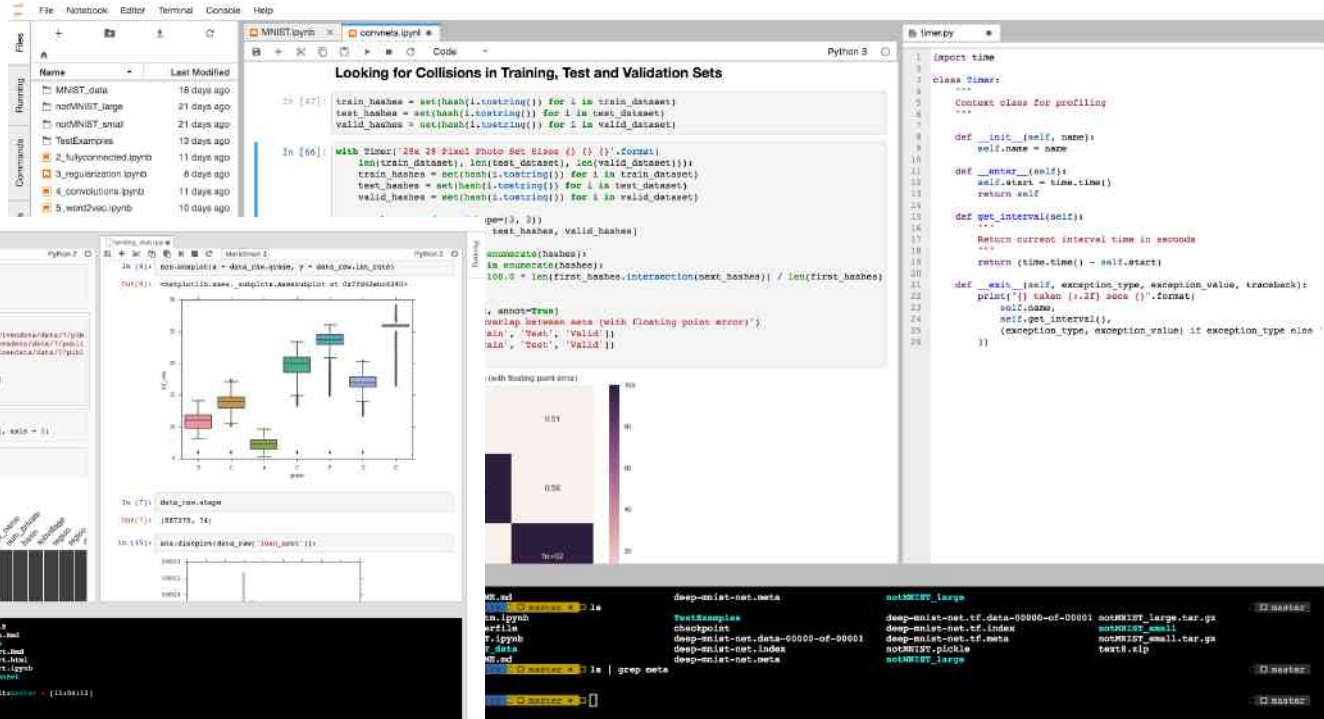


NumPy under the hood

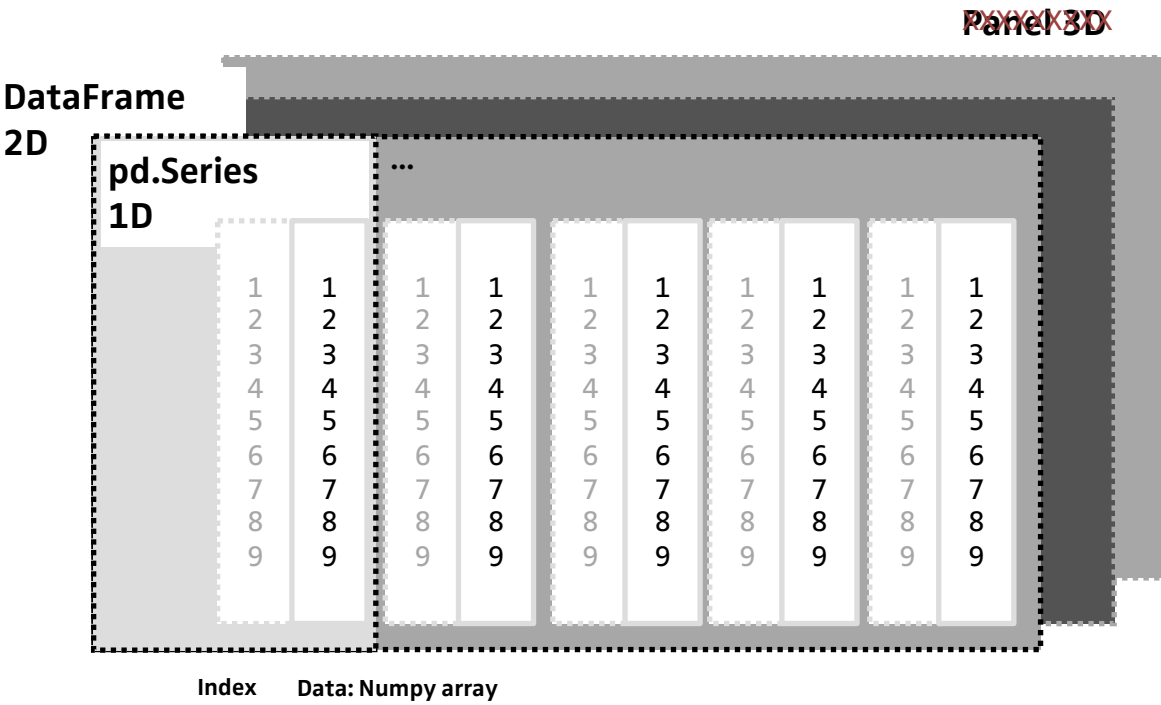


- Library for numerical operations in Python
- Typed Arrays
- Broadcasting

Hands on Pandas



Structure



The Index

- Label of a DataSeries
- Immutable but replaceable
- One or more Dimensions
- Labels are not necessarily *unique*

Index Types

- Index
- MultiIndex
- DateTimeIndex
- TimeDelta
- IntervalIndex
- CategoricalIndex

Basic Stats & Aggregation

- `describe()`
- **Aggregation**
 - sum, count, custom functions,...
 - grouping
 - pivoting
- **NaN (null) values and filler**

Visualization

- Close to the data / code
- Highly customizable
- Many tools: matplotlib, seaborn, bokeh

Automation

- Use nbcovert

Pandas Performance, Limits & Solutions

- Data sets 2-5GB
- stream processing via stepwise aggregation
- Dask for Distributed DataFrames
- Integration with pySpark, SciKit Learn
- Project Arrow

Jupyter Hub

- Jupyter as server for teams
- Collaboration
- Less overhead on local computers
- Access control



Alexander C. S. Hendorf

- Managing Partner & Principal Consultant Information Technology
Consulting on AI & Data Science
- Python Software Foundation Fellow,
Program Chair EuroPython, EuroSciPy, PyConDE & PyData, MongoDB Master
Speaker Europa & USA MongoDB World New York / San José, PyCon Italy, CEBIT
Developer World, BI Forum, IT-Tage FFM, PyData, PyParis, PyData Südwest &
Frankfurt



ah@koenigsweg.com

 @hendorf



Thank you!

Q & A



koenigsweg.com

ah@koenigsweg.com

 @hendorf

