

Datenanalyse mit IPython und Pandas

Dirk Loss, 2013-05-02

**„Data Scientist is now
the hottest job title in Silicon Valley.“**



Tim O'Reilly

Daten

Fokus: „Small/Medium Data“

- Performance-Messungen
- Logdateien
- Netzwerktraffic
- Source-Code Repositories
- ...

Pandas



„Manipulation und Analyse mehrdimensionaler Daten“

	A	B	C
2013-05-01	4	foo	0.40
2013-05-02	6	bar	1.23
2013-05-03	3	boost	0.20
2013-05-04	3	100	-6.20

Time series

Statistics

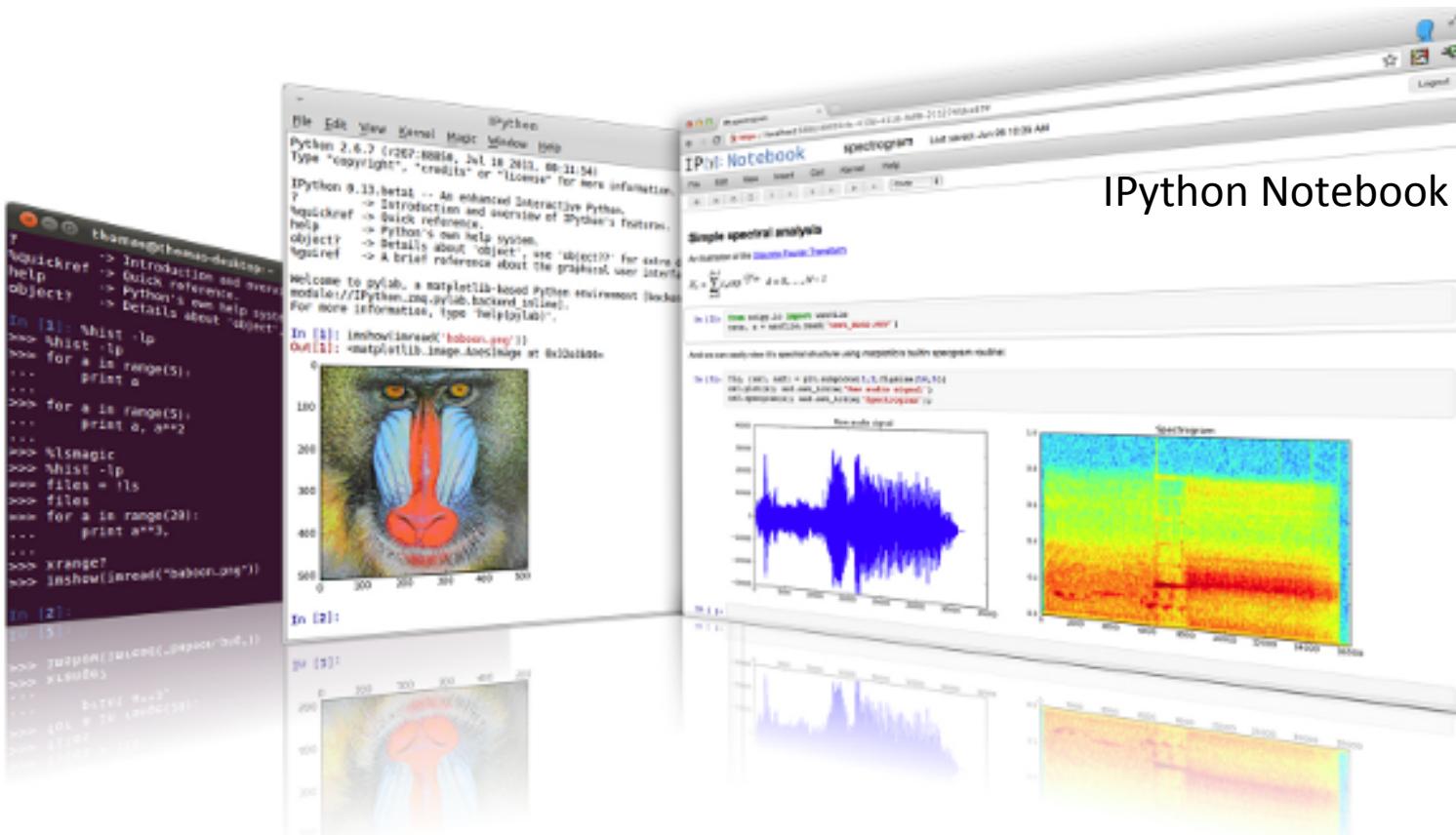
Aggregation

Missing values

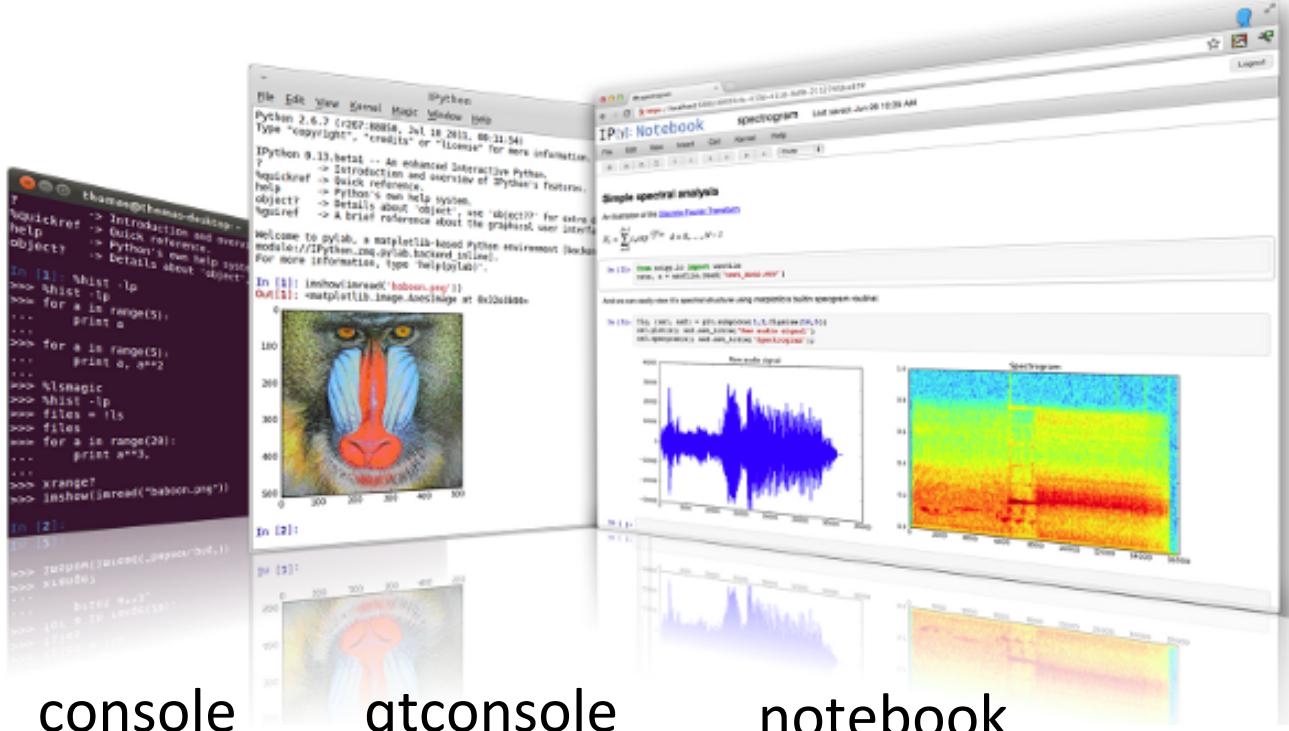
Hierarchical indexes

IPython

„Interaktive Python-Arbeitsumgebung“



<http://ipython.org>



console

qtconsole

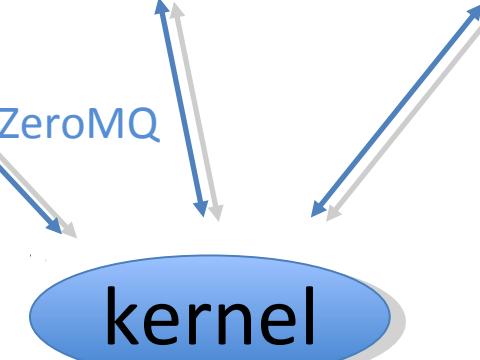
notebook

IPython



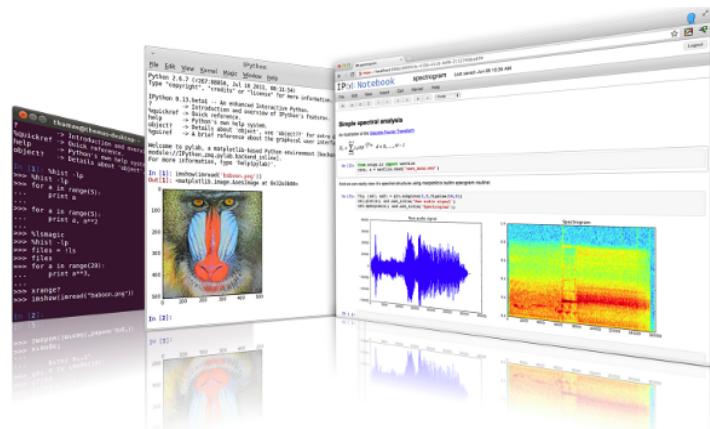
ZeroMQ

kernel



Demo

IPython



nbviewer.ipynb.org IPython Dashboard IPy scratch IPy Pandas Cookbook IPy ipython-tutorial

← → ⌂ nbviewer.ipynb.org

Home FAQ IPython Bookmarklet

IPython Notebook Viewer

A Simple way to share your IP[y]thon Notebook as Gists.

Share your own notebook, or browse others'

Enter a gist number or url

Go!

Probabilistic Programming

Why would I want samples from the posterior, anyways?

We will deal with this question for the remainder of the book, and it is an understatement to say we can perform amazingly useful things. For now, let's finish with using posterior samples to answer the follow question: what is the expected number of texts at day t , $0 \leq t \leq 70$? Recall that the expected value of a Poisson is equal to its parameter λ , then the question is equivalent to what is the expected value of λ at time t ?

In the code below, we are calculating the following: Let i index a particular sample from the posterior distributions. Given a day t , we average over all λ_i on that day t , using $\lambda_{1,i}$ if $t < \tau_1$ else we use $\lambda_{2,i}$.

```
<matplotlib.legend.Legend at 0x148bebe20>
```

Expected number of text-messages received

XKCD Plot With Matplotlib

XKCD plots in Matplotlib

Out [1]:

CHECK IT OUT!

Sometimes when showing schematic plots, this is the type of figure I want to display. But drawing it by hand is a pain in matplotlib. The problem is, matplotlib is a bit too precise. Attempting to duplicate this figure in matplotlib leads to:

<http://nbviewer.ipynb.org>

IP[y]: IPython

Interactive Computing

[Install](#) · [Docs](#) · [Videos](#) · [Notebook Viewer](#) · [News](#) · [Cite](#) · [Donate](#)

Sloan Foundation Grant

We are pleased to announce that the IPython project has received a \$1.15M grant from the Alfred P. Sloan foundation, that will support IPython development for the next two years (1/1/2013-12/31/2014). The grant, which is being made to the University of California, Berkeley and California Polytechnic State University, San Luis Obispo, will enable the project to focus on developing the IPython Notebook as a general tool for scientific and technical computing that is open, collaborative and reproducible.

Google™ Custom Search

Search

x

MEMORIAL

John Hunter

1968–2012

[J. Hunter Memorial Fund](#)

VERSIONS

Stable

0.13.2 – April 2013

Demo

Pandas

	A	B	C
2013-05-01	4	foo	0.40
2013-05-02	6	bar	1.23
2013-05-03	3	boost	0.20
2013-05-04	3	100	-6.20

Pandas DataFrame

columns	foo	bar	baz	qux
index				
A	→ 0	x	2.7	True
B	→ 4	y	6	True
C	→ 8	z	10	False
D	→ -12	w	NA	False
E	→ 16	a	18	False



Indexing in Python

0	1	2	3	4	5	6	7	8	9
H	a	I	I	o		W	e	I	t
-10	-9	-8	-7	-6	-5	-4	-3	-2	-1

`s[0] == "H"`

`s[41] => IndexError`

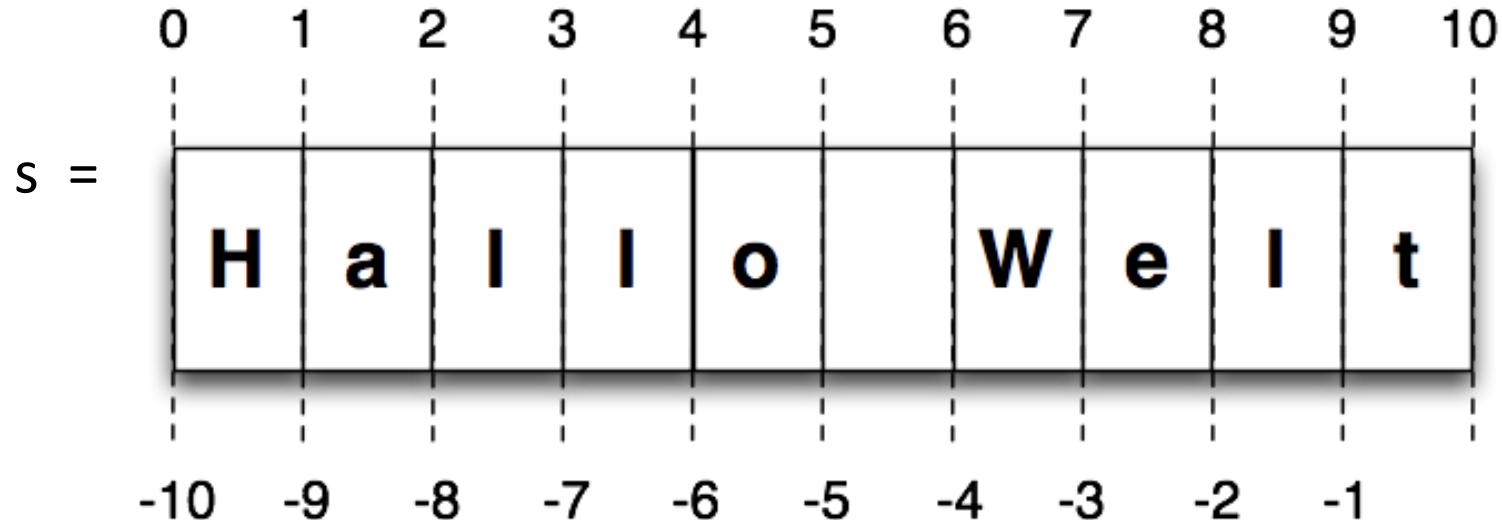
`s[4] == "o"`

`s[-12] => IndexError`

`s[-1] == "t"`



Slicing in Python



`s[2:5] == "lllo"`

`s[6:] == "Welt"`

`s[-3:] == "elt"`

`s[4:-5] == "o"`

`s[6:10000] == "Welt"`

`s[5:2] == ""`

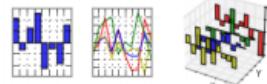


IPython



pandas

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



machine learning in Python



Weiterführende Infos

Agile Tools for Real World Data

Python for Data Analysis



O'REILLY®

Wes McKinney

IPython Videos

A screenshot of a web browser displaying the IPython website (ipython.org/videos.html). The page features a large header with the IPython logo and the tagline "Interactive Computing". Below the header, a navigation bar includes links for "Install", "Docs", "Videos", "Notebook Viewer", "News", "Cite", and "Donate". The main content area is titled "Videos and Screencasts" and contains a link to a video titled "Science And Python: retrospective of a (mostly) successful decade". A video player shows a man speaking at a conference. To the right of the video, there is a sidebar with sections for "MEMORIAL" (listing "John Hunter 1968–2012" and "J. Hunter Memorial Fund"), "VERSIONS" (listing "Stable 0.13.2 – April 2013" and "Development 1.0.dev" with links to "Download" and "Github"), and "COMMUNITY" (links to "Help Chat Room", "Stack Overflow", "Mailing list", "Wiki", and "File a bug").

IPy Videos and Screencasts — [ipython.org/videos.html](#)

IPYTHON: Interactive Computing

Install · Docs · Videos · Notebook Viewer · News · Cite · Donate

Videos and Screencasts

Science And Python: retrospective of a (mostly) successful decade

A historical view of the co-evolution of IPython and the scientific Python stack (1h), that has been very well received. Delivered by [Fernando Perez](#) as a keynote presentation at the PyCon Canada 2012 conference in Toronto, it contains multiple demos of the workflows that IPython enables ([PDF slides](#)).

Science And Python: retrospective of a (mostly) successful decade

0:00:11 / 1:00:58

Google™ Custom Search ×

MEMORIAL

John Hunter
1968–2012
[J. Hunter Memorial Fund](#)

VERSIONS

Stable
0.13.2 – April 2013
[Download](#)

Development
1.0.dev
[Github](#)

COMMUNITY

[Help Chat Room](#)
[Stack Overflow](#)
[Mailing list](#)
[Wiki](#)
[File a bug](#)

IPython Beispiel-Notebooks

This page is a curated collection of IPython notebooks that are notable for some reason. Feel free to add new content here, but please try to only include links to notebooks that include interesting visual or technical content; this should *not* simply be a dump of a Google search on every ipynb file out there.

Important contribution instructions: If you add new content, please ensure that for any notebook you link to, the link is to the rendered version using [nbviewer](#), rather than the raw file. Simply paste the notebook URL in the nbviewer box and copy the resulting URL of the rendered version. This will make it much easier for visitors to be able to immediately access the new content.

Note that [Matt Davis](#) has conveniently written a set of [bookmarklets and extensions](#) to make it a one-click affair to load a Notebook URL into your browser of choice, directly opening into nbviewer.

Entire books or other large collections of notebooks on a topic

- First things first, how to [run code in the IPython Notebook](#), this is one of IPython's official [notebook example collection](#). Another useful one from this group, an explanation of our [rich display system](#).
- A [beautiful matplotlib tutorial](#), part of the fantastic [Lectures on Scientific Computing with Python](#) by J.R. Johansson.
- A [single-atom laser model](#). This is one of a complete set of [lectures on quantum mechanics and quantum optics using QuTiP](#) by J.R. Johansson.
- An [introduction to Compressed Sensing](#), part of [Python for Signal Processing](#): an entire book (and [blog](#)) on the subject by Jose Unpingco.
- An [introduction to Bayesian inference](#), this is just chapter 1 in an ongoing book titled [Probabilistic Programming and Bayesian Methods for Hackers Using Python and PyMC](#) by Cameron Davidson-Pilon

10-Minuten pandas Überblick

The screenshot shows a Vimeo video player interface. At the top, there are three tabs: "10-minute tour of pandas", "A gallery of interesting IPy", and "10-minute tour of pandas". Below the tabs, the URL "vimeo.com/59324550" is visible. The main content area displays a Jupyter Notebook cell. The cell contains Python code using the pandas library. The code generates a DataFrame with columns "Momentum", "Value", and "ShortInterest", and then performs groupby operations to calculate means and z-scores across different industries.

```
index=tickers, name='ccy')

In [276]: df = DataFrame({'Momentum' : np.random.randn(1000) / 200 + 0.03,
                      'Value' : np.random.randn(1000) / 200 + 0.08,
                      'ShortInterest' : np.random.randn(1000) / 200 - 0.02},
                      index=tickers.take(np.random.permutation(Nfull)[:1000]))

df.head()

Out[276]:
      Momentum  ShortInterest    Value
ALDGA  0.017770   -0.023712  0.083033
ARMVB  0.015804   -0.029218  0.081803
CCRDJ  0.032402   -0.028089  0.080964
ZBEBX  0.028301   -0.022696  0.086967
YIIBI  0.026148   -0.021818  0.075221

In [ ]: means = df.groupby(industries).mean()

In [ ]: means.plot(kind='bar')

In [ ]: means = df.groupby([industries, ccy]).mean()
means

In [ ]: keys = [industries, ccy]
zscore = lambda x: (x - x.mean()) / x.std()
summed = df.groupby(keys).apply(zscore)

In [ ]: summed.groupby(keys).agg(['mean', 'std'])

10:28
```

10-minute tour of pandas

from Wes McKinney 2 months ago NOT YET RATED

<http://vimeo.com/59324550>

35-Minuten pandas Vortrag

The screenshot shows a Vimeo video player interface. At the top, the URL vimeo.com/63295598 is visible. The video title is "Data Wrangling Kung Fu With pandas" by Wes McKinney. The video duration is 36:55. The video content displays the PyData logo and the text "SILICON VALLEY 2013". Below the video, there are logos for NumFOCUS and Continuum Analytics. The Continuum Analytics logo includes the text "Platinum Sponsor". The video player has standard controls like play, pause, and volume. Below the player, a summary box provides a brief description of the talk.

Data Wrangling Kung Fu With pandas

from PyData PRO 3 weeks ago / via ContinuumUps NOT YET RATED

In this talk I'll show how a number of tools from the pandas library can be used to quickly wrangle raw data into shape for analysis. Techniques for structured and semi-structured data manipulation, cleaning and preparation, reshaping, and other common tasks will be the main focus.

<http://vimeo.com/63295598>

3-Stunden IPython und Pandas

The screenshot shows a YouTube video player with the following details:

- Title:** Data analysis in Python with pandas
- Uploader:** NextDayVideo
- Views:** 29.413
- Published:** 09.03.2012
- Description:** Wes McKinney
The tutorial will give a hands-on introduction to manipulating and analyzing large and small structured data sets in Python using the pandas library. While the focus will be on learning the nuts and bolts of the library's features, I als
- Thumbnail:** A screenshot of a Jupyter Notebook showing a line plot and some Python code.
- Player Controls:** Shows the video is at 0:20:00 of 3:16:06.

To the right of the video player, there is a sidebar displaying related videos:

- Bjarne Stroustrup: The 5 Programming Languages You Need (2:02)
- What makes Python so AWESOME (1:13:21)
- Boston Algorithmic Finance Meetup with Wes McKinney (50:30)
- Python in Big Data (PyData Workshop, March 2-3, 2012) (49:00)
- Guido van Rossum on the History of Python (1:50:21)
- Programming (von igoronline) (119 Videos)
- IPython: Python at your fingertips (von NextDayVideo) (39:59)
- Tutorial: scikit-learn - Machine Learning in Python with Contributor (von MarakanaTechTV) (1:15:21)

<http://www.youtube.com/watch?v=w26x-z-BdWQ>

Fazit

- IPython als Arbeitsumgebung
(nicht nur für Python)
- Pandas für die Datenanalyse
(insbesondere für Zeitreihen)

Dokumentation

ipython.org

pandas.pydata.org

nbviewer.ipython.org

matplotlib.org

sympy.org

Installation unter Ubuntu

- sudo apt-get install
ipython-notebook
- sudo apt-get install
python-pandas
python-matplotlib
python-scipy
python-sympy
python-nose