

PythonとStapy5年のあゆみ

辻 真吾 (@tsjshg)

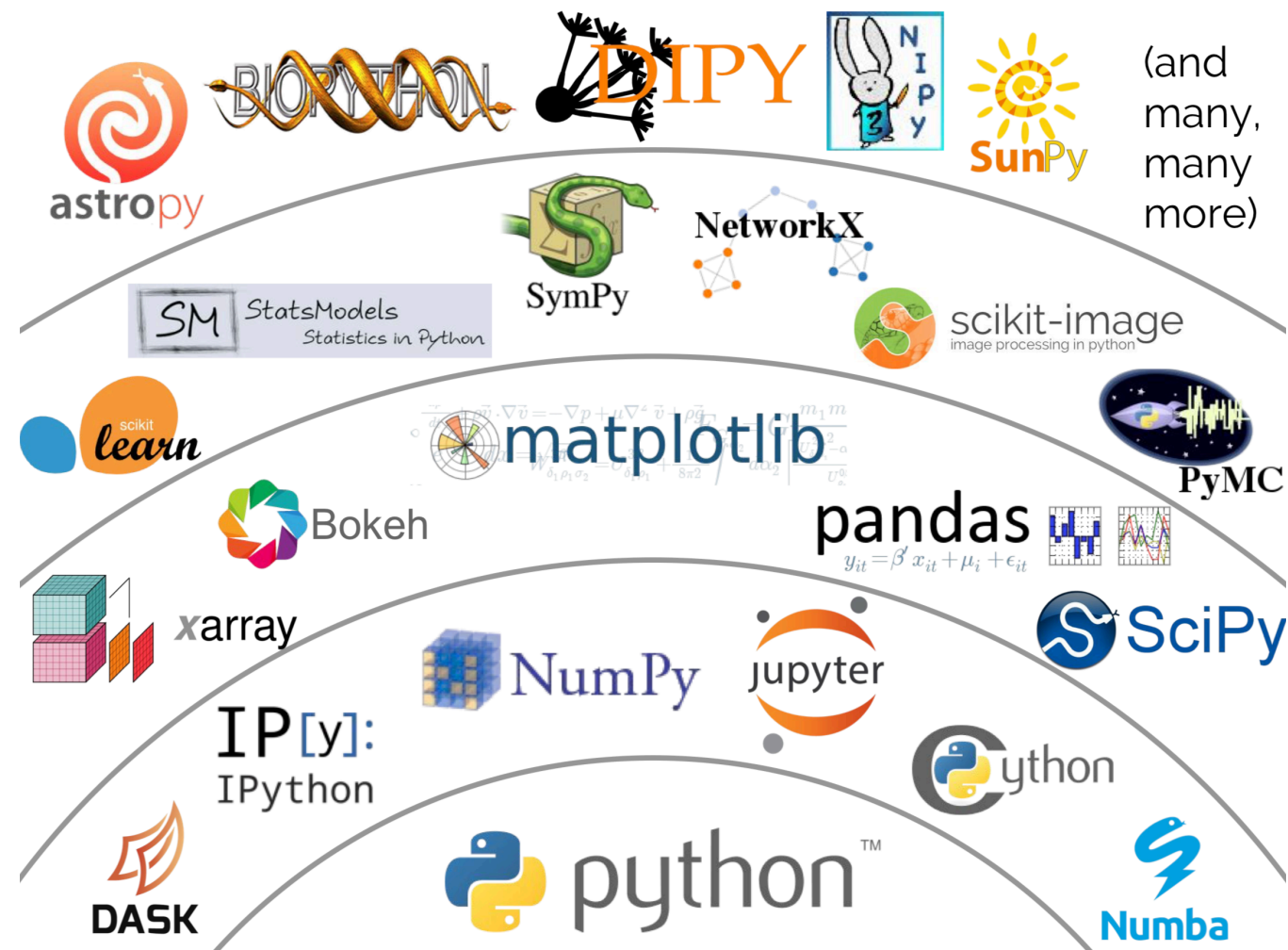
みんなのPython勉強会 # 57 (3度目のオンライン開催)
5/14, 2020

お前、誰よ？

- 東京大学先端科学技術研究センターに勤務
 - 生命科学とコンピュータサイエンスの両方にそこそこ詳しいです
- 詳しくは、www.tsjshg.info をご覧ください
- 今日は、この5年のPython界隈の話題と、stapyを絡めた話をしようと思っ
ていましたが・・・
- ただPythonデータサイエンスに関連するライブラリの話になってしまいました。
た。ごめんなさい。

Python データサイエンス

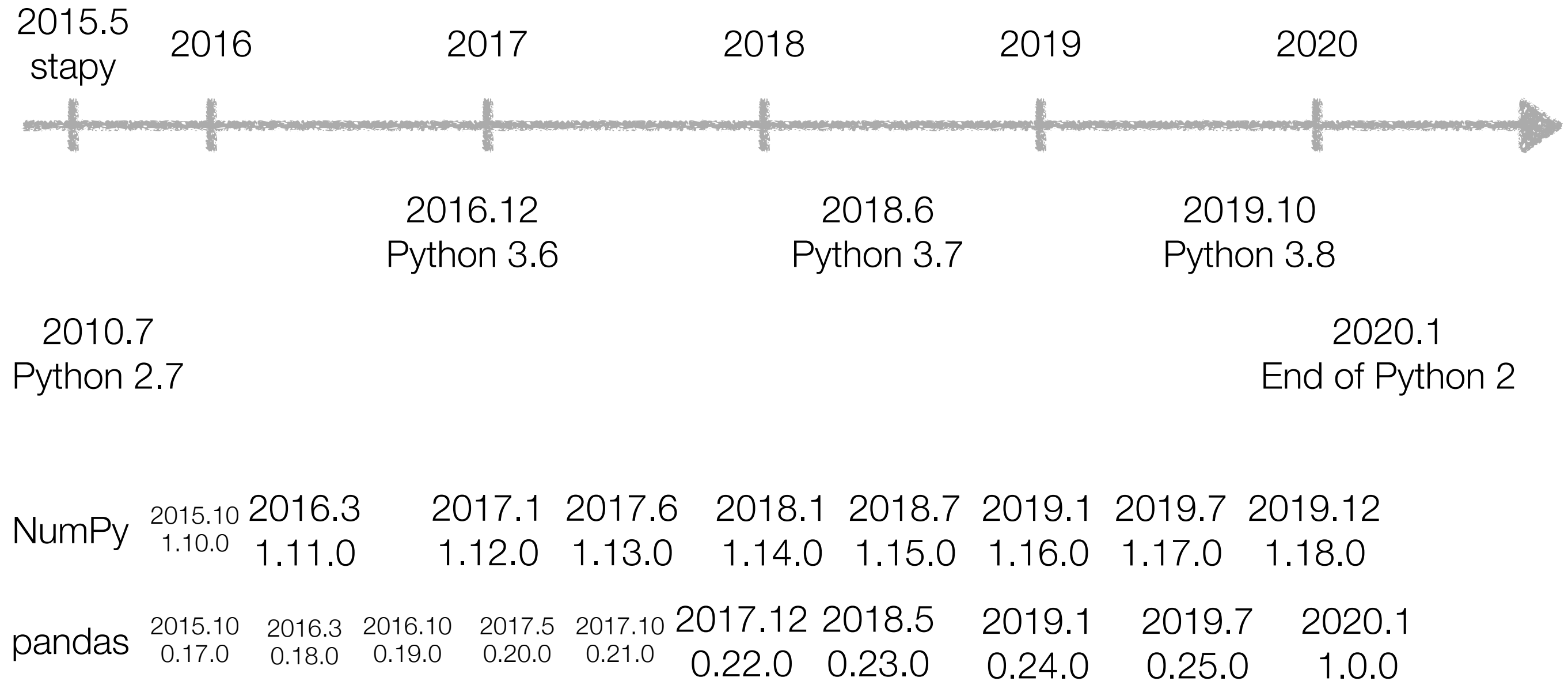
- Pythonは汎用プログラミング言語
- python.org
- データサイエンスを可能にするのはライブラリ群
- pipで個別にインストール
- Anacondaディストリビューション
- condaコマンドも便利



Jake vanderPlas, "The Unexpected Effectiveness of Python in Science", PyCon 2017



Python, NumPy, pandasバージョンの変遷 (5年)



- Python本体の2から3への完全移行
- 外部ライブラリの活発な開発状況

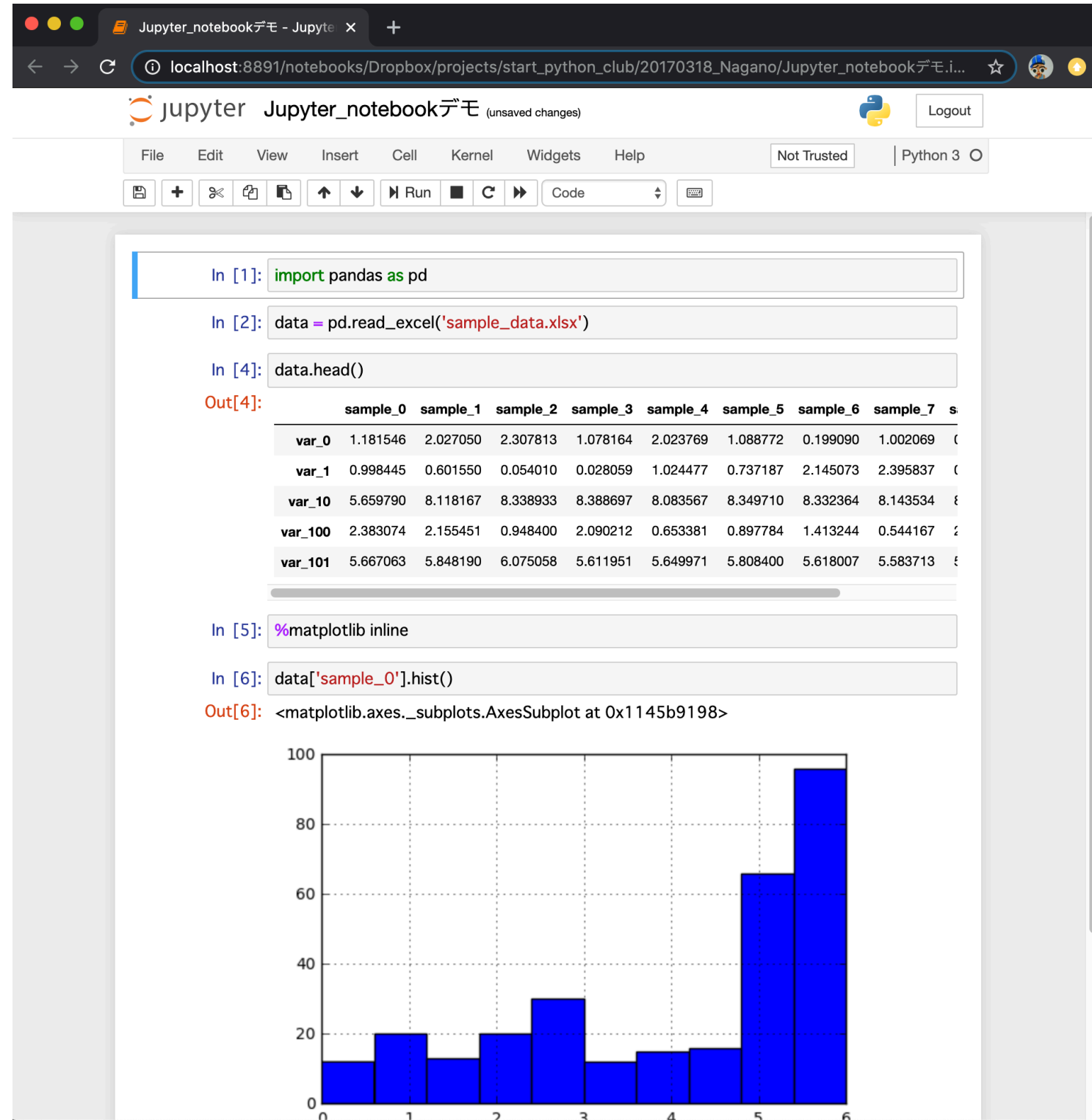
Pythonとライブラリの進化でみるこれまでとこれから

- Jupyter NotebookからJupyter Labへ
- NumPyを支えるテクノロジー
- pandasとZen of Python
- PyCaretの登場

Jupyter NotebookからJupyter Labへ

Jupyter Notebook

- データサイエンスに限らず、Pythonコードの実行環境として人気
- Webブラウザ内でコードの記述、実行、結果の保存が可能
 - 拡張子.ipynb
- Markdownを使ったドキュメントの記述
- カーネルを追加することで、さまざまな言語を利用可能



The screenshot displays a Jupyter Notebook interface in a web browser. The browser address bar shows the URL: localhost:8891/notebooks/Dropbox/projects/start_python_club/20170318_Nagano/Jupyter_notebookデモ.i... The notebook title is "Jupyter_notebookデモ (unsaved changes)". The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help) and a toolbar with icons for file operations and execution. The code area shows the following input cells:

```
In [1]: import pandas as pd
In [2]: data = pd.read_excel('sample_data.xlsx')
In [4]: data.head()
```

The output for In [4] is a table with 10 columns (sample_0 to sample_7) and 5 rows (var_0 to var_101). The output for In [5] is `%matplotlib inline`. The output for In [6] is a histogram of the 'sample_0' column, showing a distribution of values from 0 to 6. The histogram has a y-axis from 0 to 100 and an x-axis from 0 to 6. The bars are blue.

	sample_0	sample_1	sample_2	sample_3	sample_4	sample_5	sample_6	sample_7	sample_8
var_0	1.181546	2.027050	2.307813	1.078164	2.023769	1.088772	0.199090	1.002069	1.002069
var_1	0.998445	0.601550	0.054010	0.028059	1.024477	0.737187	2.145073	2.395837	2.395837
var_10	5.659790	8.118167	8.338933	8.388697	8.083567	8.349710	8.332364	8.143534	8.143534
var_100	2.383074	2.155451	0.948400	2.090212	0.653381	0.897784	1.413244	0.544167	0.544167
var_101	5.667063	5.848190	6.075058	5.611951	5.649971	5.808400	5.618007	5.583713	5.583713

Jupyter Notebookのサンプル画面

The screenshot displays the JupyterLab interface with the following components:

- File Browser:** Shows a directory structure with files like `data.csv`, `OS_plot.png`, `RFS_plot.png`, `依頼_200511.xlsx`, and `生存時間曲線など.ipynb`.
- Code Editor (Top):** Contains Python code for loading data from an Excel file and displaying a table. The output shows a table with columns: `patient`, `OS`, `sensor_os`, `RFS`, and `sensor_rfs`.
- Data Preview (Right):** Shows a preview of the `data.csv` file with columns: `patient`, `OS`, `sensor_os`, and `RFS`.
- Code Editor (Bottom):** Contains Python code for plotting and displays a line plot. The plot shows a linear relationship between two variables, with the y-axis ranging from 1.25 to 3.00.

The bottom of the interface shows a status bar with the text "data.csv" and a small icon.

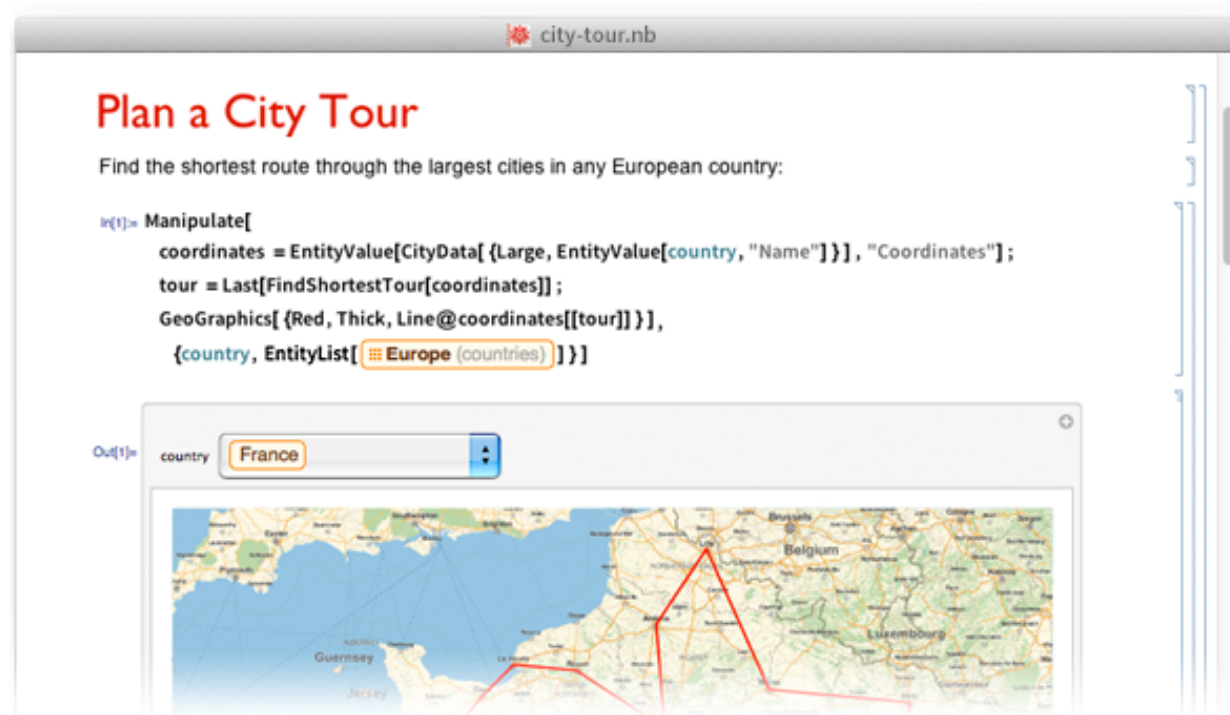
Jupyter Labのサンプル画面

ノートブックの分割表示、CSVや画像のプレビュー、セルのノートブック間でのコピーなどが可能

ノートブックの元祖

ノートブックドキュメント

Wolframノートブックは、テキスト、グラフィックス、インターフェース等とコードを組み合わせることができ、デスクトップでもWebでも使えます：



- おもに数式処理ソフト Mathematicaで使われる Wolfram 言語はノートブックの先駆け
- この Wolfram エンジンは現在無償公開されていて、Jupyter に組み込むことも可能 (らしい)

(拡大)

Wolfram ノートブックは Jupyter のようなノートブックライブラリのきっかけとなっています。

PYTHON のプログラマー向けの注意事項

Wolfram ノートブックは 25 年以上に渡り常に開発され続けており、Wolfram 言語のデスクトップ版とクラウド版のどちらにおいても完全に統合されています。Wolfram ノートブックは Jupyter のようなノートブックライブラリのきっかけとなっています。

NumPyを支えるテクノロジー

NumPyとは？

- Pythonで数値計算をするときの必須ライブラリ
- ベクトルや行列の計算が効率良くできるnumpy.ndarray型
 - Pythonのリストと違い、一度定義するとサイズを変更できず、格納できるデータ型も1種類
- 要素のアクセスや演算など便利な機能がたくさん

```
>>> a[0,3:5]
array( [3,4] )

>>> a[4:, 4:]
array( [ 28, 29],
       [ 34, 35] )

>>> a[:, 2]
array( [2, 8, 14, 20, 26, 32] )

>>> a[2::2, ::2]
array( [ 12, 14, 16],
       [ 24, 26, 28] )
```

0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
18	19	20	21	22	23
24	25	26	27	28	29
30	31	32	33	34	35

とにかく速い

この速度比較はほとんど意味ないですが・・・

```
In [1]: import numpy as np
```

```
In [5]: vec_a = np.random.randn(200_0000)
vec_b = np.random.randn(200_0000)
```

200万次元のベクトル2つ

```
In [6]: %%timeit
vec_a @ vec_b
```

内積を計算

877 μ s \pm 12.8 μ s per loop (mean \pm std. dev. of 7 runs, 1000 loops each)

```
In [ ]:
```

```
In [7]: %%timeit
sum([a * b for a, b in zip(vec_a, vec_b)])
```

609 ms \pm 12.1 ms per loop (mean \pm std. dev. of 7 runs, 1 loop each)

CPUコアがいくつかあると並列化してくれる

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

```
In [3]: # 10000次元の説明変数を持つサンプルが4万個  
X = np.random.randn(4_0000, 10000)
```

```
In [*]: # Mは4万×4万の正方行列 i行j列はX[i]とX[j]の内積  
M = X @ X.T
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

アクティビティモニタ (自分のプロセス)

プロセス名	% CPU	CPU時間	スレッド	アイドル・ウェ...	% GPU
python3.7	591.4	1:35.89	14	1	0.0
アクティビティモニタ	10.5	54:54.60	6	4	0.0
Google Chrome Helper (Renderer)		58.35	25	26	0.0
Google Chrome		8:22:21.99	34	25	0.0
Google Chrome Helper (GPU)		2:52.42	13	47	4.1
Google Chrome Helper (Renderer)		2:22.83	27	2	0.0
Google Chrome Helper (Renderer)		1:11:43.97	21	17	0.0
Google Chrome Helper (Renderer)		16:06.99	31	255	0.0
Google Chrome Helper (Renderer)		1:18:34.87	22	13	0.0
Google Chrome Helper (Renderer)		22:11.95	24	1	0.0
Google Chrome Helper (Renderer)		5:15:05.01	22	2	0.0
Google Chrome Helper (Renderer)		5:41:28.88	19	1	0.0
Google Chrome Helper (Renderer)		3:08.59	17	1	0.0
Google Chrome Helper (Renderer)		10:46.72	16	2	0.0
Google Chrome Helper (Renderer)		20:46.73	17	2	0.0
Google Chrome Helper (Renderer)		40.26	21	9	0.0
Google Chrome Helper (Renderer)		2:33.22	15	1	0.0

4万行×1万列の行列Xを用意

行列Mのi, j成分はX[i]とX[j]の内積になる

(単純には) 16万回の内積計算

100%以上なので
並列化してくれている

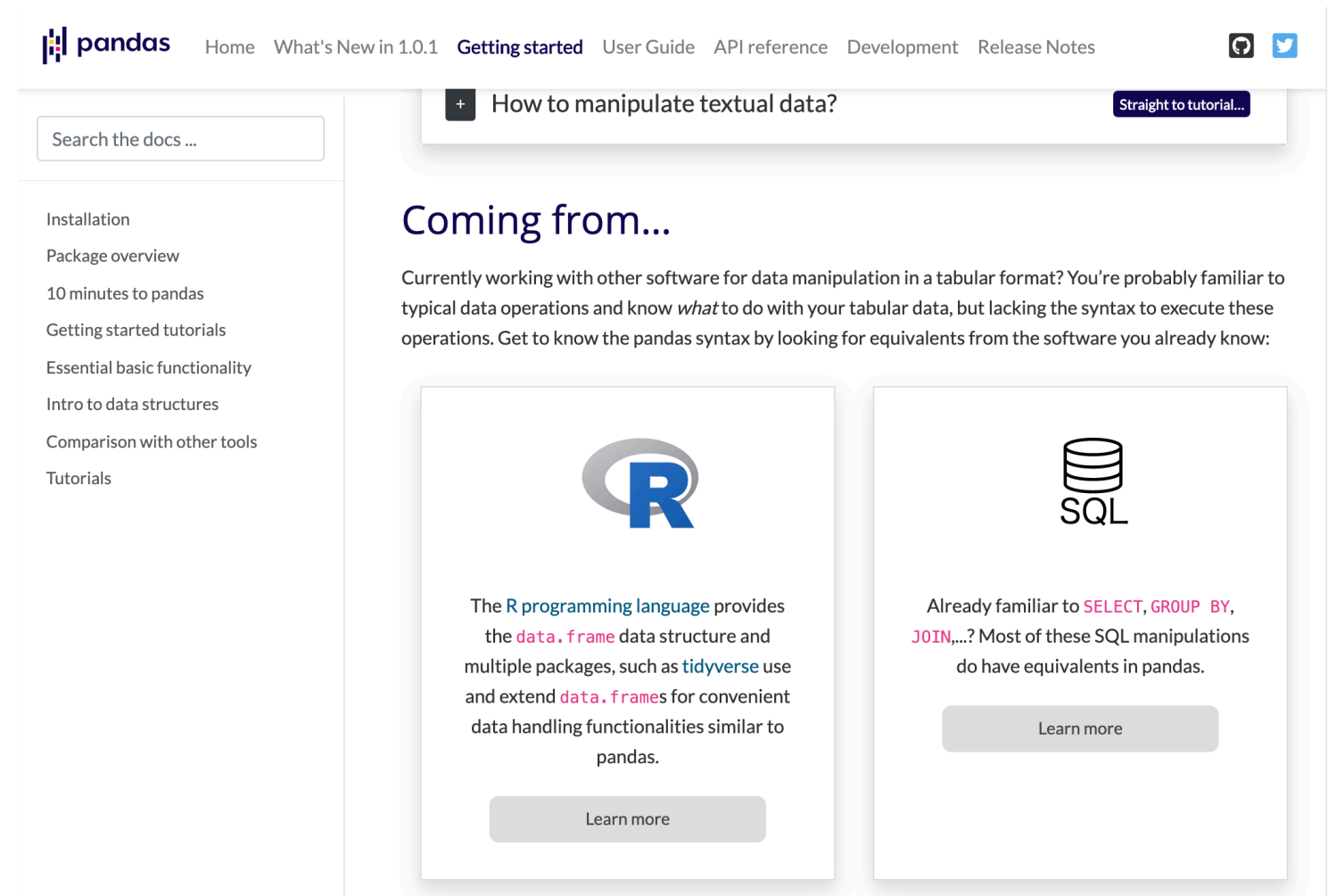
Basic Linear Algebra Subprograms (BLAS)

- ベクトルや行列など線形代数の演算を実行するAPIのデファクトスタンダード (Reference BLAS)
 - 1970年代からの歴史がある*
- 最適化されたOpenBLASやIntelが無料で配布しているMath Kernel Library (MKL) などがある
- `numpy.show_config()`で確認できる
 - AnacondaディストリビューションはMKLを利用

pandas と Zen of Python

pandas

- 行と列に名前が付いた表形式の2次元データ（DataFrame型）と1次元のSeries型を提供する超高機能ライブラリ
- Pythonデータサイエンスに欠かせない存在
- 機能が多すぎるので、最低限これだけという記事*もある



The screenshot shows the pandas documentation website. The main content area is titled "Coming from..." and explains that users familiar with other data manipulation software like R or SQL can find equivalents in pandas. It features two cards: one for R, stating that R provides the `data.frame` structure and that packages like `tidyverse` extend it for pandas-like functionality; and one for SQL, stating that many SQL operations like `SELECT`, `GROUP BY`, and `JOIN` have equivalents in pandas. Both cards include a "Learn more" button.

RのDataFrameやSQLのテーブルからの影響

* <https://medium.com/dunder-data/minimally-sufficient-pandas-a8e67f2a2428>

やり方がいろいろあるすぎる

```
import pandas as pd

df = pd.DataFrame({'足': [2,4,6], 'genome(Gbp)':[3.2, 2.4, 0.5], 'size': [1.65, 1.7, 0.005]},
                  index=['ヒト', 'パンダ', 'アブラムシ'])

df
```

	足	genome(Gbp)	size
ヒト	2	3.2	1.650
パンダ	4	2.4	1.700
アブラムシ	6	0.5	0.005

```
df.足
```

```
ヒト    2
パンダ  4
アブラムシ  6
Name: 足, dtype: int64
```

```
df.size
```

列へのアクセスはこちらがおすすめ

```
df['size']
```

```
ヒト    1.650
パンダ  1.700
アブラムシ  0.005
Name: size, dtype: float64
```

Zen of Pythonに反する？

```
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!

ixメソッドの廃止

- DataFrameの要素へのアクセス
 - 行や列の名前ならloc
 - 行や列の位置ならiloc
- これらを混在できるixメソッドが0.20まで使えた
- 賛否はあった*が廃止を決定
- “Explicit is better than implicit.”の哲学に沿った変更と言える（実際は別の理由かも）

```
df
```

	足	genome(Gbp)	size
ヒト	2	3.2	1.650
パンダ	4	2.4	1.700
アブラムシ	6	0.5	0.005

```
df.loc['パンダ', '足']
```

 行・列名で指定

```
4
```

```
df.iloc[1, 0]
```

 行・列の位置で指定

```
4
```

PyCaretの登場

caret 機械学習を便利にするRのライブラリ

- 機械学習モデルを作る時によくある作業を自動化
 - データの分割、クロスバリデーションやハイパーパラメータの調整など
- scikit-learnを使っても10行以上になるような作業を1～数行で実行

The `caret` Package

Max Kuhn

2019-03-27

1 Introduction

The `caret` package (short for **C**lassification **A**nd **R**Egression **T**raining) is a set of functions that attempt to streamline the process for creating predictive models. The package contains tools for:

- data splitting
- pre-processing
- feature selection
- model tuning using resampling
- variable importance estimation

as well as other functionality.

There are many different modeling functions in R. Some have different syntax for model training and/or prediction. The package started off as a way to provide a uniform interface the functions themselves, as well as a way to standardize common tasks (such parameter tuning and variable importance).

The current release version can be found on [CRAN](#) and the project is hosted [on github](#).

Some resources:

でた！ (2020年4月)

Home About Install Guide Contribute Git

Getting Started Functions Preprocessing Modules Tutorials

An open source **low-code** machine learning library.

**POWERFUL SOLUTION
LOW-CODE**

Productivity tool Easy to use Business ready solution

Why PyCaret

PyCaret is an open source, **low-code** machine learning library in Python that allows you to go from preparing your data to deploying your model within seconds in your choice of notebook environment.

<https://pycaret.org/>

scikit-learnを含む既存ライブラリのラッパー
caretと名前は似ているが、作者Moez Ali氏の1人プロジェクト？

まとめ

- Pythonデータサイエンスは外部ライブラリとその進化に支えられている
 - たとえば、IPython→Jupyter→JupyterLab
- Pythonの哲学（Zen of Python）とライブラリの変化
- 少ないコードで大きな効果
 - Numpyを支えるOpenBLASやMKLによる並列化
 - PyCaretに代表される機械学習自動化の流れ

当日時間切れで対応できなかった質問への回答



Anonymous

7:54 PM

5

pycaretに限らずpythonはscikit-learnなど機械学習用の便利なライブラリが多いです。正直良く中身がわかってなくても使えてしまいます。辻さんはどこまで外部ライブラリの中身・原理を理解しておくべきだと思いますか。

機械学習モデルでしたら、手法の概要くらいは理解するべきだと思っています。ロジスティクス回帰とランダムフォレストの違いがまったくわからないのに、ライブラリを使うのはすこし危険な気がします。一方で、実際の実装を細かく見るのは難しいので、そこまでは気にしなくてもよい気がします。自分で実装したり、だれかの実装を読むのはすごく勉強になるのでおすすですが、時間がかかるので、連休中とか盆暮れ正月あたりの課題ですかね。



Anonymous

7:37 PM

3

stapyの5年間で数値計算ライブラリに起こった変化のうち、辻さんがイチオシのものはなんですか？



本編でもお話ししたIntel Math Kernel Library (MKL) が Anacondaにデフォルトで組み込まれたのが、（調べたら） Anaconda 2.5（2016年2月リリース）だったので、stapy5年の範囲に入っているのをこれを挙げたいと思います。MKLはすごい値段が高いライブラリだというイメージがあったので、このニュースを聞いた時は驚きました。マイクロソフトのVSCodeもそうですが、最近IT業界の老舗巨大企業がオープン化（MKLはオープンソースではないですが）へ舵を切っているように見えて、一般ユーザには嬉しい流れです。



Anonymous

7:55 PM



pandas だから

df.loc['パンダ', '足']

だったんですか・・・？

はい。



Anonymous

7:58 PM new

1

今後、Just In Timeコンパイルが今後どれだけ高速化を達成するのか、気になっています。

PythonのJITと言えばNumbaが有名ですよ。昔はpyscoなんかもあって、手軽に速くなるので、すこし使っていた時期もありましたが、残念ながらpyscoは2012年にメンテされなくなってしまいました。PyPyはいまでも活発に開発されているようで、使ってみたい気もしますが、結局作るプログラムが外部ライブラリに依存することが多く、それほど高速化の恩恵を受けられないこともあったりするので、難しいですよ。



Anonymous

7:41 PM

0

`np.random.randn(200_0000)` で200万個の乱数を計算するんですね？

大きな整数はアンダースコアで区切って表記できます。普通3桁ごとですが、200万だとこの方が分かり易いかなと思ひまして。

また、`np.random.randn(3, 4)`とすると、3行4列の行列になって返ってくるので便利です。



ほとはら

7:46 PM

0 

Numpyが自動で並列処理してくれるなら自前で並列処理を組むのは無駄なことだろうか

これはそうではないと思います。ちょっとした計算をするときに並列化してくれるのは便利ですが、独立した手続きにできる処理を自分で並列化する方が全体として高速化できる可能性は十分あると思います。

ただ、こうしたチューニングは作業としては一番最後にもってくるべきかと思います。まず動くものを作って、その後、自分が書いたコードの範囲で無駄なことをやっていないかチェックした後という感じです。（だいたい、このあたりで力尽きて、まあいっかとなりそうですが・・・）