BPStudy#121~地に足をつけて学ぶ機械学習、データサイエンス

今データサイエンスが必要 とされる理由とPythonの役割

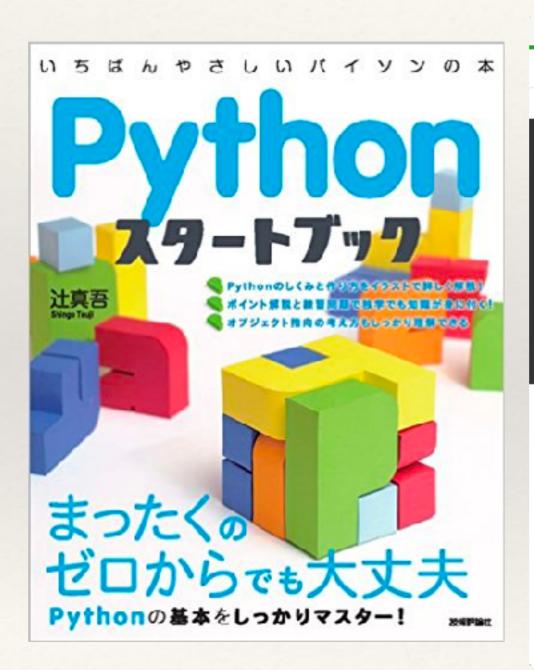
辻真吾@tsjshg

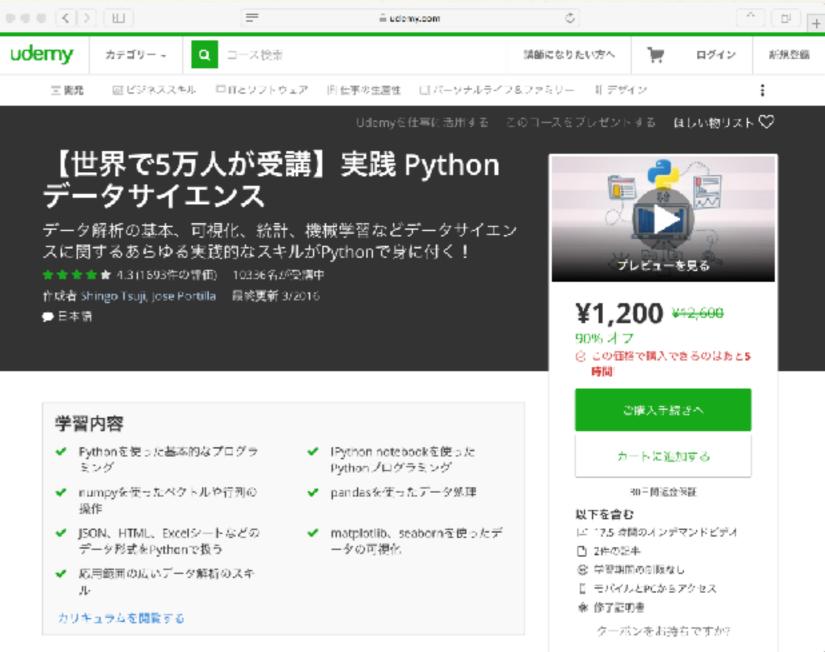
2017.9.26

## 自己紹介

- \* 1975年生まれ東京都足立区出身
- \* コンピュータと数学でなんでもやる数理工学を専攻
- \* 修士終了後、創業間もない(株)いい生活に就職
  - \* Javaを使ったWebアプリ開発
- \* 3年弱で会社を辞めて、現在の勤め先(東大先端研)に博士課程の学生 として復学
  - \* バイオインフォマティクス (生命情報科学)
- \* 42才ですが、人生の次の一手を模索中です

## 本とWeb





https://www.udemy.com/python-jp/

## お知らせ1



https://startpython.connpass.com/

## お知らせ2



https://edgeai.connpass.com/

### 本日の目次

- \* データサイエンスとその必要性
- \* Pythonの役割
  - \* 汎用言語でデータ解析
  - \* 解析の再現性とオープンソース
- \* データサイエンスのこれから

### データサイエンスとは?

### データ駆動型サイエンス



Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia
Wikipedia store

Interaction

Help About Wikipedia Community portal Recent changes Contact page Article Talk Read Edit View history

Not logge

#### Data science

From Wikipedia, the free encyclopedia

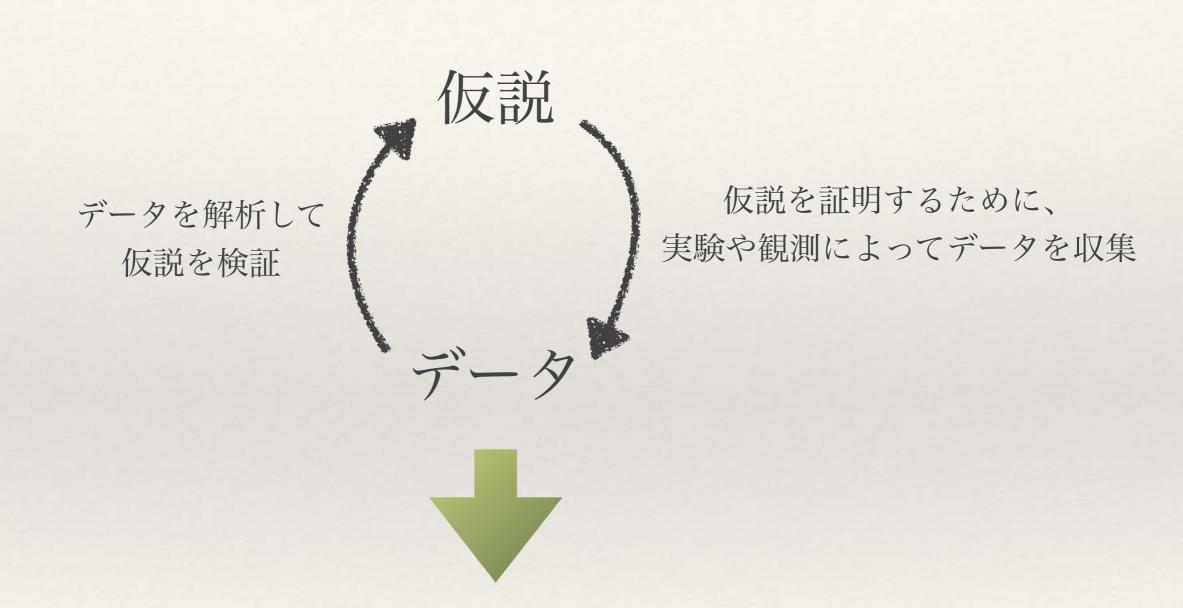
Not to be confused with information science.

**Data science**, also known as **data-driven science**, is an interdisciplinary field about scientific methods, processes, and systems to extract knowledge or insights from data in various forms, either structured or unstructured,<sup>[1][2]</sup> similar to data mining.

Data science is a "concept to unify statistics, data analysis and their related methods" in order to "understand and analyze actual phenomena" with data.<sup>[3]</sup> It employs techniques and theories drawn from many fields within the broad areas of mathematics, statistics, information science, and computer science, in particular from the subdomains of machine learning, classification, cluster analysis, data mining, databases, and visualization.

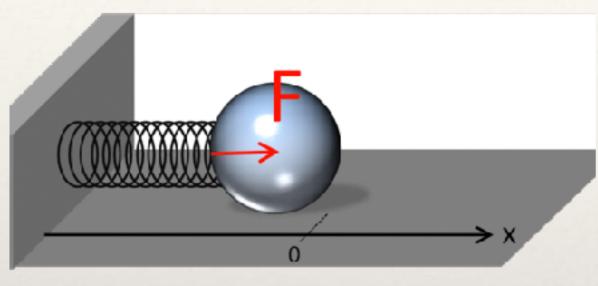
Turing award winner Jim Gray imagined data science as a "fourth paradigm" of science (empirical, theoretical, computational and now data-driven) and asserted that "everything about science is changing because of the impact of information technology" and the data deluge. [4][5]

### サイエンス

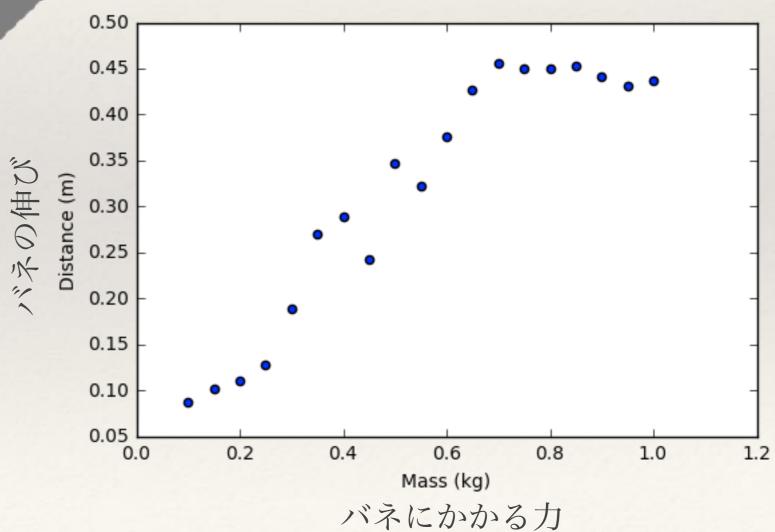


普遍的な原理

### たとえば



フックの法則 F = -k x



### DNA配列だと

- \* ヒトの細胞には、約30億文字 分のDNA情報が格納されている
- \* これが最近の装置の進歩によって1週間弱で解読可能
- \* 100人分やったら、3,000億文字

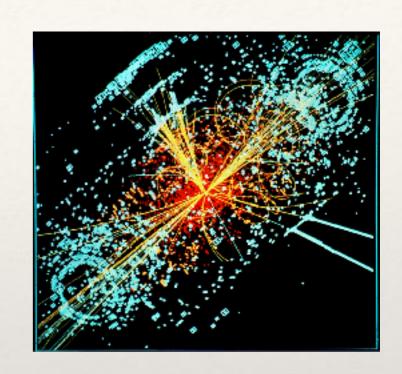




https://jp.illumina.com/systems/sequencing-platforms/hiseq-2500.html

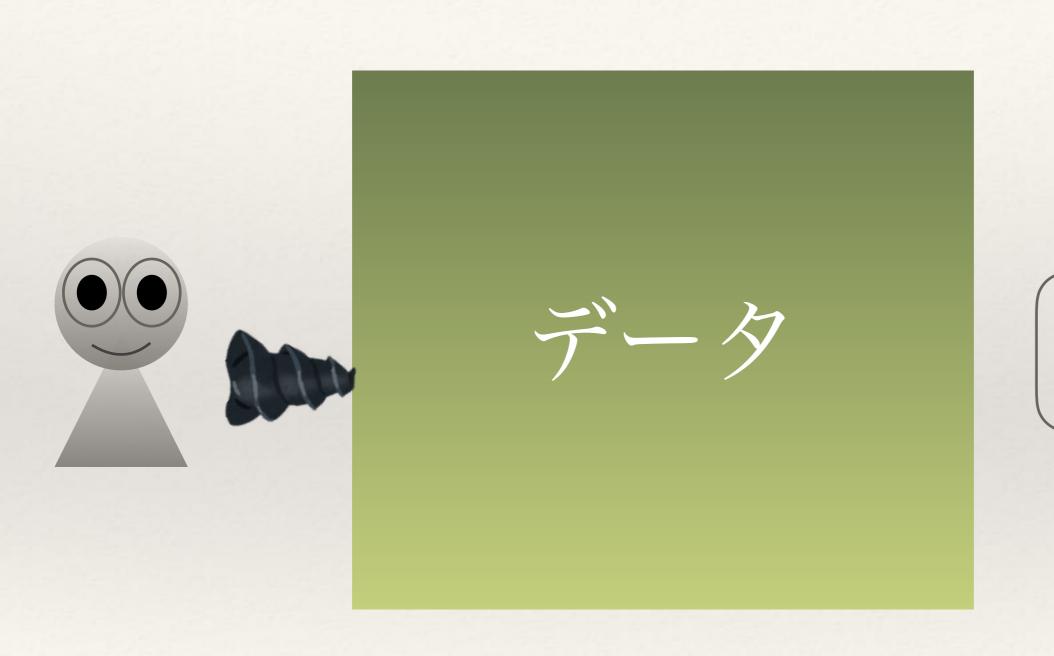
## とにかくビッグデータ

- \* ヒッグス粒子の発見
  - \* データ解析にはPythonが 使われたらしい
- \* Web上のお客さんの行動履歴
- \* IoT
  - \* 地理情報を含む行動履歴
  - \* エネルギー消費量





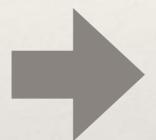
## データサイエンスの必要性



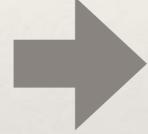
普遍的な何か

## データサイエンスの実際

データの収集と前処理



統計解析 機械学習



- データを集める (CSV, Excel, DB...)
- Webスクレイピング (通信とデータの整形)
- 欠損値や異常値の検出と処理

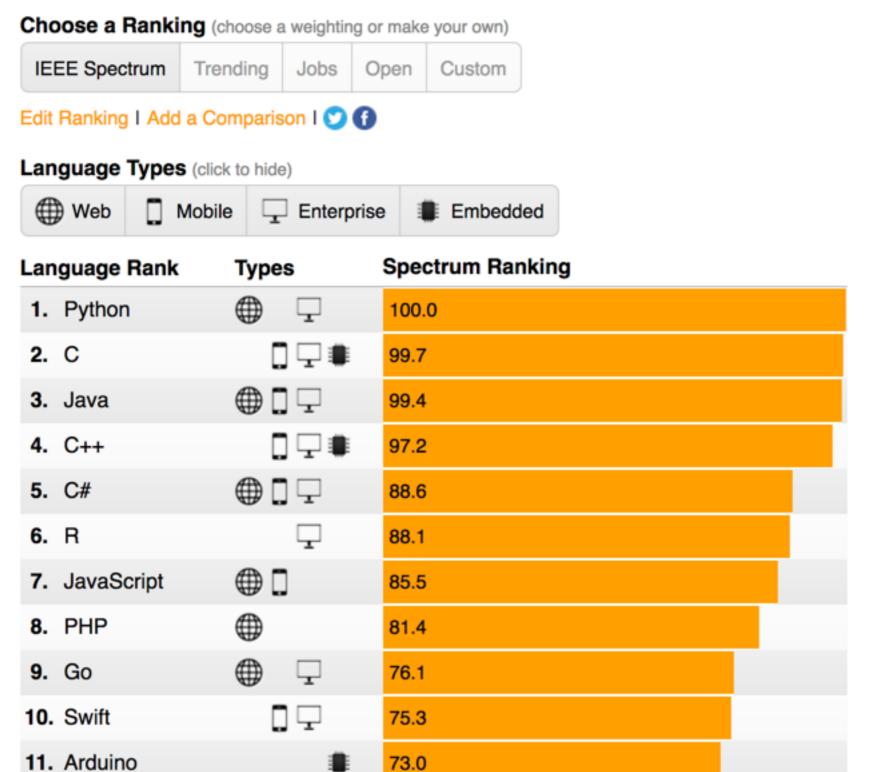
後処理と次の計画

完成品のWebアプリ化

汎用言語でデータサイエンスを実践する意味は大きい

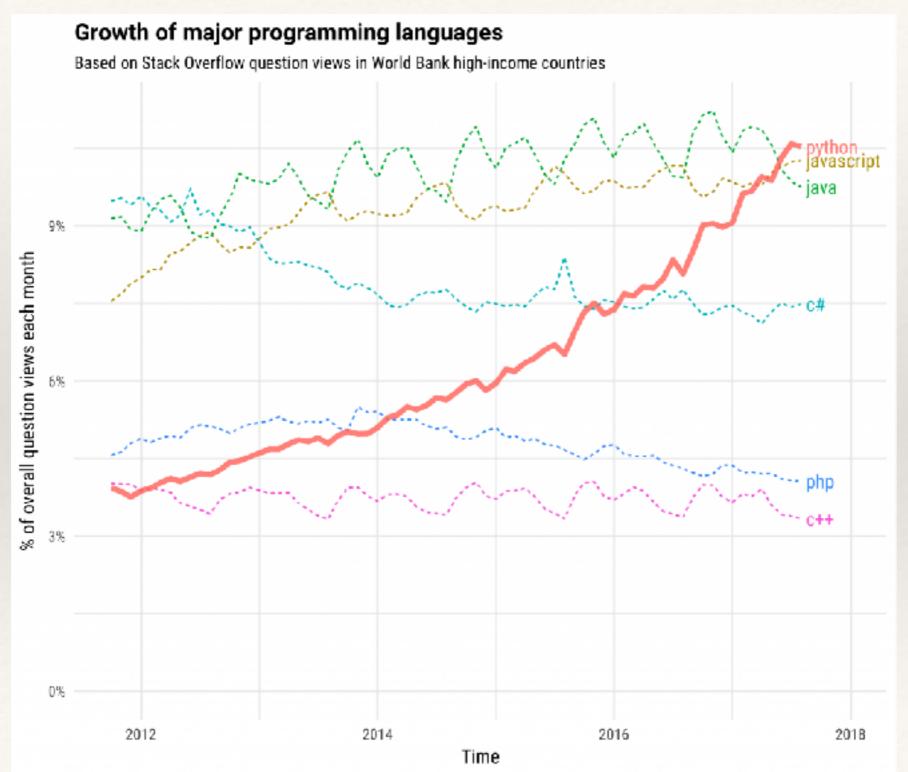
# Python急成長の真相

## IEEE Spectrum





### Stack Overflow



### いくつかの選択肢

- \* SAS (www.sas.com)
  - \* BIツールの代表格
- \* Mathematica (www.wolfram.com)
  - \* 天才ウルフラムが作った老舗数式処理ソフト
- Matlab (jp.mathworks.com)
  - \* 数値シミュレーションなどに利用されている
- \* R
  - \* オープンソースのデータ解析専用言語
- \* Python
  - \* この中では唯一、オープンかつ汎用





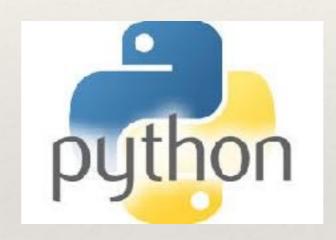


## Pythonはglue (のり) 言語





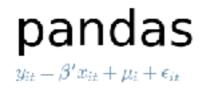




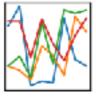
















### Anacondaがおすすめ!

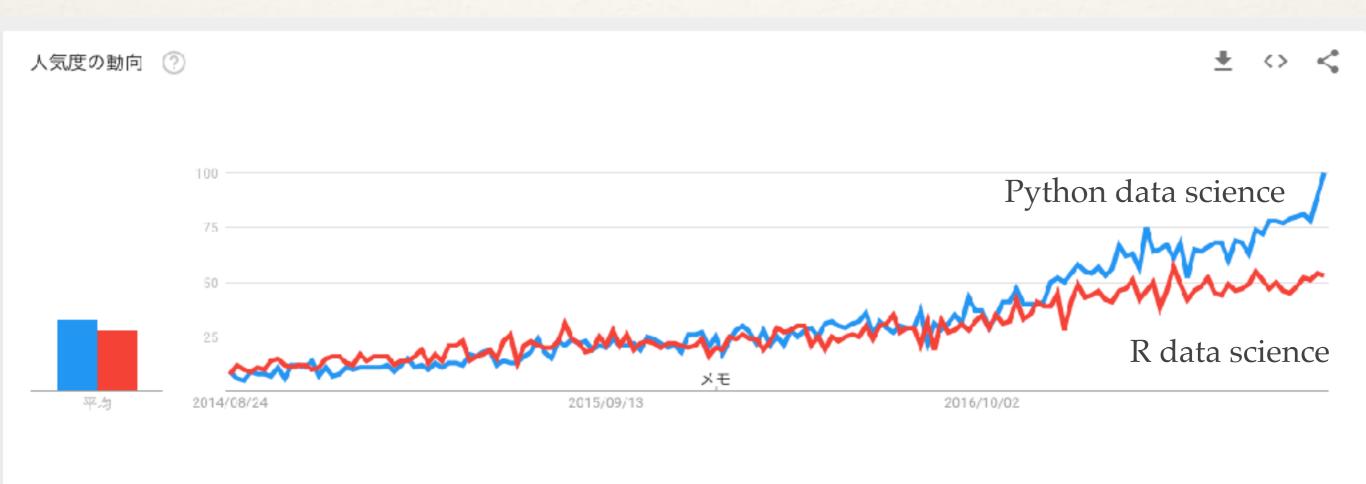
- \* Continuum Analytics社が配布 するPython
  - \* 標準のPythonにcondaをはじ めとして多くの外部ライブラ リ(データ解析用が中心) を同梱
  - \*無料
  - \* データ解析環境の構築に最適



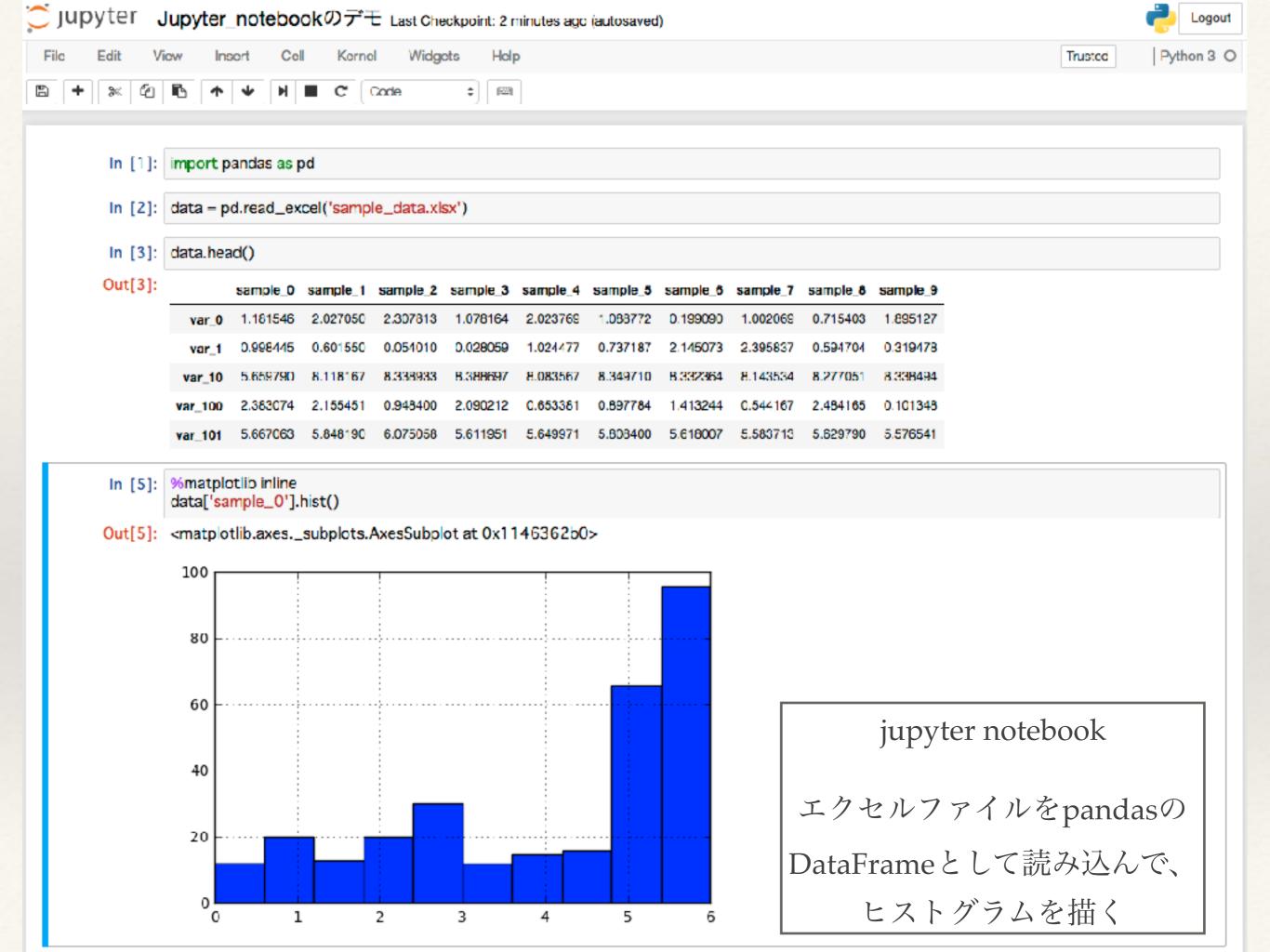
https://www.continuum.io

# R & Python

過去3年間のGoogle Trendsのデータ



データサイエンス自体が伸びているので、もちろんRも中心的な存在



In [6]: import seaborn as sns sns.pairplot(data)

Out[6]: <seaborn.axisgrid.PairGrid at 0x1146b7c88>



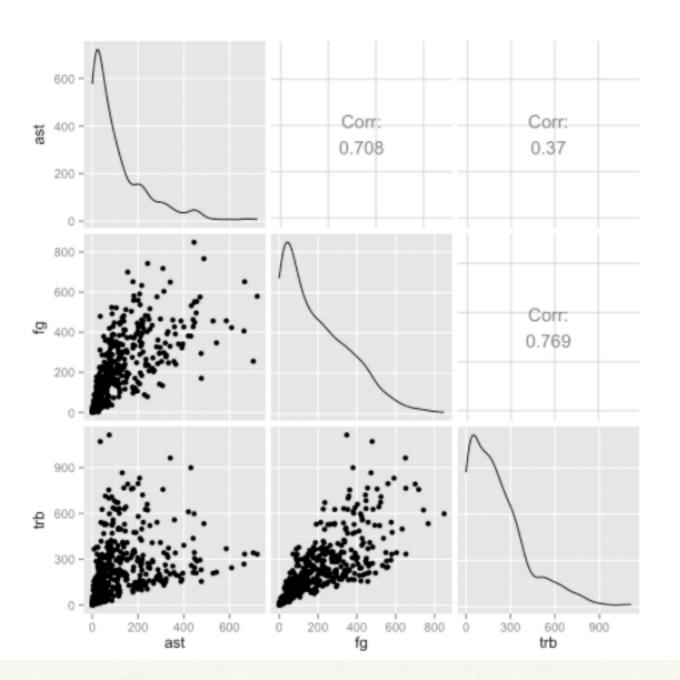
ペアプロットを描く。なんか、データが1つおかしそう・・・

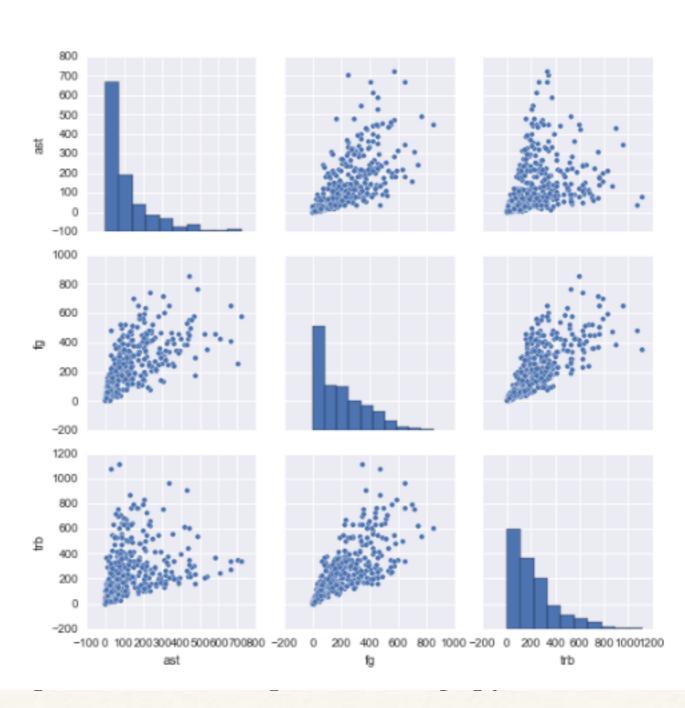
```
R
```

```
library(GGally)
ggpairs(nba[,c("ast", "fg", "trb")])
```

#### Python

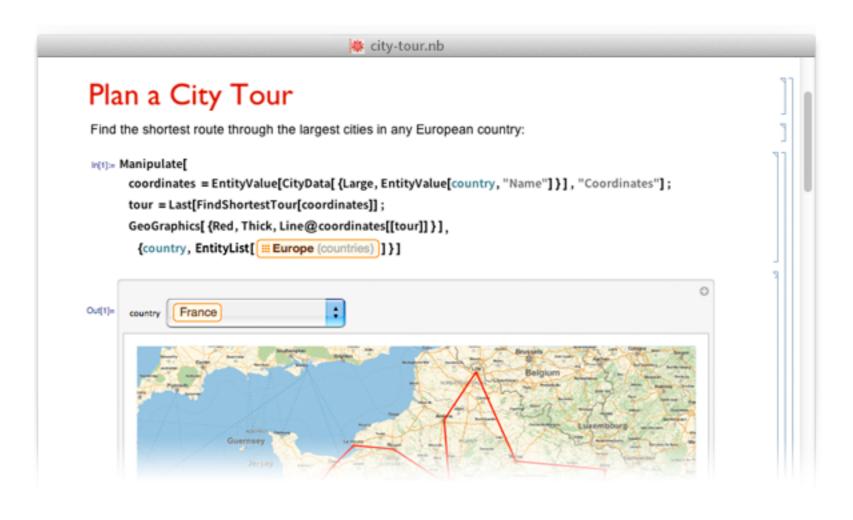
```
import seaborn as sns
import matplotlib.pyplot as plt
sns.pairplot(nba[["ast", "fg", "trb"]])
plt.show()
```





### 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せることができ、デスクトップでもWebでも使えます:



#### 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。

## Pythonが台風の目に

\* Pythonとそれを支えるデータサイエンスのエコシステムは、良いものをどんどん取り入れ進化が加速

- \* もちろん、それぞれの特技はある
  - \* たとえば、Rは統計に強い

### マニアックな統計関数

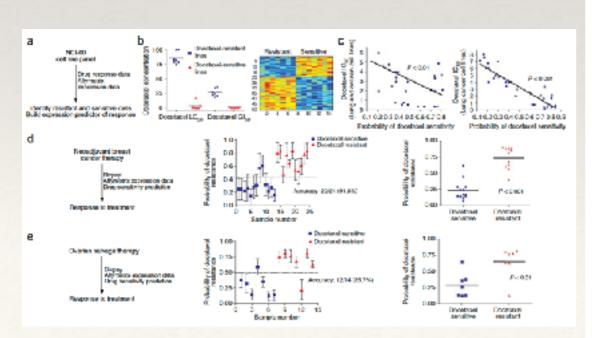
	R	Python
Kolmogorov-Smirnov検定	0	scipy.stats
Shapiro-Wilk検定	0	scipy.stats
Hosmer-Lemeshow検定	ResourceSelection	?

適切なツールを選べる知識と柔軟性が重要 (だと思う)

## サイエンスと再現性

## 本当にあった怖い話

- \* 2006年、名門Nature Medicine誌に、1本の論文が掲載される
- \* 抗がん剤の効きを、細胞株 (試験管で培養できる細胞) のデータを使って予測できるとする当時としては画期的な内容
- \* 教授にこの論文の再現を命じられる私



medicine

### Genomic signatures to guide the use of chemotherapeutics

Anil Potti<sup>1,2</sup>, Holly K Dressman<sup>1,3</sup>, Andrea Bild<sup>1,3</sup>, Richard F Riedel<sup>1,2</sup>, Gina Chan<sup>4</sup>, Robyn Sayer<sup>4</sup>, Janiel Cragun<sup>4</sup>, Hope Cottrill<sup>4</sup>, Michael J Kelley<sup>2</sup>, Rebecca Petersen<sup>5</sup>, David Harpole<sup>5</sup>, Jeffrey Marks<sup>5</sup>, Andrew Berchuck<sup>1,6</sup>, Geoffrey S Ginsburg<sup>1,2</sup>, Phillip Febbo<sup>1–3</sup>, Johnathan Lancaster<sup>4</sup> & Joseph R Nevins<sup>1–3</sup>

Using *in vitro* drug sensitivity data coupled with Affymetrix microarray data, we developed gene expression signatures that predict sensitivity to individual chemotherapeutic drugs. Each signature was validated with response data from an independent set of cell line studies. We further show that many of these signatures can accurately predict clinical response in individuals treated with these drugs. Notably, signatures developed to predict response to individual agents, when combined, could also predict response to multidrug regimens. Finally, we integrated the chemotherapy response signatures with signatures of oncogenic pathway descendables to individual agents. The development of cape expression

### なにかおかしい・・・

- \* データ解析だけの論文だが、複雑な計算方法の詳細が書かれていない
- \* 業界の常識的な方法論でいろいろ試してみる
- \* まったく再現しない・・・
- \* 世界中の研究者が疑いはじめる

### 全部うそ

ABOUT TOL IP LOGIN NDIVIDUAL LOGIN SUBSCRIBE **MAILING LIST CONTACT US DOCUMENTS** 







Inside information on cancer research and drug development

texaschildrens.org/cancer

publication date: May. 22, 2015.

#### The Price of Deception: How a Duke Patient was Harmed in Potti's Fraudulent Trials







About Wikipedia Community portal Recent changes Contact page

Tools

What links here

Article Talk

#### Anil Potti

From Wikipedia, the free encyclopedia

Anil Potti is a physician and former Duke oncogenomics. He, along with Joseph Nevi scandal at Duke University[1][2][3] and many retracted. On November 9, 2015, the Offic engaged in research misconduct.[4] Accord with ORI, Potti can continue to perform res year 2020, while he "neither admits nor de

Contents [hide]

- 1 Biography
- 2 Scientific misconduct
- 3 Consequences
- 4 Research questions and Institute of Medi
- 5 Continuing controversy
- 6 Retracted papers

### 再現性とオープンなこと

- \* データサイエンスと言うからには、再現性が必要
  - \* もちろん、ビジネスの現場でも
- \*解析を完全に再現するには、オープンな基盤は重要
  - \* Pythonが果たす役割
- \* オープンサイエンスという試みが世界的に動き出している
  - \* STAP細胞ってあるんですかね?

### なにがどこまで出来るのか?

## 自然言語処理

- \* 人が使う言語を計算機に理解させるための技術
- \* 人工知能研究の黎明期(1950年代)からある
- \* そんな簡単じゃない
  - 1.職場の上司とW不倫中です。
  - 2.職場の上司と付き合っています。私は既婚者です。

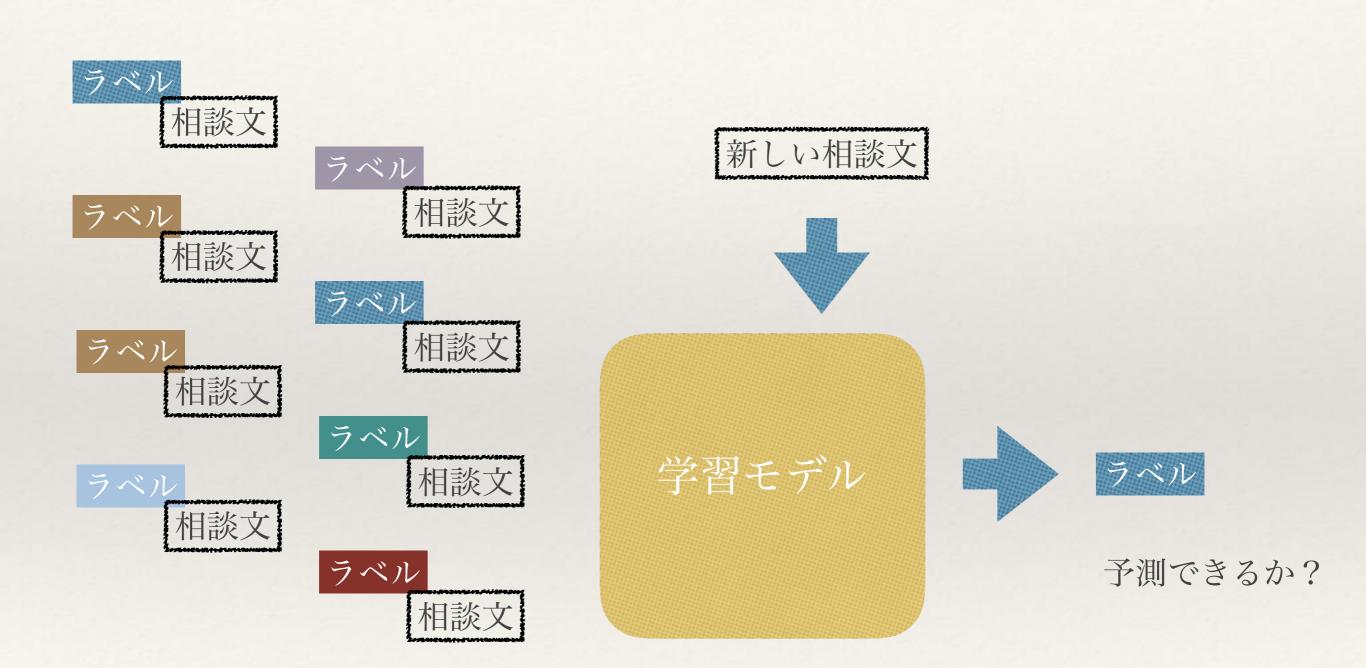
\*占いコンテンツに強い会社さんとの協業事例なので、例文が若干人間味溢れておりますが、ご了承下さい。

### たとえば・・・

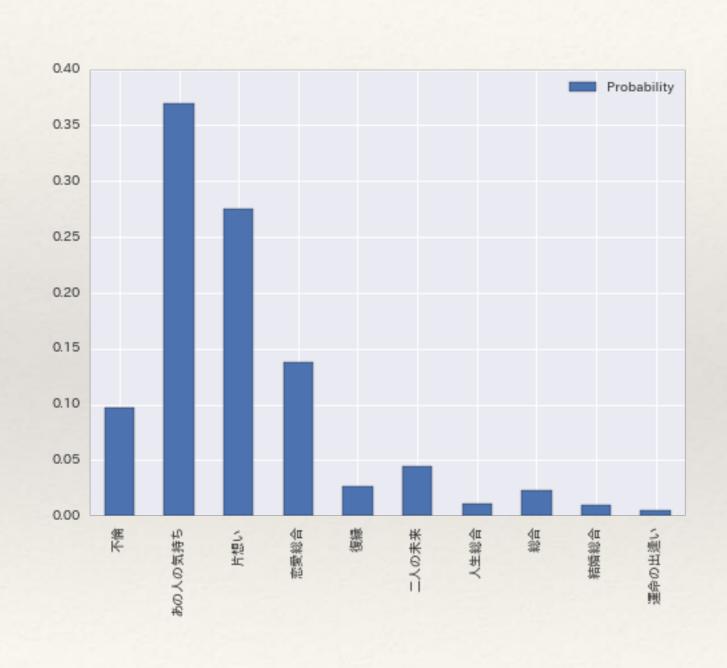
実際の相談内容です。公開出来ないのですみませんが、割愛

これは「片想い」の相談 (ユーザーの入力)

## 大量のデータを使った学習



### 意味を漠然と捉えることはできる

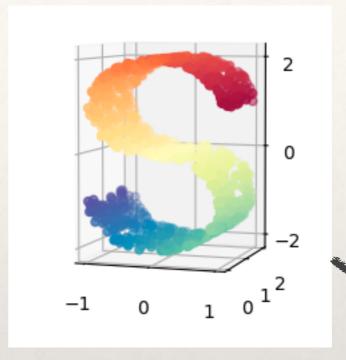


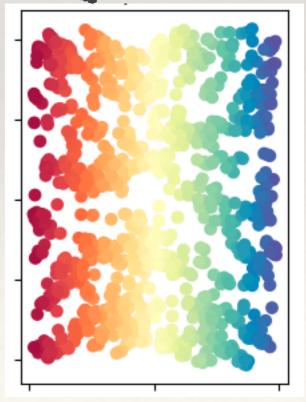
10カテゴリ、3万件の教師データを 使ったRandom Forestsモデルの予測

あの人の気持ちが気になる片想い。 ちょっと不倫の要素も。

## 多様体仮説 (学習)

- \* 単語の数は沢山あるけど、トピックの数は、はるかに少ない
- \* 高次元データでも、それを生成 する本質は低次元で表現でき る(はず)
- \* 計算してデータを加工するのは 計算機、推測するのは人間
  - \* PCA, Autoencoder





### まとめ

- \* データサイエンスはデータ駆動型サイエンス
- \* 前処理負担が重いので、汎用言語が役に立つ
- \* 再現性は重要。オープンな環境がそれを支える
- \* 大量データの背後に何があるのか意識する

### ご静聴ありがとうございました