



university of
groningen



umcg

A python pipeline for data analysis in psychology

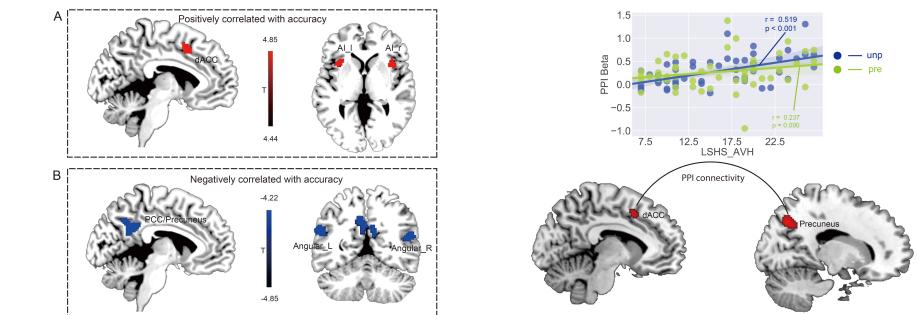
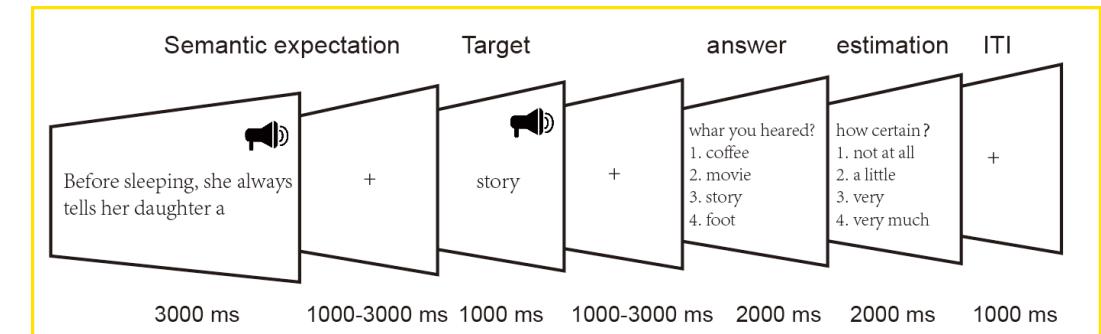
Haiyang Geng

15th Jan 2020

@CNC, UMCG, Groningen

What types of data we analyze by using this python pipeline

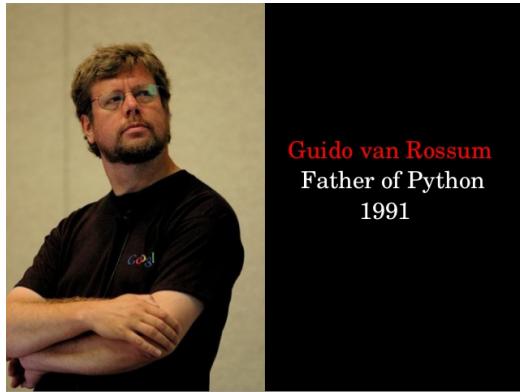
- python pipeline
1. Reaction times and response
 2. Beta values of brain activation and connectivity
 3. Questionnaires and demographic data



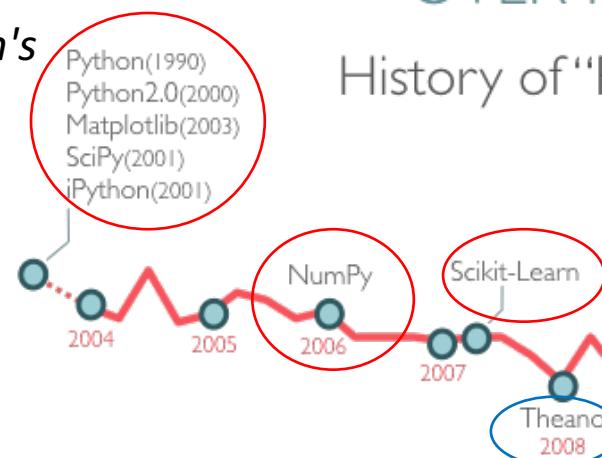


(PhD) **Life is short,
you need python (in psychology)**
Bruce Eckel (Haiyang Geng)

A brief history of python

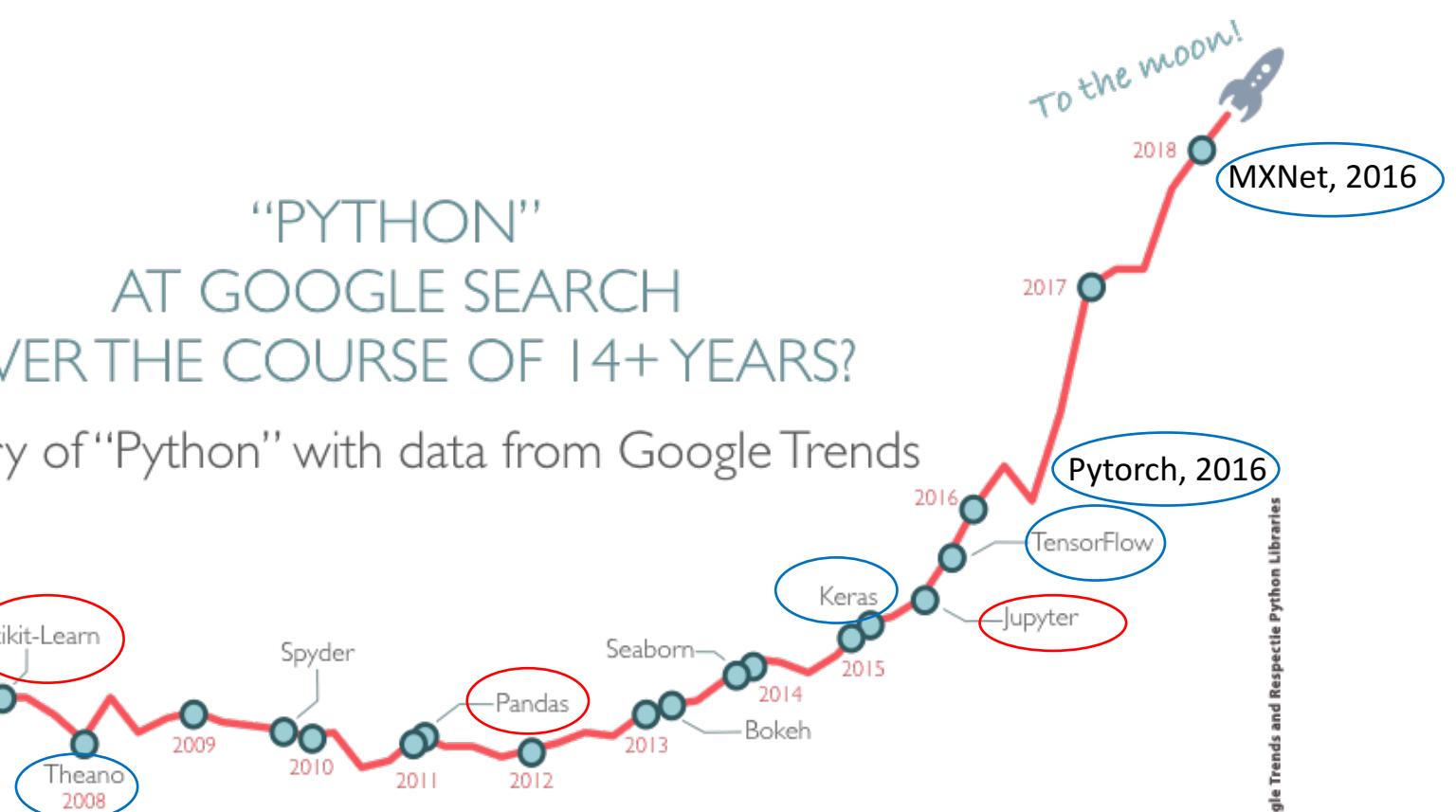


Monty Python's Flying Circus



“PYTHON”
AT GOOGLE SEARCH
OVER THE COURSE OF 14+ YEARS?

History of “Python” with data from Google Trends



Data science in python: <https://www.youtube.com/watch?v=DifMYH3iuFw>

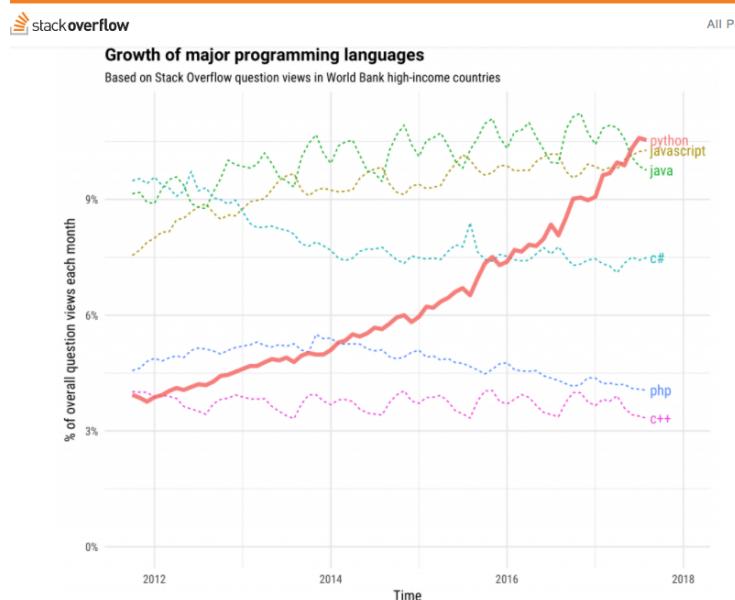
@atillaguzel

Source: Google Trends and Respective Python Libraries

Why do we need python?

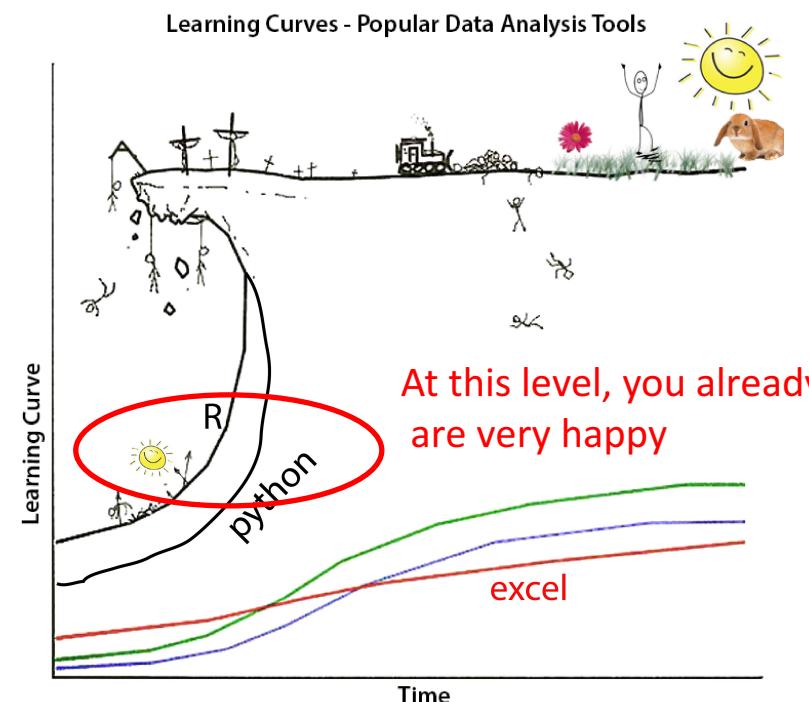
1. Big and active community:

- Abundant open source codes
- Easy to solve problems



2. Python is easier to start:

- Not spend much time to learn
- Save you lots of time



<https://nceas.github.io/oss-lessons/spatial-data-gis-law/1-mon-spatial-data-intro.html>

3. Python is friendly for version control:

- Excel is a mess for recording history
- Everything can be written in python code
 - (tracked by Git)



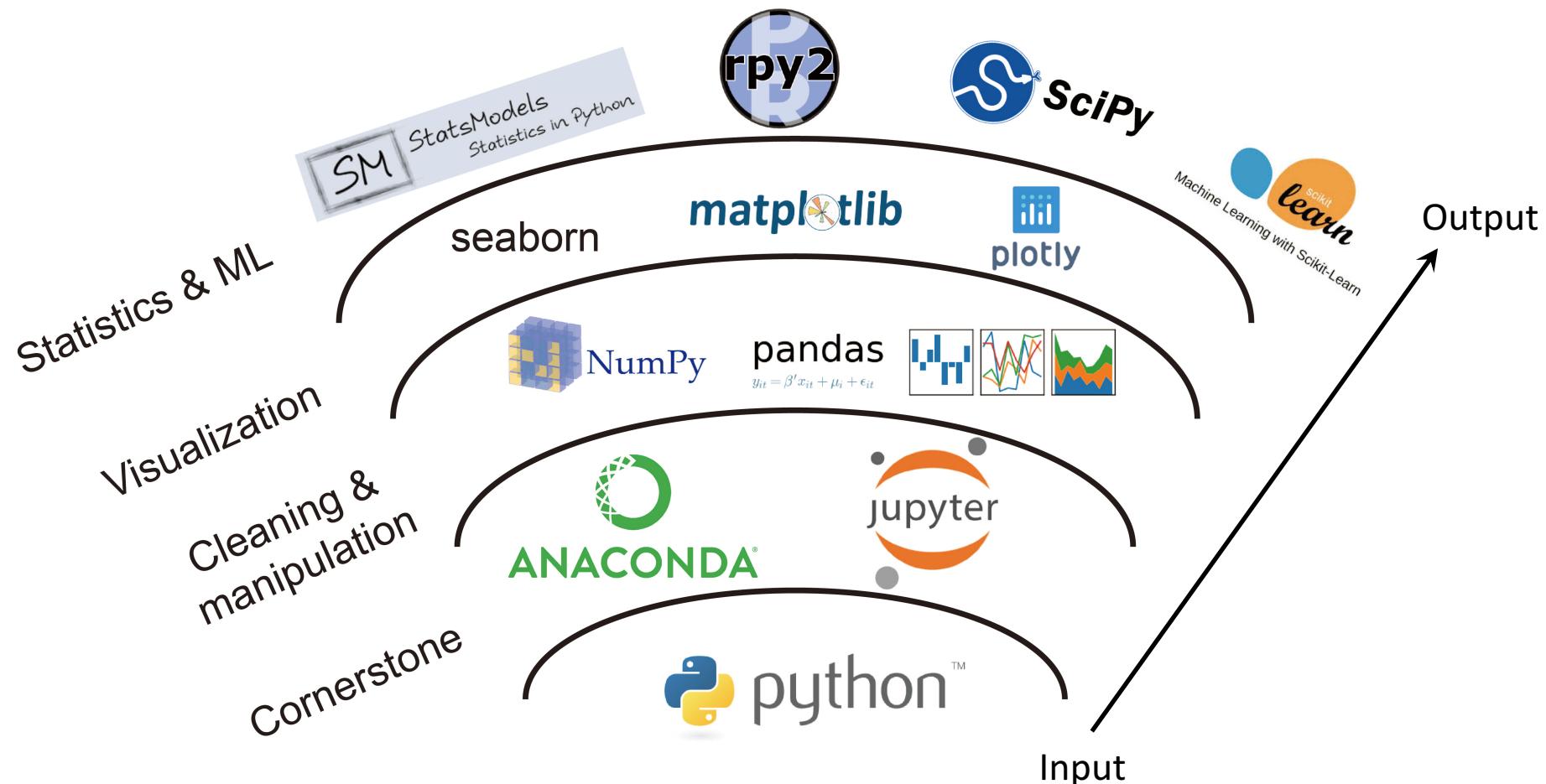
```
1 # combine with questionnaire data
2 se_scale = pd.read_excel('SE_scale.xlsx')
3 se_scale = se_scale.set_index('Subject')
4 #print(se_scale)
5 se_scale = se_scale.set_index(['Subject'])
6 se_scale = se_scale.iloc[:,7:76]
7 opt_all = opt_all.join(se_scale)
8 se_qpe = pd.read_excel('SE-QPE_40.xlsx')
9 se_qpe = se_qpe.set_index('Subject')
10 QPE_severity = pd.DataFrame(se_qpe.iloc[:,12:16].sum(axis = 1),columns = ['QPE_severity'])
11 opt_all = opt_all.join(QPE_severity)
```

4. Excellent ecosystem (packages) for data science

- Numpy, pandas, scipy and scikit-learn...

Why do we need python?

- Excellent ecosystem (packages) for data science

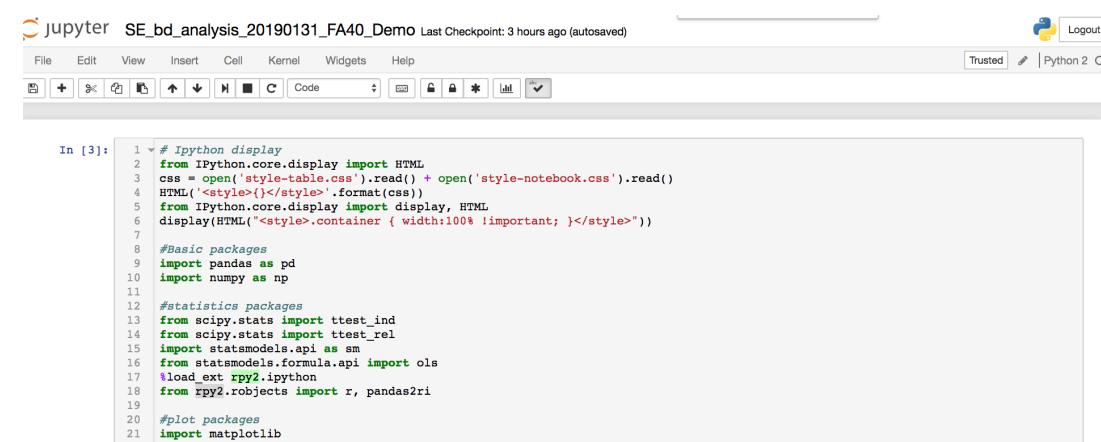
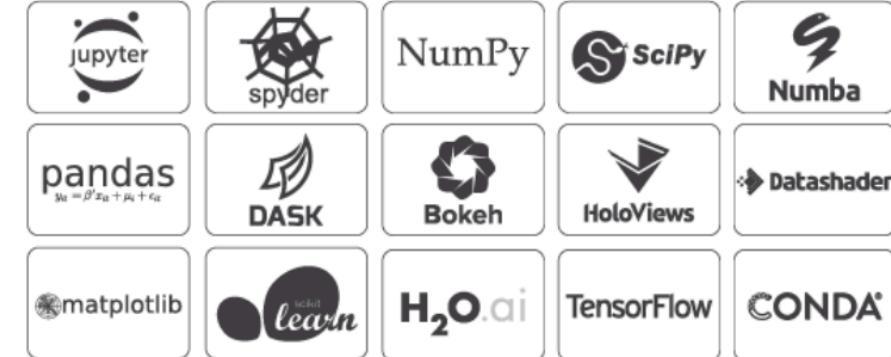


Cornerstone in python ecosystem

conda create
conda install



The World's Most Popular Python/R Data Science Platform



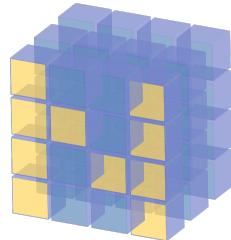
A screenshot of a Jupyter Notebook interface. The title bar shows "jupyter SE_bd_analysis_20190131_FA40_Demo Last Checkpoint: 3 hours ago (autosaved)". The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help. The toolbar has icons for file operations like new, open, save, etc. A code cell titled "In [3]" contains the following Python code:

```
1 # In[3]:  
2 # IPython.core.display import HTML  
3 css = open('style-table.css').read() + open('style-notebook.css').read()  
4 HTML(<style>{}</style>'.format(css))  
5 from IPython.core.display import display, HTML  
6 display(HTML("<style>.container { width:100% !important; }</style>"))  
7  
8 #Basic packages  
9 import pandas as pd  
10 import numpy as np  
11  
12 #statistics packages  
13 from scipy.stats import ttest_ind  
14 from scipy.stats import ttest_rel  
15 import statsmodels.api as sm  
16 from statsmodels.formula.api import ols  
17 %load_ext rpy2.ipython  
18 from rpy2.robjects import r, pandas2ri  
19  
20 #plot packages  
21 import matplotlib  
22 import matplotlib.pyplot as plt
```

An open-source, interactive, light and web browser-based python GUI

Data manipulation and cleaning

For matrix (Matlab)

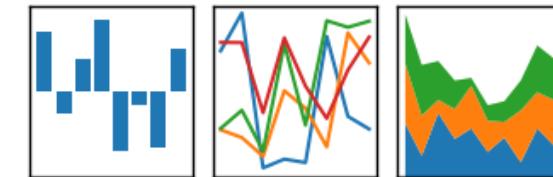


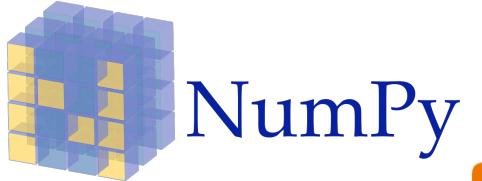
NumPy

For table (R)

pandas

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





NumPy

For matrix data
Neuroimaging data

Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science Interactively at www.DataCamp.com



NumPy

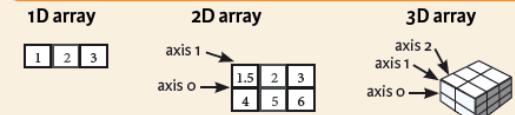
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```



NumPy Arrays



Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]),
      dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3,4))
>>> np.ones((2,3),dtype=np.int16)
>>> d = np.arange(10,25,5)
>>> np.linspace(0,2,9)
>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2x2 identity matrix
Create an array with random values
Create an empty array

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savetxt('array.npy', a, b)
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("myfile.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter=" ")
```

Data Types

<code>>>> np.int64</code>	Signed 64-bit integer types
<code>>>> np.float32</code>	Standard double-precision floating point
<code>>>> np.complex</code>	Complex numbers represented by 128 floats
<code>>>> np.bool</code>	Boolean type storing TRUE and FALSE values
<code>>>> np.object</code>	Python object type
<code>>>> np.string_</code>	Fixed-length string type
<code>>>> np.unicode_</code>	Fixed-length unicode type

Inspecting Your Array

<code>>>> a.shape</code>	Array dimensions
<code>>>> len(a)</code>	Length of array
<code>>>> b.ndim</code>	Number of array dimensions
<code>>>> e.size</code>	Number of array elements
<code>>>> b.dtype</code>	Data type of array elements
<code>>>> b.dtype.name</code>	Name of data type
<code>>>> b.astype(int)</code>	Convert an array to a different type

Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

Array Mathematics

Arithmetic Operations

<code>>>> g = a - b</code>	Subtraction
<code>>>> np.subtract(a,b)</code>	Subtraction
<code>>>> b + a</code>	Addition
<code>>>> np.add(b,a)</code>	Addition
<code>>>> a / b</code>	Division
<code>>>> np.divide(a,b)</code>	Division
<code>>>> a * b</code>	Multiplication
<code>>>> np.multiply(a,b)</code>	Exponentiation
<code>>>> np.exp(b)</code>	Square root
<code>>>> np.sqrt(b)</code>	Print sines of an array
<code>>>> np.sin(a)</code>	Element-wise cosine
<code>>>> np.cos(b)</code>	Element-wise natural logarithm
<code>>>> np.log(a)</code>	Dot product
<code>>>> e.dot(f)</code>	
<code>>>> array([[1., 2., 3.], [4., 5., 6.]])</code>	

Comparison

<code>>>> a == b</code>	Element-wise comparison
<code>>>> array([[False, True, True], [False, False, False]], dtype=bool)</code>	
<code>>>> a < 2</code>	Element-wise comparison
<code>>>> array[[True, False, False], dtype=bool]</code>	
<code>>>> np.array_equal(a, b)</code>	Array-wise comparison

Aggregate Functions

<code>>>> a.sum()</code>	Array-wise sum
<code>>>> a.min()</code>	Array-wise minimum value
<code>>>> b.max(axis=0)</code>	Maximum value of an array row
<code>>>> b.cumsum(axis=1)</code>	Cumulative sum of the elements
<code>>>> a.mean()</code>	Mean
<code>>>> b.median()</code>	Median
<code>>>> a.corrcoef()</code>	Correlation coefficient
<code>>>> np.std(b)</code>	Standard deviation

Copying Arrays

<code>>>> h = a.view()</code>	Create a view of the array with the same data
<code>>>> np.copy(a)</code>	Create a copy of the array
<code>>>> h = a.copy()</code>	Create a deep copy of the array

Sorting Arrays

<code>>>> a.sort()</code>	Sort an array
<code>>>> c.sort(axis=0)</code>	Sort the elements of an array's axis

Subsetting, Slicing, Indexing

Also see Lists

Subsetting

```
>>> a[2]
3
>>> b[1,2]
6.0
```

Select the element at the 2nd index

Slicing

```
>>> a[0:2]
array([1, 2, 3])
>>> b[0:2,1]
array([ 1.,  5.])
>>> b[1,:]
array([1.5, 2., 3.])
```

Select items at index 0 and 1

Boolean Indexing

```
>>> a[a<2]
array([1, 2, 3])
```

Select elements from a less than 2

Fancy Indexing

```
>>> b[[1, 0, 1, 0], [0, 1, 2, 1]]
array([ 1.,  4.,  5.,  6.])
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]
array([ 1.,  4.,  5.,  6.,  1.5,  2.,  3.,  1.5])
```

Select elements (1,0),(0,1),(1,2) and (0,0)
Select a subset of the matrix's rows and columns

Array Manipulation

Transposing Array

```
>>> i = np.transpose(b)
>>> i.T
```

Permute array dimensions
Permute array dimensions

Changing Array Shape

```
>>> b.ravel()
>>> g.reshape(3,-2)
```

Flatten the array
Reshape, but don't change data

Adding/Removing Elements

```
>>> h.resize(2,6)
>>> np.append(h,g)
>>> np.insert(a, 1, 5)
>>> np.delete(a,[1])
```

Return a new array with shape (2,6)
Append items to an array
Insert items in an array
Delete items from an array

Combining Arrays

```
>>> np.concatenate((a,d),axis=0)
array([ 1.,  2.,  3., 10., 15., 20.])
>>> np.vstack((a,b))
array([[ 1.,  2.,  3.], [ 1.5,  2.,  3.], [ 4.,  5.,  6.]])
>>> np.r_[e,f]
>>> np.hstack((e,f))
array([ 7.,  7.,  1.,  0.,  7.,  7.,  0.,  1.])
>>> np.column_stack((a,d))
array([[ 1.,  10.], [ 2.,  15.], [ 3.,  20.]])
>>> np.c_[a,d]
```

Concatenate arrays

Stack arrays vertically (row-wise)

Stack arrays vertically (row-wise)
Stack arrays horizontally (column-wise)

Create stacked column-wise arrays

Create stacked column-wise arrays

Splitting Arrays

```
>>> np.hsplit(a,3)
[array([1, 2, 3]), array([4, 5, 6]), array([1.5, 2., 3.])]
>>> np.vsplit(c,2)
[array([[ 1.,  2.,  3.], [ 4.,  5.,  6.]]), array([[ 1.5,  2.,  3.], [ 4.,  5.,  6.]])]
```

Split the array horizontally at the 3rd index
Split the array vertically at the 2nd index



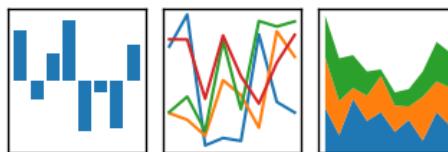


pandas

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

For table data

read_csv
df.head, describe
df.join
df['clm'], df.iloc
pd.concat
groupby



<http://pandas.pydata.org>

Syntax – Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(
    {"a" : [4 ,5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.
```

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

	a	b	c
n			
d	1	4	7
e	2	5	8
	6	9	12

```
df = pd.DataFrame(
    {"a" : [4 ,5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=[ 'n', 'v']))
Create DataFrame with a MultiIndex
```

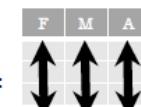
Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={
          'variable' : 'var',
          'value' : 'val'})
      .query('val >= 200')
      )
```

Tidy Data – A foundation for wrangling in pandas

In a tidy data set:



&



Each variable is saved in its own column
Each observation is saved in its own row

Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

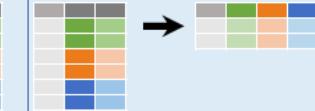


M * A

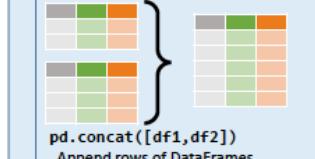
Reshaping Data – Change the layout of a data set



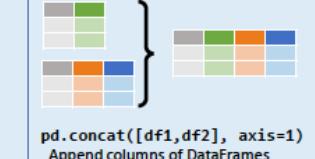
pd.melt(df)
Gather columns into rows.



df.pivot(columns='var', values='val')
Spread rows into columns.



pd.concat([df1, df2])
Append columns of DataFrames



pd.concat([df1, df2], axis=0)
Append rows of DataFrames

df.sort_values('mpg')
Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)
Order rows by values of a column (high to low).

df.rename(columns = { 'y' : 'year'})
Rename the columns of a DataFrame

df.sort_index()
Sort the index of a DataFrame

df.reset_index()
Reset index of DataFrame to row numbers, moving index to columns.

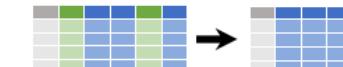
df.drop(columns=['Length', 'Height'])
Drop columns from DataFrame

Subset Observations (Rows)



```
df[df.Length > 7]  
Extract rows that meet logical criteria.  
df.drop_duplicates()  
Remove duplicate rows (only considers columns).  
df.head(n)  
Select first n rows.  
df.tail(n)  
Select last n rows.
```

Subset Variables (Columns)



```
df[['width', 'length', 'species']]  
Select multiple columns with specific names.  
df['width'] or df.width  
Select single column with specific name.  
df.filter(regex='regex')  
Select columns whose name matches regular expression regex.
```

regex (Regular Expressions) Examples

'.'	Matches strings containing a period ''.
'Length\$'	Matches strings ending with word 'Length'
'^Sepal'	Matches strings beginning with the word 'Sepal'
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^?!Species\$.+'	Matches strings except the string 'Species'

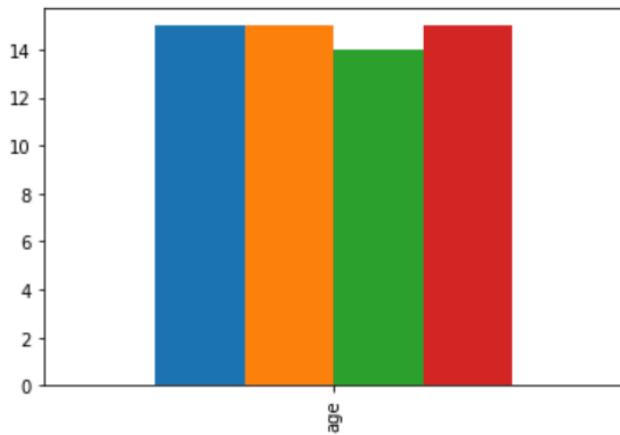
```
df.loc[:, 'x2' : 'x4']  
Select all columns between x2 and x4 (inclusive).  
df.iloc[:, [1,2,5]]  
Select columns in positions 1, 2 and 5 (first column is 0).  
df.loc[df['a'] > 10, ['a', 'c']]  
Select rows meeting logical condition, and only the specific columns .
```

Logic in Python (and pandas)

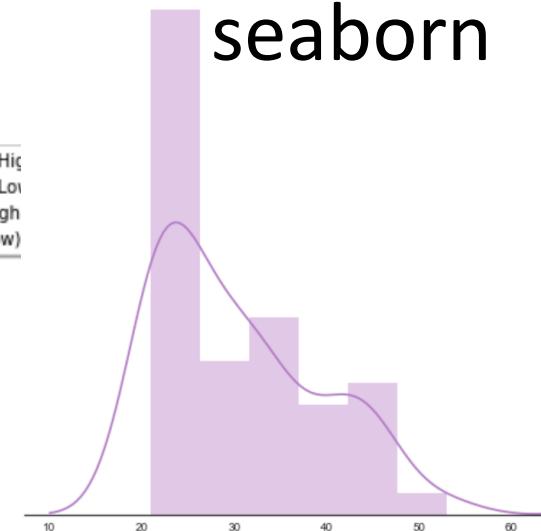
<	Less than	!=	Not equal to
>	Greater than	df.column.isin(values)	Group membership
==	Equals	pd.isnull(obj)	Is NaN
<=	Less than or equals	pd.notnull(obj)	Is not NaN
>=	Greater than or equals	df._.~_, df.any(), df.all()	Logical and, or, not, xor, any, all

Data visualization

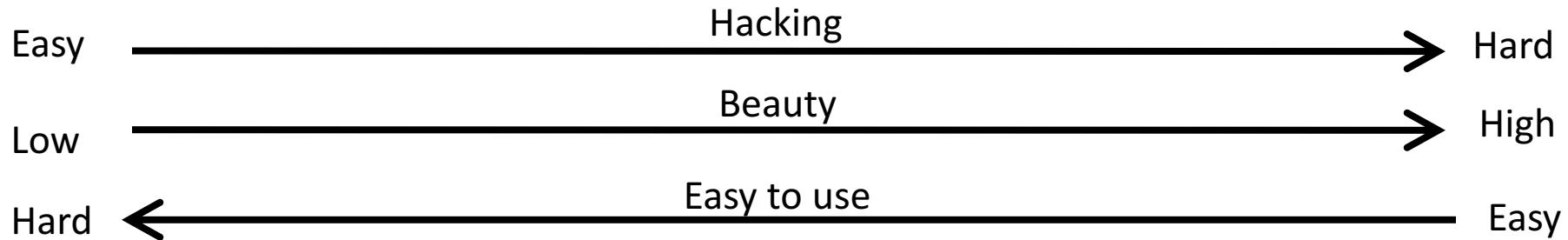
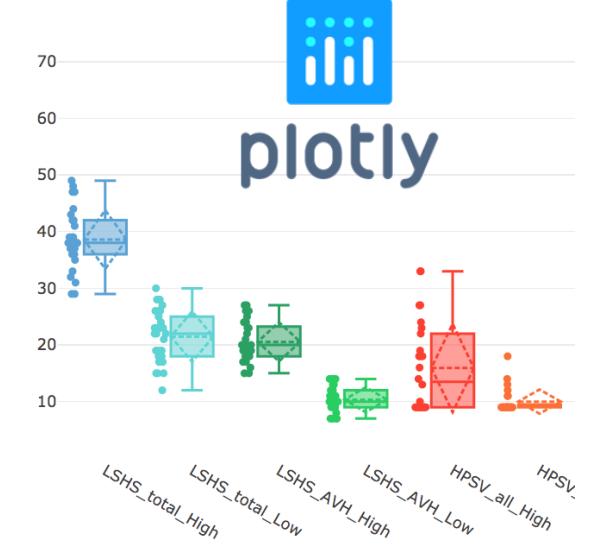
matplotlib



seaborn



plotly



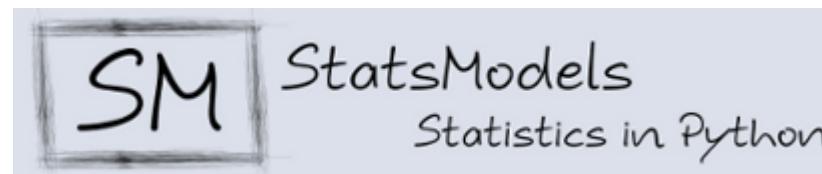
Statistics

Basic statistics (t test, Pearson correlation): [scipy.stats](#)



```
ttest_ind(opt_all_fa_High_unp,opt_all_fa_Low_unp).statistic
```

Advanced statistics (GLM and AVOVA): [statsmodels](#)



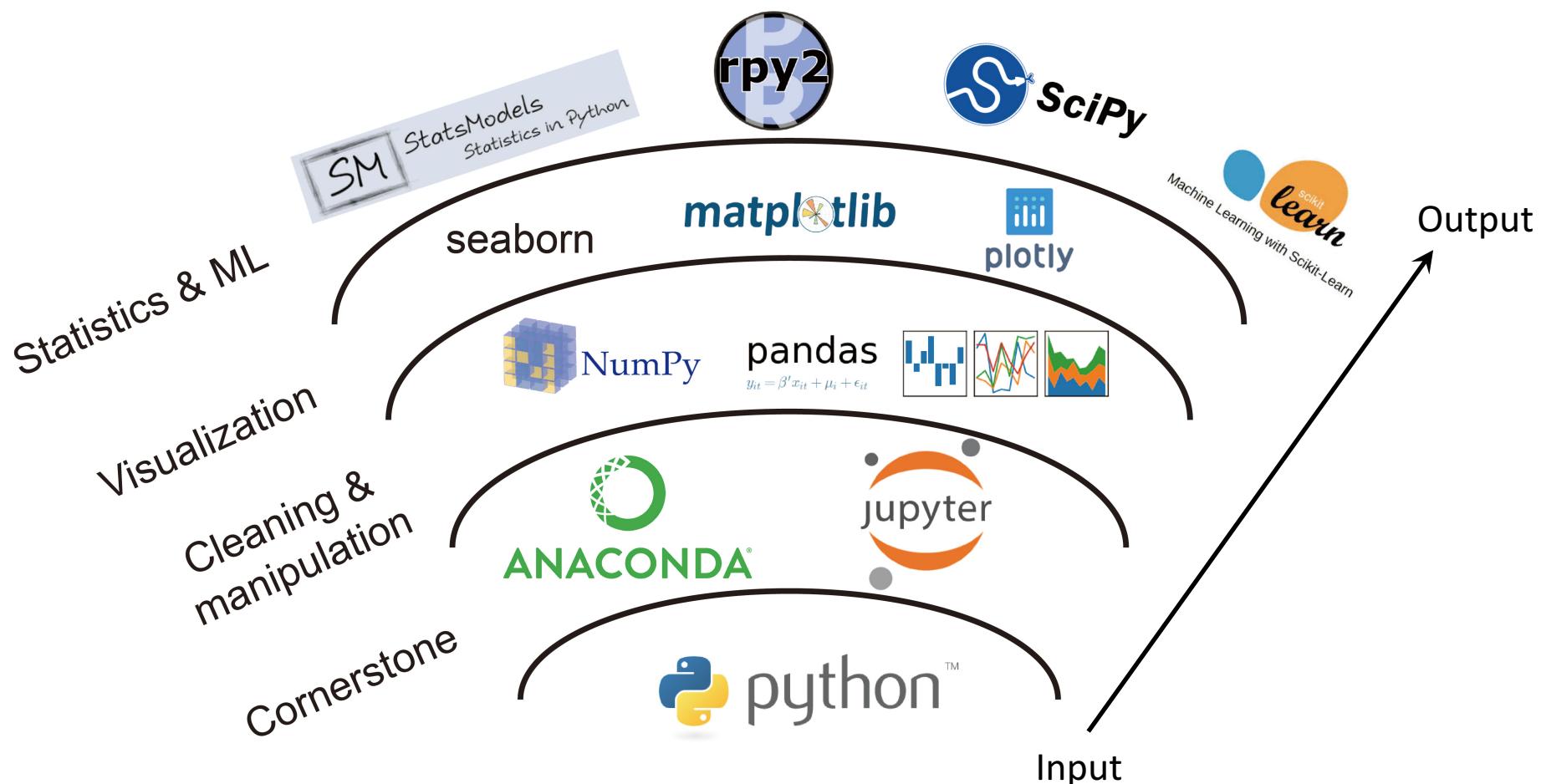
```
results = smf.ols('Lottery ~ Literacy + np.log(Pop1831)',  
data=dat).fit()
```

More advanced statistics (everything from R): [rpy2](#)



```
%R -i r_df_bd_val r_df_bd_val$Subject <-  
factor(r_df_bd_val$Subject);aov_group_conditon <- aov(conf ~  
group*condition + Error(Subject/condition),  
data=r_df_bd_val);print(summary(aov_group_conditon));print(model.t  
ables(aov_group_conditon, "means"))
```

A pipeline for psychological data analysis



Programming saves your time to do more important things

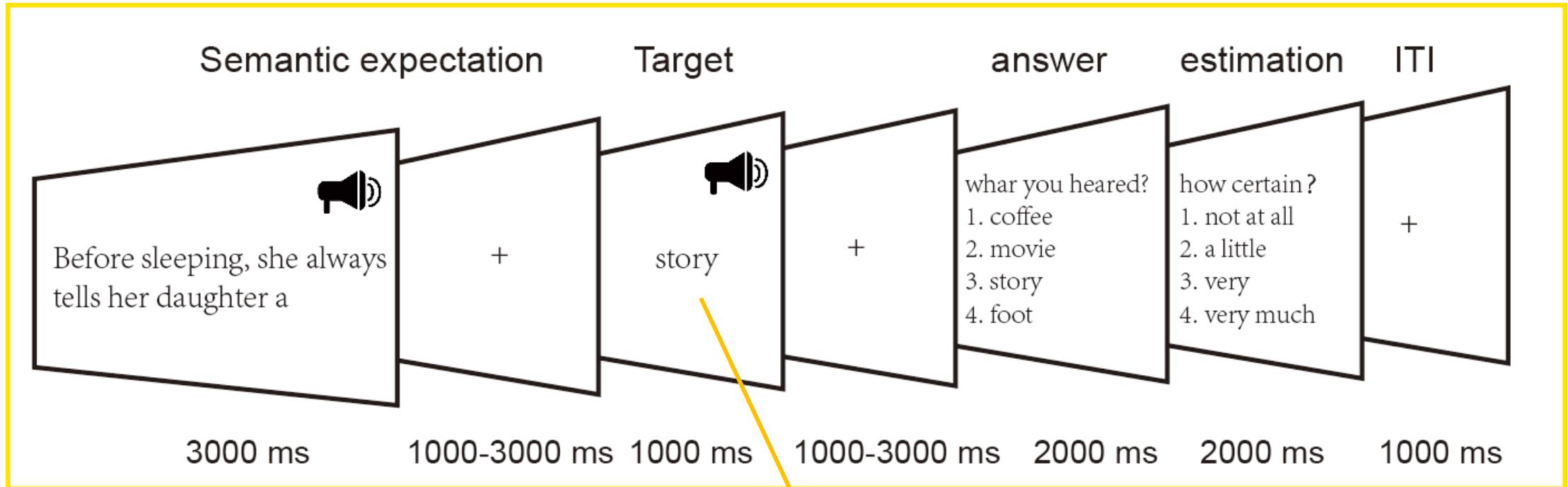
- Standardize your pipeline by programming, make your and other's life easy
- Put all the data and codes into one box, avoid confusion, data missing and mistakes
- Look into the details of your data by visualizing, remove bad (noisy) data to make your data clean



An example: get your hand dirty

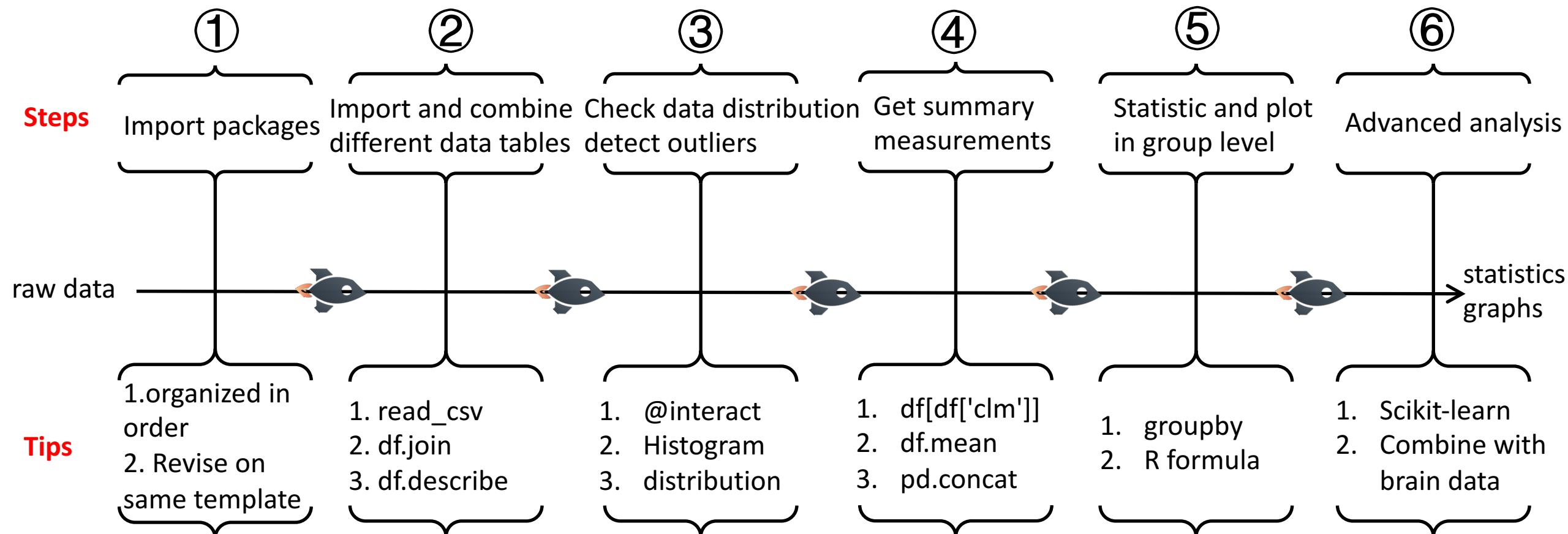
Task design

Semantic Expectation Task



Predictable condition: story, 60 trials in four sessions
Unpredictable condition : coffee, 60 trials in four sessions

Summary for the demo pipeline





Several useful tips (suggestions) for using this pipeline

1. During programming

- Use **Dash** and **Stack Overflow**, which can almost solve all the (programing) questions.
- Use **Anaconda** to manage different **python environments**
- Collect **optimal tools** (packages) to **polish** your pipeline **regularly**

2. During data analysis

- **Check data** before and after manipulation by using visualization and summary statistic
- **Validate** your code by using 'benchmark' software (tool) (spss, jasp, R) when doing complex statistic

3. About algorithms

- Implementation is easy, understanding is hard, put more time on figuring out underlying algorithms



Thanks for your attention