

資料處理與視覺化 - numpy, pandas

2020 / 9 / 21
PRESENTED BY
AI Foundation

CREATED FOR

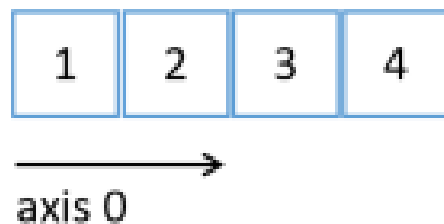


NumPy 簡介

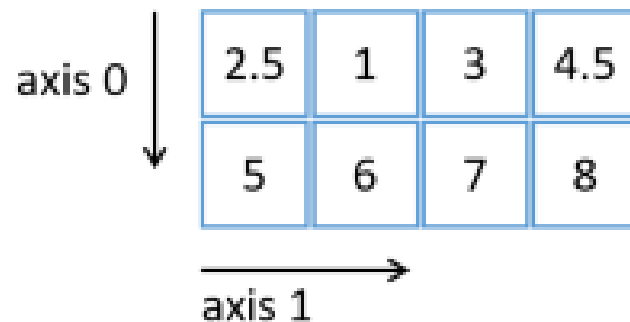


POWERFUL N-DIMENSIONAL ARRAYS

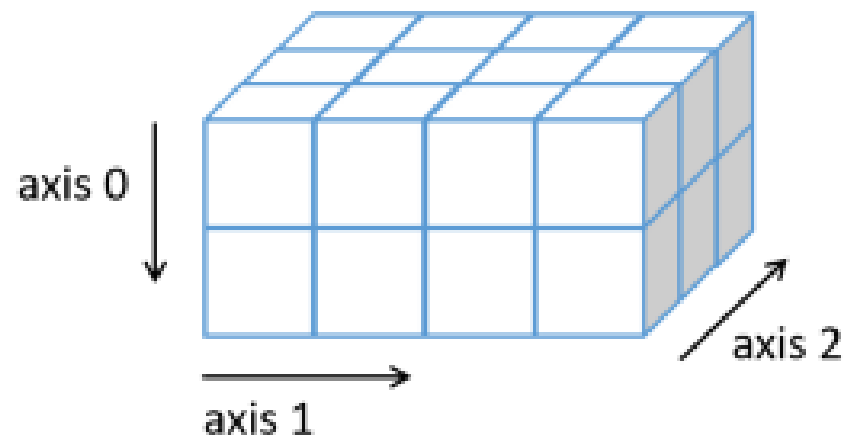
一維陣列(1D array)



二維陣列(2D array)



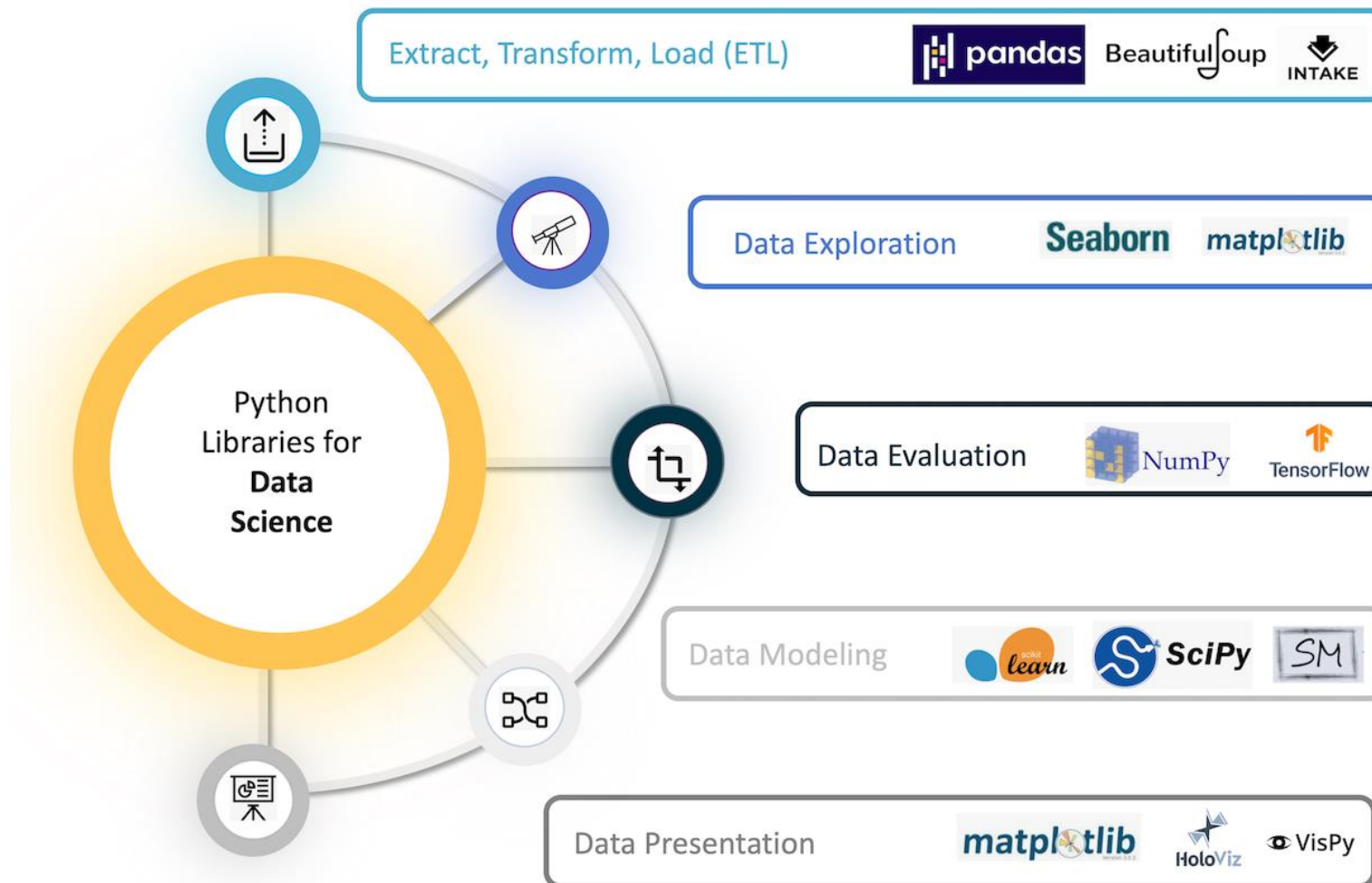
三維陣列(3D array)



Fast and versatile !

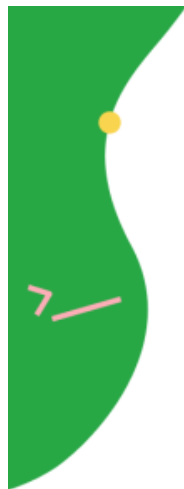
<https://medium.com/python4u/hello-numpy-b5ebe67a1ada>

資料科學生態系



<https://numpy.org/>

各套件在 Github 上被引用的比例



Packages Imported by Machine Learning Projects on GitHub

1	numpy	74%
2	scipy	47%
3	pandas	41%
4	matplotlib	40%
5	scikit-learn	38%
6	six	31%
7	tensorflow	24%
8	requests	23%
9	python-dateutil	22%
10	pytz	21%



<https://venturebeat.com/2019/01/24/github-numpy-and-scipy-are-the-most-popular-packages-for-machine-learning-projects/>

各種使用到 NumPy 的套件

Quantum Computing



QuTiP
PyQuil
Qiskit

Statistical Computing



Pandas
statsmodels
Seaborn

Signal Processing



SciPy
PyWavelets

Image Processing



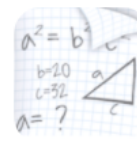
Scikit-image
OpenCV

3-D Visualization



Mayavi
Napari

Symbolic Computing



SymPy

Astronomy Processes



AstroPy
SunPy
SpacePy

Cognitive Psychology



PsychoPy

Bioinformatics



BioPython
Scikit-Bio
PyEnsembl

Bayesian Inference



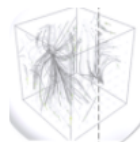
PyStan
PyMC3

Mathematical Analysis



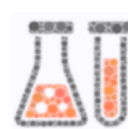
SciPy
SymPy
cvxpy

Simulation Modeling



PyDSTool

Multi-variate Analysis



PyChem

Geographic Processing



Shapely
GeoPandas
Folium

Interactive Computing



Jupyter
IPython
Binder

<https://numpy.org/>

NumPy安裝

CONDA

If you use `conda` , you can install it with:

```
conda install numpy
```

PIP

If you use `pip` , you can install it with:

```
pip install numpy
```

如果你是裝 **Anaconda** ，內建就會有 **NumPy** ，不需要額外安裝！

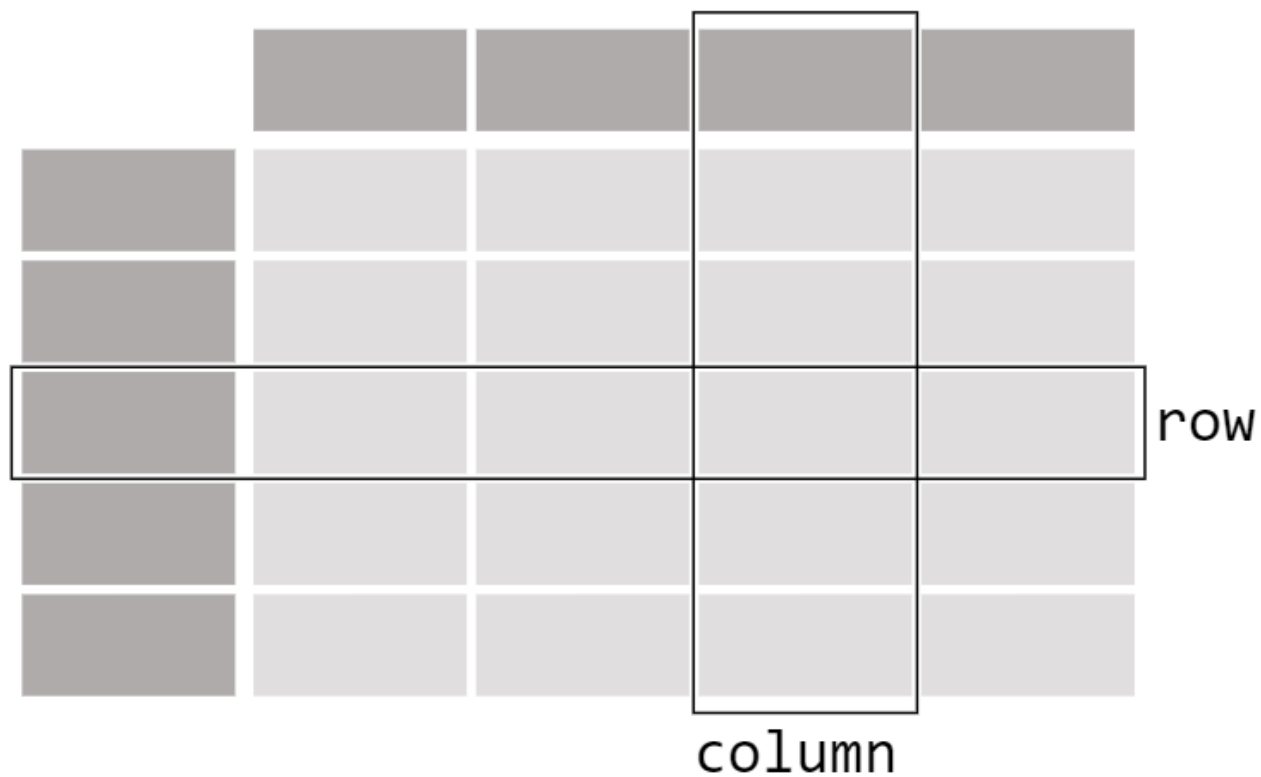
<https://numpy.org/install/>

Pandas 簡介



專門處理表格型資料

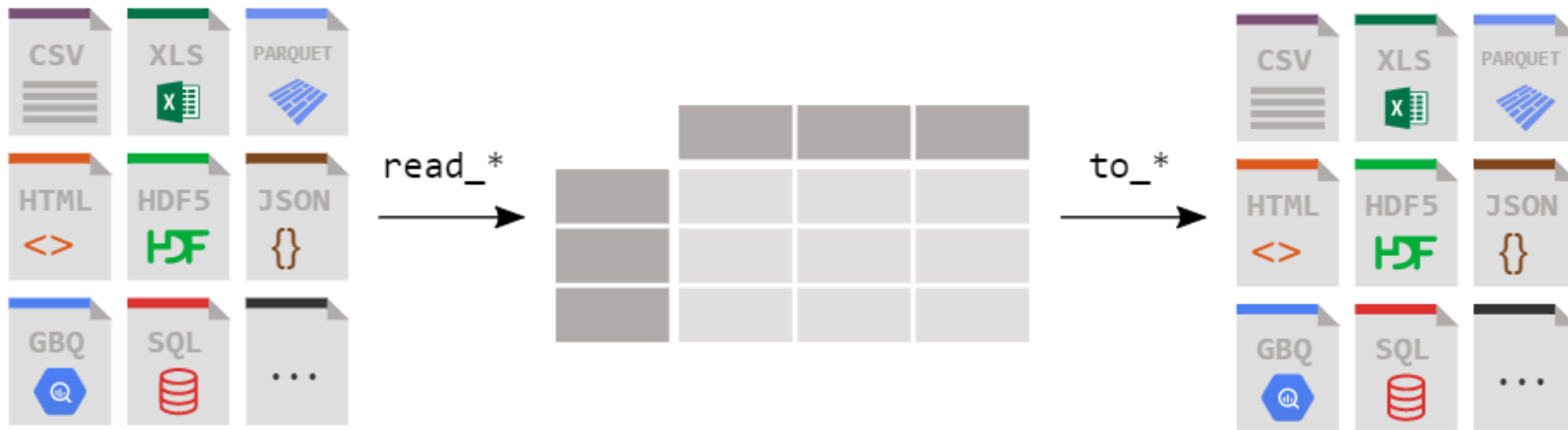
DataFrame



	A	B	C
1	Name	Age	Sex
2	jerry	22	male
3	tom	35	male
4	ellie	58	female

https://pandas.pydata.org/docs/getting_started/intro_tutorials/01_table_oriented.html#min-tut-01-tableoriented

支援各種格式的表格型資料讀寫



https://pandas.pydata.org/docs/getting_started/intro_tutorials/02_read_write.html#min-tut-02-read-write

合併表格



		key

		key		



		key

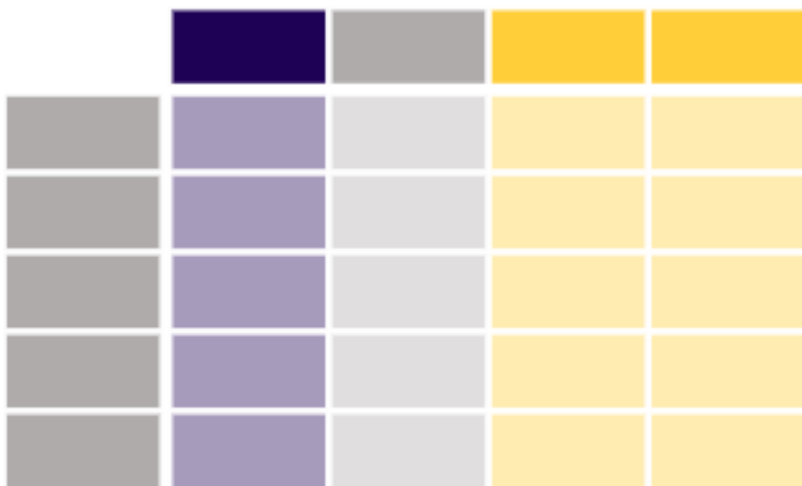
		key		



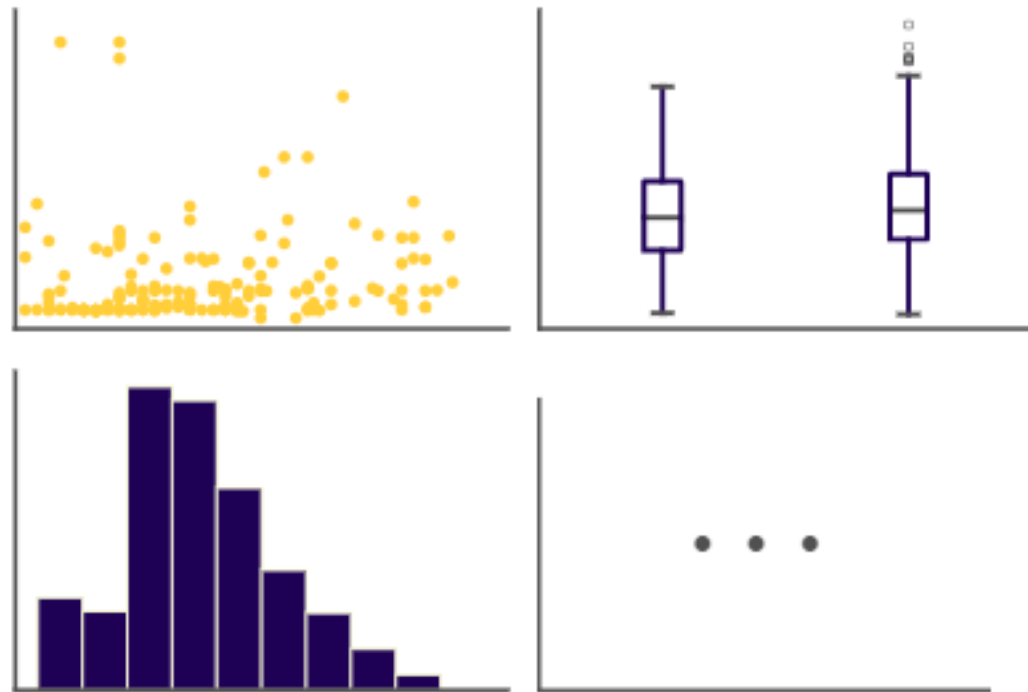
		key		

https://pandas.pydata.org/docs/getting_started/intro_tutorials/08_combine_dataframes.html#min-tut-08-combine

資料視覺化



`.plot.*`

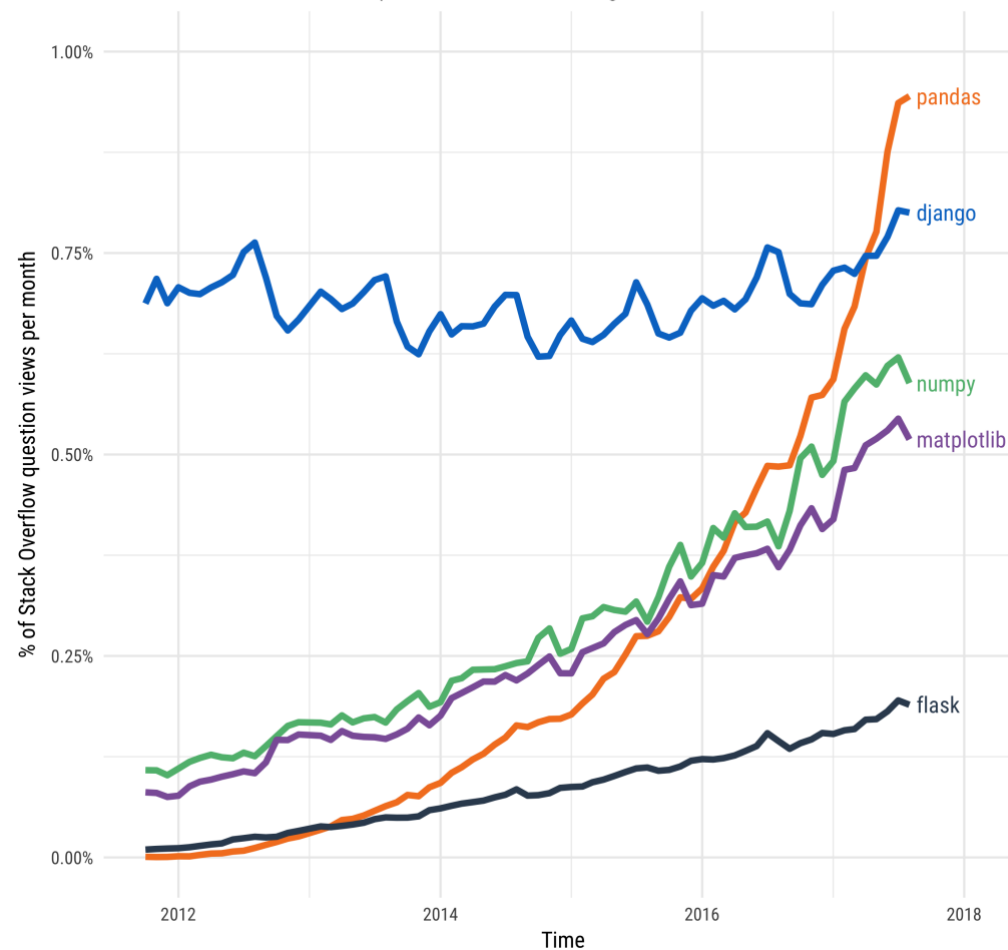


https://pandas.pydata.org/docs/getting_started/intro_tutorials/04_plotting.html#min-tut-04-plotting

在 Stack Overflow 上被討論到的次數

Stack Overflow Traffic to Questions About Selected Python Packages

Based on visits to Stack Overflow questions from World Bank high-income countries



<https://www.sqlshack.com/getting-started-with-pandas-in-python/>

Pandas 安裝

Working with conda?

pandas is part of the [Anaconda](#) distribution and can be installed with Anaconda or Miniconda:

```
conda install pandas
```

Prefer pip?

pandas can be installed via pip from [PyPI](#).

```
pip install pandas
```

如果你是裝 **Anaconda**，內建就會有 **Pandas**，不需要額外安裝！

https://pandas.pydata.org/docs/getting_started/index.html#installation

Pandas Cheat Sheet

Python For Data Science Cheat Sheet

Pandas Basics

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



Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.



Use the following import convention:
>>> import pandas as pd

Pandas Data Structures

Series

A one-dimensional labeled array capable of holding any data type

a	3
b	-5
c	7
d	4

Index

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

DataFrame

A two-dimensional labeled data structure with columns of potentially different types

	Country	Capital	Population
0	Belgium	Brussels	11190846
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],  
           'Capital': ['Brussels', 'New Delhi', 'Brasilia'],  
           'Population': [11190846, 1303171035, 207847528]}
```

```
>>> df = pd.DataFrame(data,  
                      columns=['Country', 'Capital', 'Population'])
```

I/O

Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)  
>>> df.to_csv('myDataFrame.csv')
```

Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')  
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')  
  
Read multiple sheets from the same file  
>>> xls = pd.ExcelFile('file.xls')  
>>> df = pd.read_excel(xls, 'Sheet1')
```

Asking For Help

```
>>> help(pd.Series.loc)
```

Selection

Also see NumPy Arrays

Getting

```
>>> s['b']  
-5  
  
>>> df[1:]  
Country Capital Population  
1 India New Delhi 1303171035  
2 Brazil Brasilia 207847528
```

Get one element

Get subset of a DataFrame

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0], [0]]  
'Belgium'  
  
>>> df.iat([0], [0])  
'Belgium'
```

Select single value by row & column

By Label

```
>>> df.loc[[0], ['Country']]  
'Belgium'  
  
>>> df.at[[0], ['Country']]  
'Belgium'
```

Select single value by row & column labels

By Label/Position

```
>>> df.ix[2]  
Country Brazil  
Capital Brasilia  
Population 207847528
```

Select single row of subset of rows

```
>>> df.ix[:, 'Capital']  
0 Brussels  
1 New Delhi  
2 Brasilia
```

Select a single column of subset of columns

```
>>> df.ix[1, 'Capital']  
'New Delhi'
```

Select rows and columns

Boolean Indexing

```
>>> s[~(s > 1)]  
>>> s[(s < -1) | (s > 2)]  
>>> df[df['Population'] > 1200000000]
```

Series s where value is not >1
s where value is <-1 or >2
Use filter to adjust DataFrame

Setting

```
>>> s['a'] = 6
```

Set index a of Series s to 6

Dropping

```
>>> s.drop(['a', 'c'])  
>>> df.drop('Country', axis=1)
```

Drop values from rows (axis=0)
Drop values from columns (axis=1)

Sort & Rank

```
>>> df.sort_index()  
>>> df.sort_values(by='Country')  
>>> df.rank()
```

Sort by labels along an axis
Sort by the values along an axis
Assign ranks to entries

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape  
>>> df.index  
>>> df.columns  
>>> df.info()  
>>> df.count()
```

(rows, columns)
Describe index
Describe DataFrame columns
Info on DataFrame
Number of non-NA values

Summary

```
>>> df.sum()  
>>> df.cumsum()  
>>> df.min()/df.max()  
>>> df.idxmin()/df.idxmax()  
>>> df.describe()  
>>> df.mean()  
>>> df.median()
```

Sum of values
Cumulative sum of values
Minimum/maximum values
Minimum/Maximum index value
Summary statistics
Mean of values
Median of values

Applying Functions

```
>>> f = lambda x: x*2  
>>> df.apply(f)  
>>> df.applymap(f)
```

Apply function
Apply function element-wise

Data Alignment

Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])  
>>> s + s3  
a 10.0  
b NaN  
c 5.0  
d 7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)  
a 10.0  
b -5.0  
c 5.0  
d 7.0  
  
>>> s.sub(s3, fill_value=2)  
>>> s.div(s3, fill_value=4)  
>>> s.mul(s3, fill_value=3)
```

DataCamp

Learn Python for Data Science Interactively



<http://datacamp-community-prod.s3.amazonaws.com/dbed353d-2757-4617-8206-8767ab379ab3>

來看程式碼吧！