

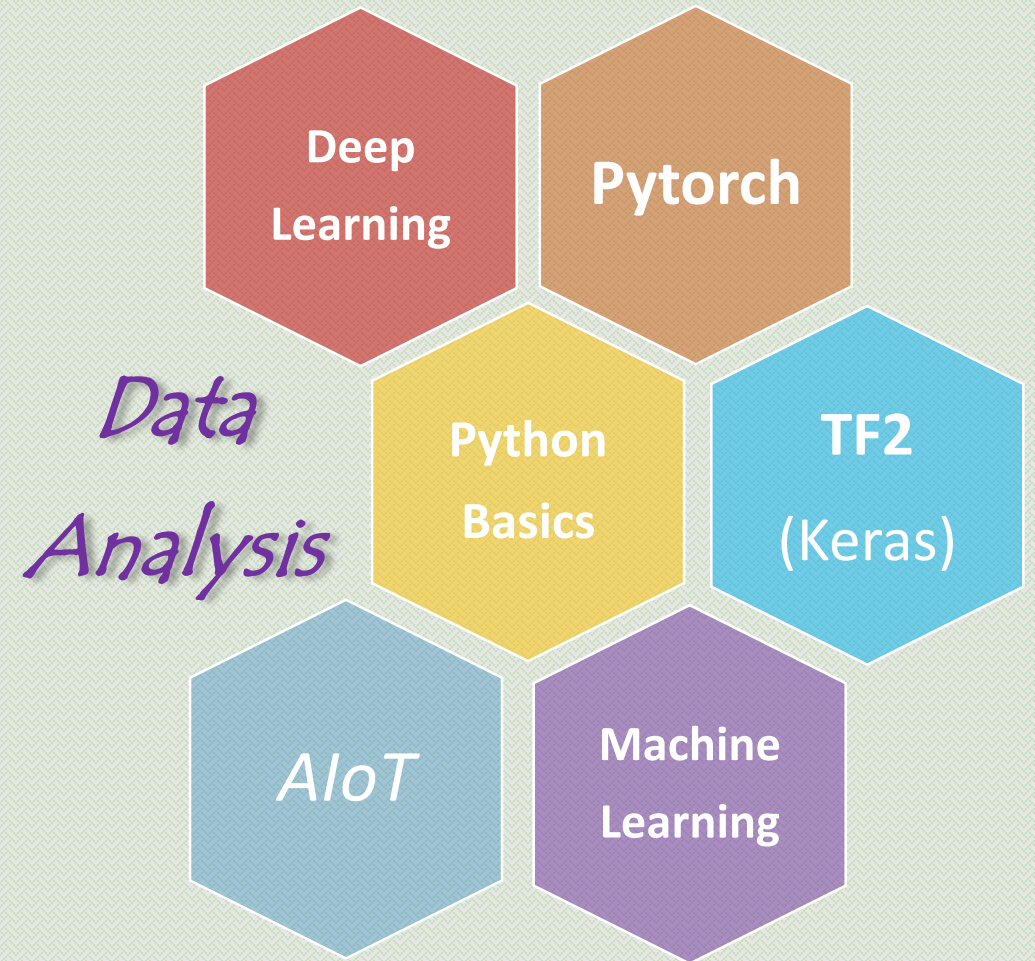
< PART IV >

Feature Engineering

(特徵工程)

胡中興 (C. Alex Hu, PhD)

Courses on “*Python Data Analysis* (Python 資料分析)”

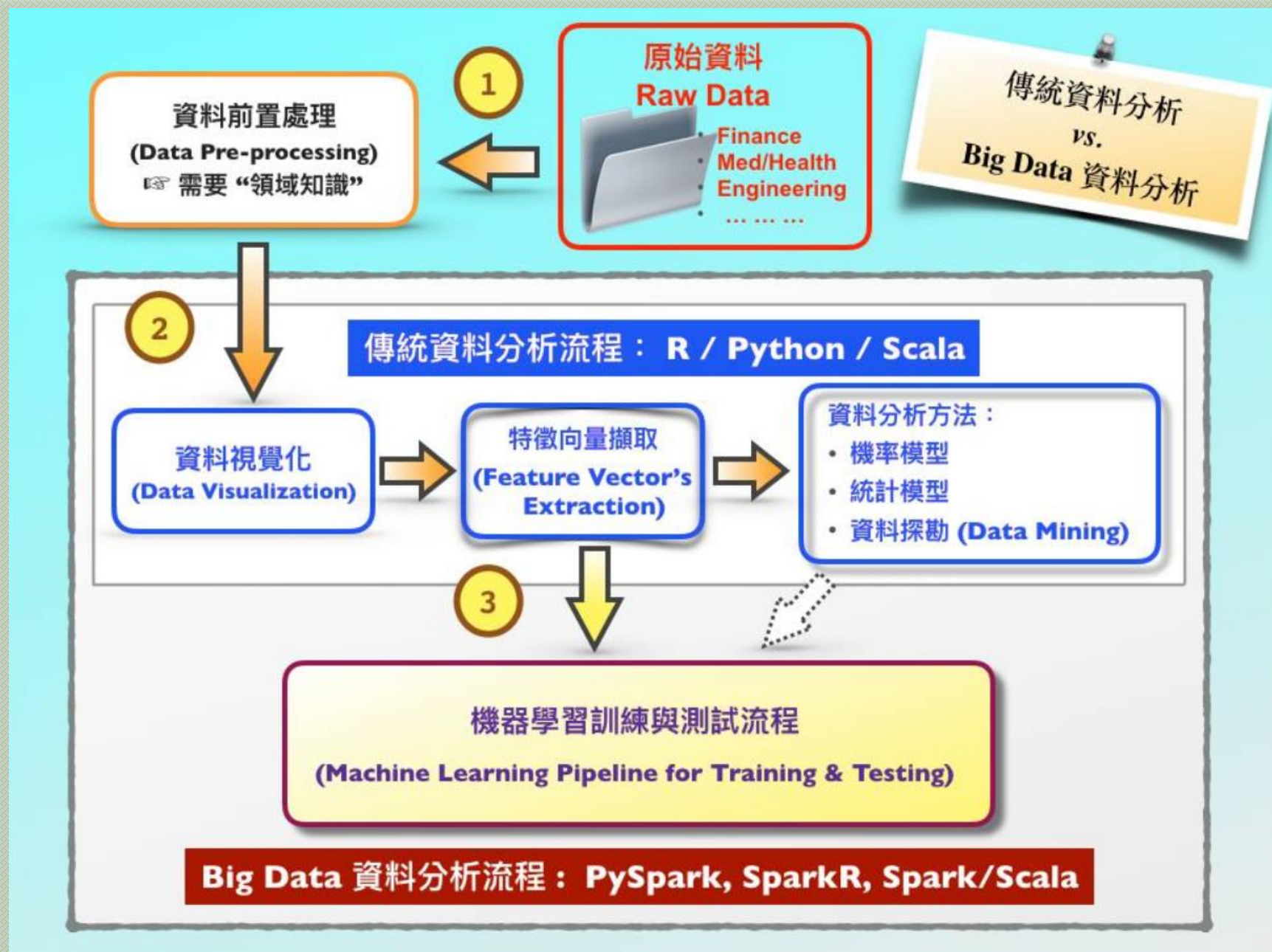


OUTLINE

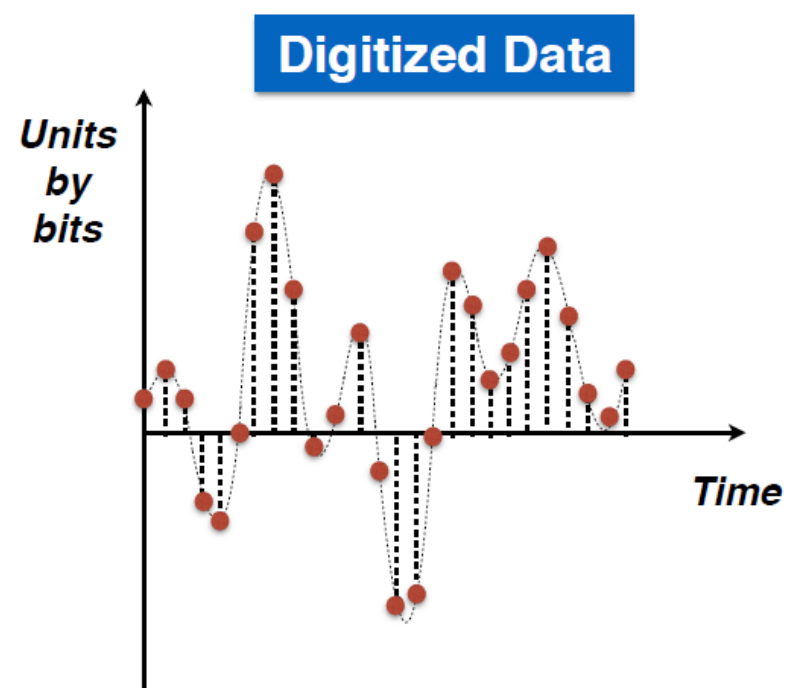
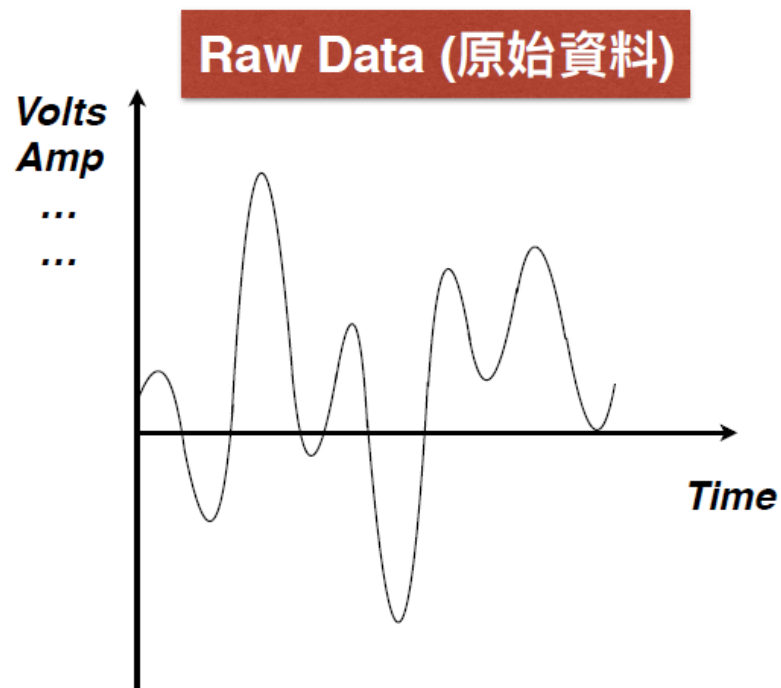
- Data Analysis & Machine Learning Pipeline
- Feature Engineering
 - ❖ Feature Extraction
 - ❖ Feature Selection
- Scikit-Learn Workshop for Feature Engineering

Data Analysis & Machine Learning Pipeline

(資料分析與機器學習流程)



1



Data Acquisition
(數據擷取)

- Sampling Rate (取樣速率)
- Quantization (數據量化)

Analog Signal
(類比訊號)

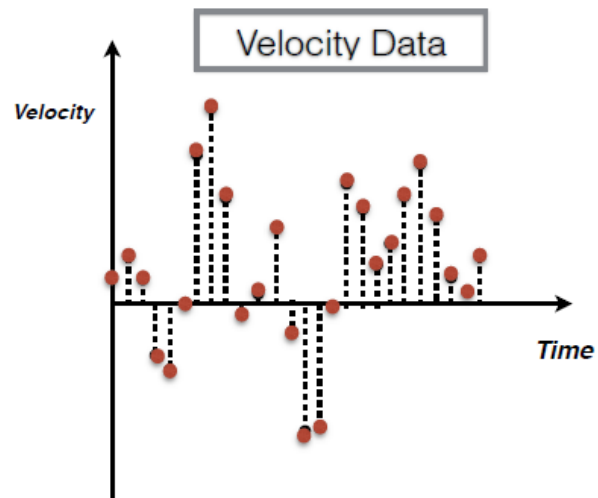


A/D C



Digital Signal
(數位訊號)

Digitized Dataset



無人自駕車、
無人機、
手機定位、
智慧型機器人、
智慧感測 ...

GPS Data

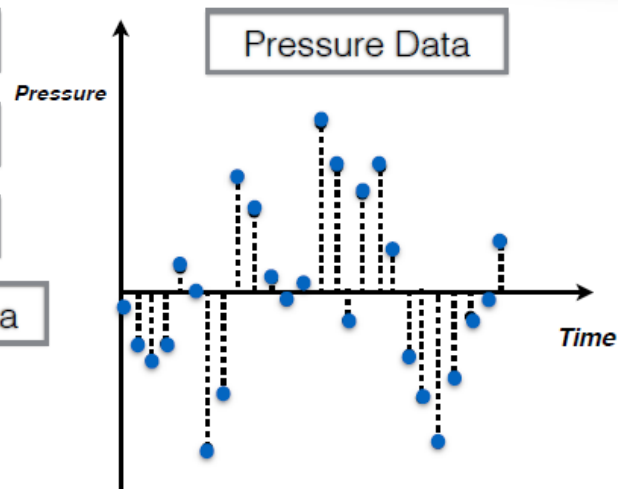
Pitching Data

Rolling Data

Yawing Data

Temperature Data

.....
.....



1

Data Pre-processing

(資料前置處理)

需要領域知識

Data Tables

Time	GPS	Velocity	Pitching
t_1
t_2
t_3
...
...

1

Data Pre-processing : 產線需建立 Data Tables

[需要用到的 Python 技術]

NumPy & Pandas

[One Example] : biopsy data

ID	area	shape	texture	...
id1
id2
id3
...
...

[Another Example] : AAPL 股票

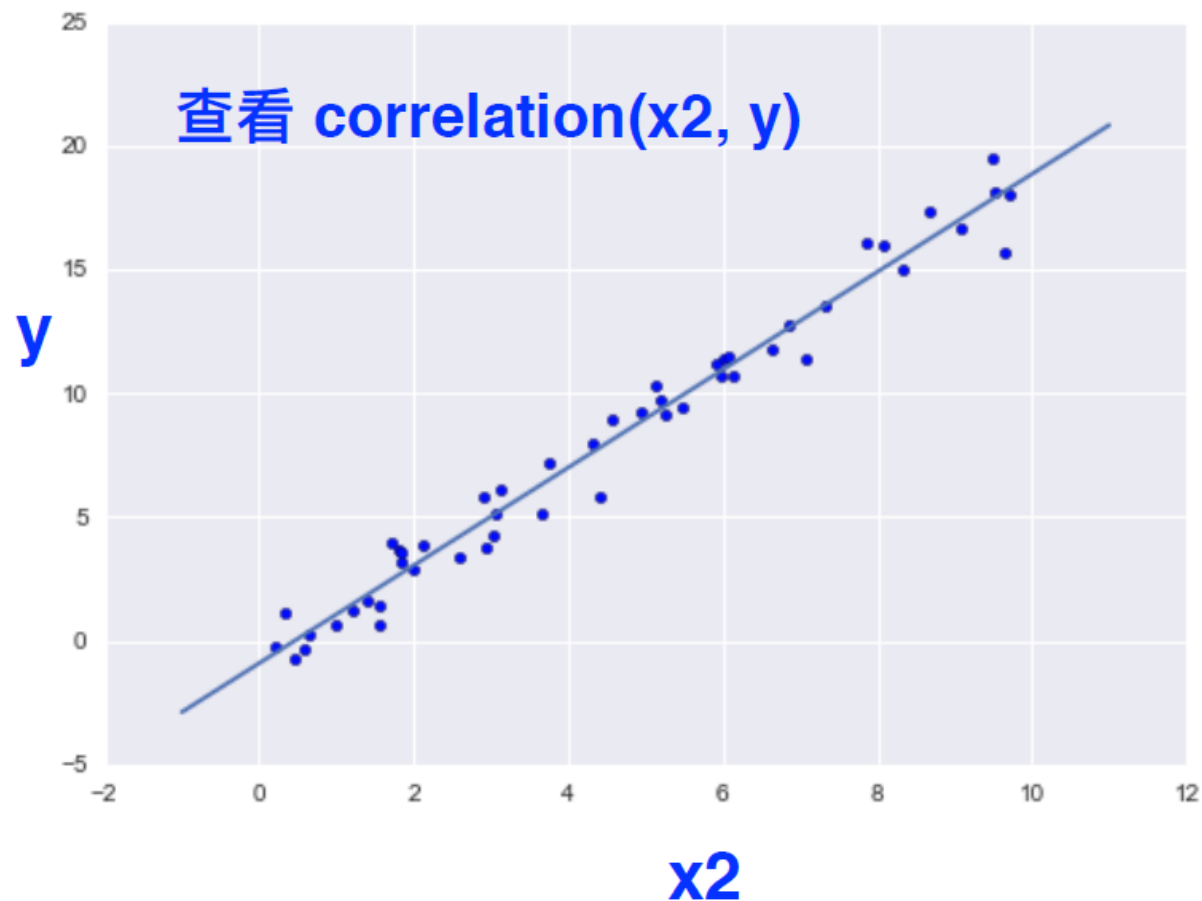
Date	Open	High	Low	Closed	...
d ₁
d ₂
d ₃
...
...

2

Data Visualization : 從 Data Tables 繪圖

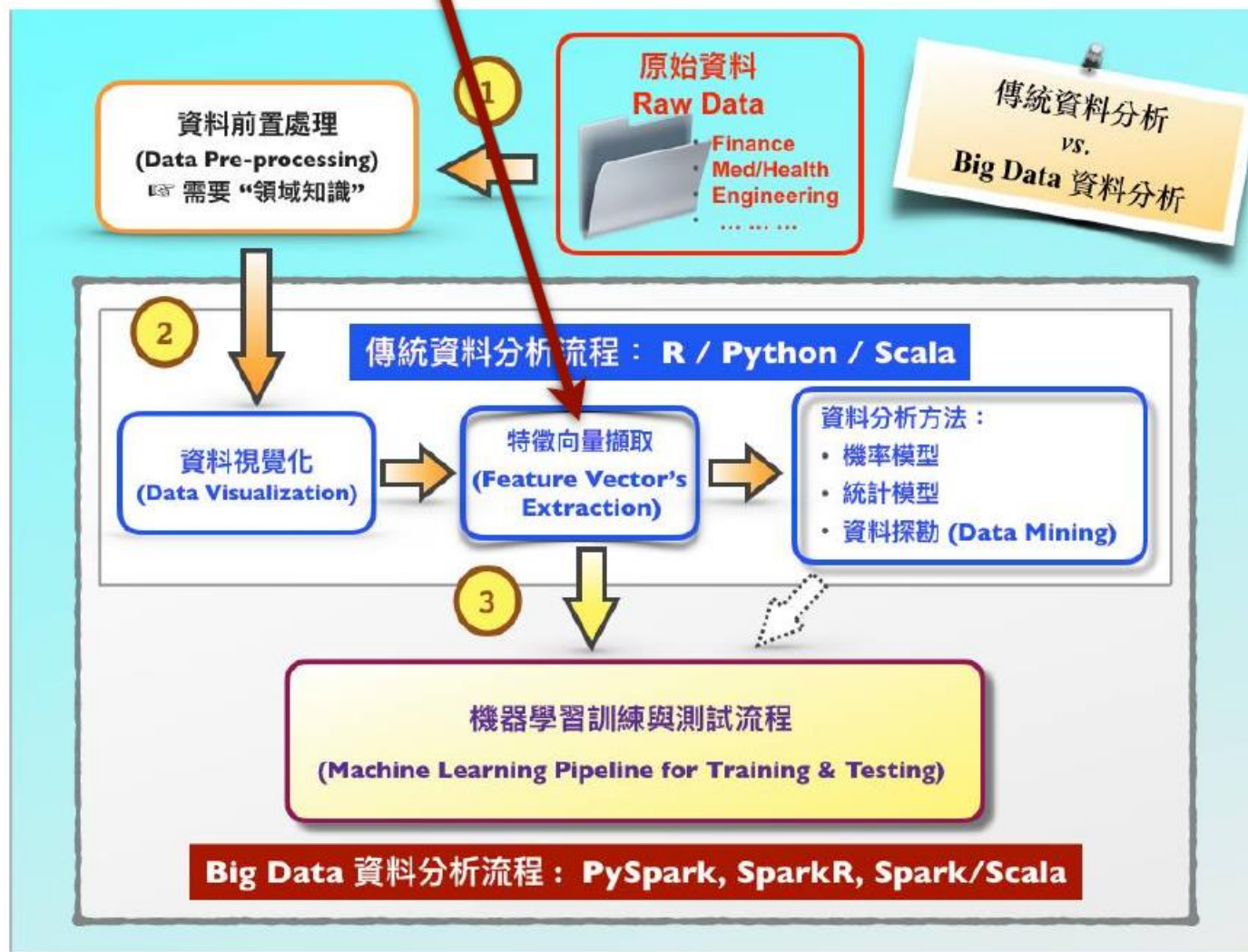
[需要用到的 Python 技術] : Matplotlib

No.	x1	x2	x3	y
1
2
3
4
...



2

Feature Engineering : 決定 Feature Variables



2

Feature Engineering : 決定 Feature Variables

Supervised Learning — Predictive Model

Data
(features)

Labeled Data
(target)

X → Algorithms → y

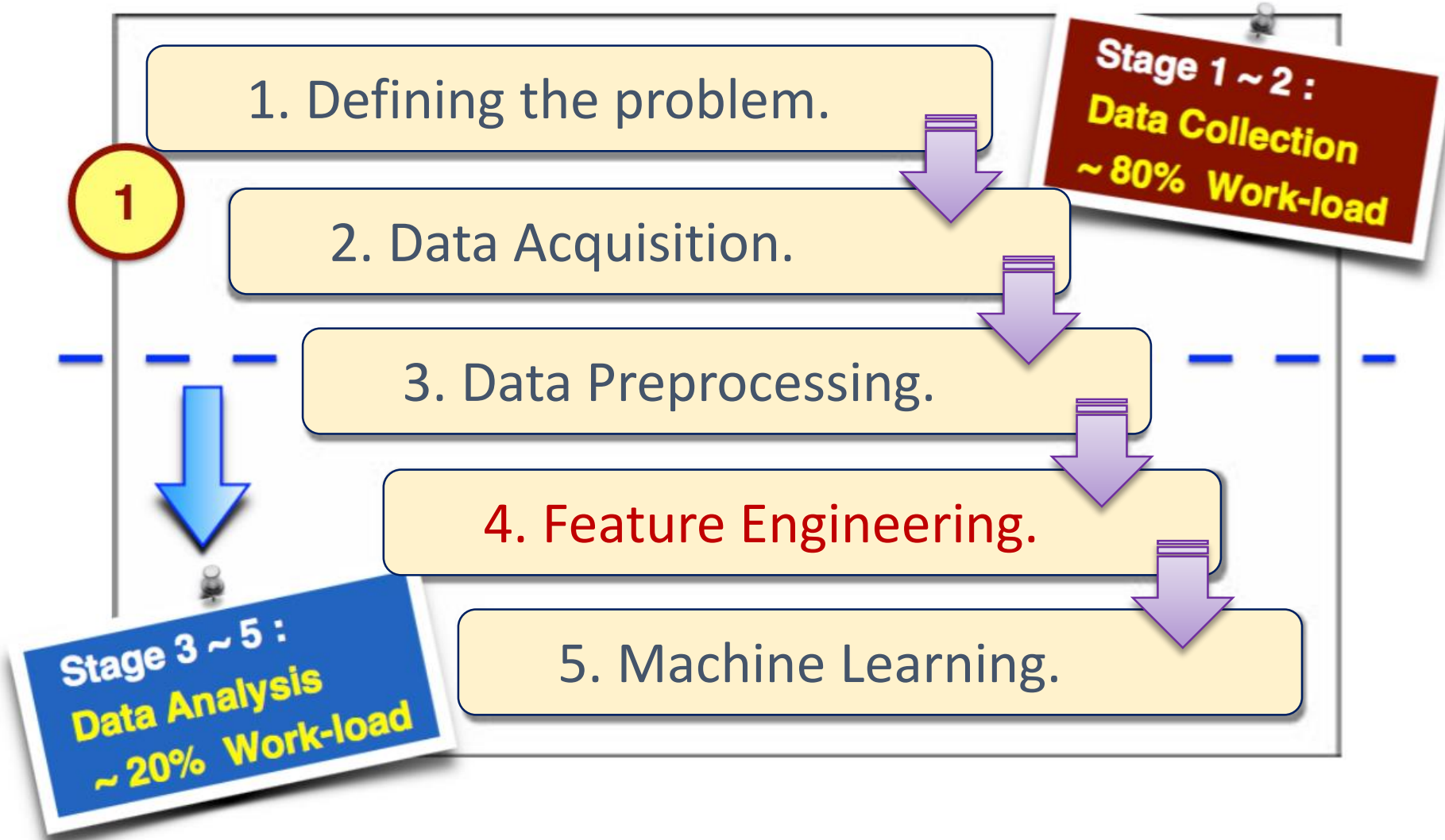
$\{ X, y \}$: training dataset for algorithms



Supervised Learning

⇒ **Machine Learning Training & Testing Pipeline**

Lifecycle of Data-driven Machine Learning Projects



Feature Engineering

Data Preprocessing

- Vectorization
- Normalization
- Missing Values
-
- Data Visualization

Feature Engineering

- Feature Extraction
- **Feature Selection**
 - Feature Importance
 - Dimensionality Reduction
 -

Machine Learning

- Supervized Learning
- Unsupervised Learning

STEP 1

Feature Extraction

- Obtain a pool of appropriate features from the pre-processed data for *feature selection*.

STEP 2

Feature Selection

- Select a set of useful features out of the feature pool for *Machine Learning*.

Feature Importance

Dimensionality Reduction

Data Preprocessing

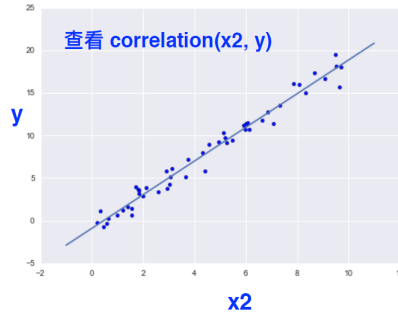
- Vectorization
- Normalization
- Missing Values
-
- **Data Visualization**

2

Data Visualization : 從 Data Tables 繪圖

[需要用到的 Python 技術] : Matplotlib

No.	x1	x2	x3	y
1
2
3
4
...



Feature Engineering

- Feature Extraction
- **Feature Selection**

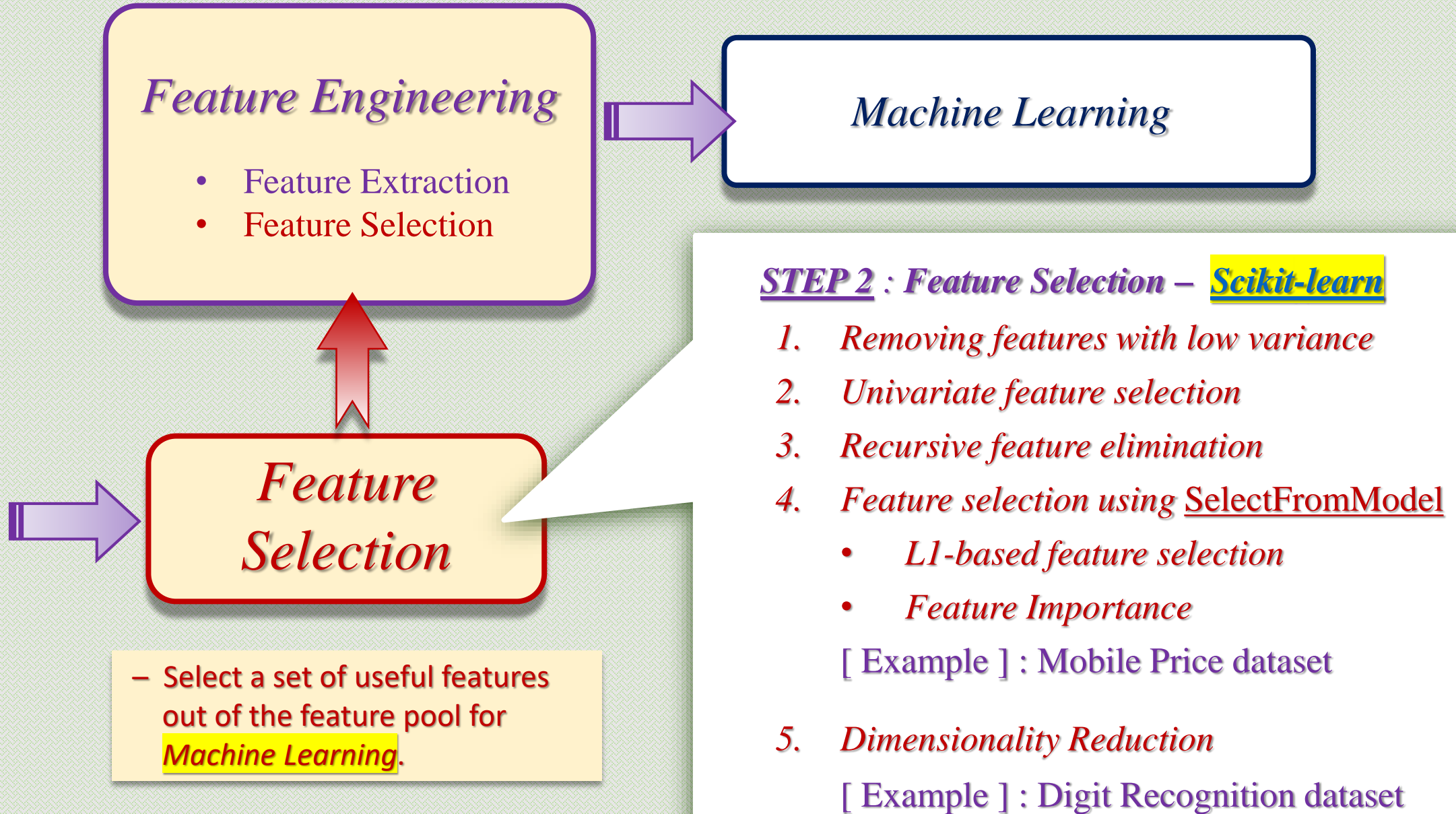
STEP 1 : Feature Extraction by Data Visualization

- *Pair plots*
 - *Correlation matrix (by heatmap)*
- [Example] : Iris dataset

STEP 1

Feature Extraction

- Obtain a pool of appropriate features from the pre-processed data for *feature selection*.



Scikit-learn Workshop

Scikit-Learn Workshop
for
Feature Engineering

CONTENT:

- [1. Introduction](#)
- [2. Feature Extraction by *Data Visualization*](#)
- [3. Feature Selection with `Scikit-Learn`](#)
 - [3.1 Mobile Price dataset](#)
 - [3.1.1 Exploratory Data Analysis \(EDA\)](#)
 - [3.1.2 Feature Extraction for Mobile Price dataset](#)
 - [3.2 Feature Selection Methods](#)
 - [3.2.1 Removing features with low variance](#)
 - [3.2.2 Univariate feature selection](#)
 - [3.2.3 Recursive feature elimination](#)
 - [3.2.4 Feature selection using `SelectFromModel`](#)
 - [L1-based feature selection](#)
 - [Feature Importance](#)
 - [3.3 Classification of `Mobile Price Dataset`](#)
- [4. Feature Selection with Dimensionality Reduction](#)
 - [4.1 Digit Recognition dataset](#)
 - [4.2 Feature Selection with `PCA`](#)
 - [4.3 Classification of `Digit Recognition Dataset`](#)

Iris dataset

Mobile Price

Digit
Recognition