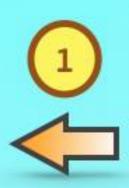
10. AI 代理人與機器學習理論基礎

< PART III Python 機器學習實務 >

OUTLINE

- 資料分析流程
- Raw Data (原始資料) 前置處理 NumPy, Pandas
- 資料視覺化 Matplotlib
- 特徵工程 (Feature Engineering)
- 機器學習 Training & Testing Pipeline Scikit-Learn
- Al Agent & Learning Agent

資料前置處理 (Data Pre-processing) ☞ 需要"領域知識"



原始資料 Raw Data



傳統資料分析 vs. Big Data 資料分析

2

傳統資料分析流程: R / Python / Scala

資料視覺化 (Data Visualization)



特徵向量擷取

(Feature Vector's Extraction)



資料分析方法:

- 機率模型
- 統計模型
- · 資料探勘 (Data Mining)

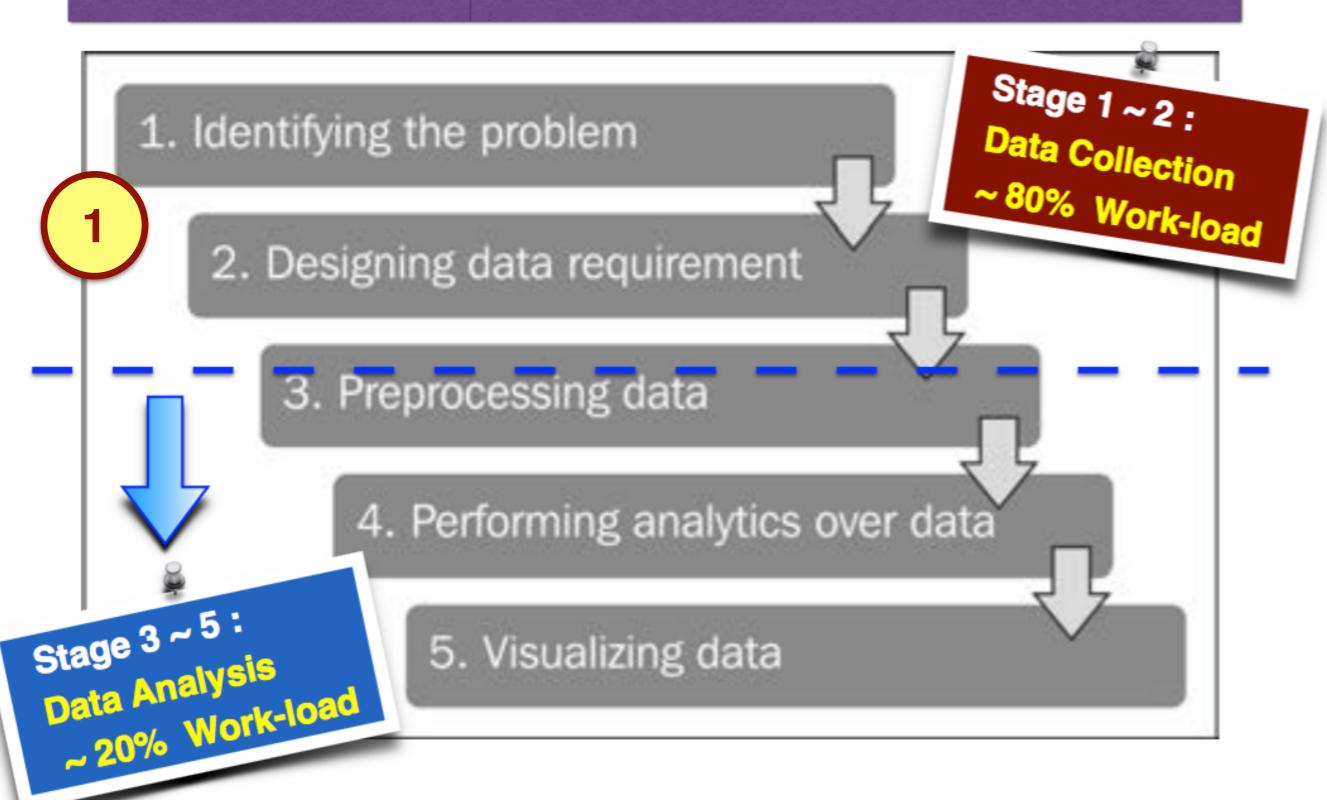
3 4

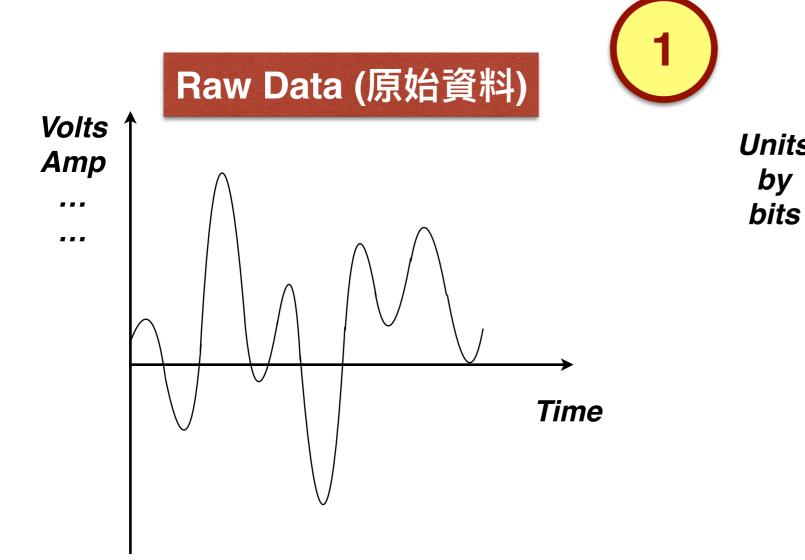
機器學習訓練與測試流程

(Machine Learning Pipeline for Training & Testing)

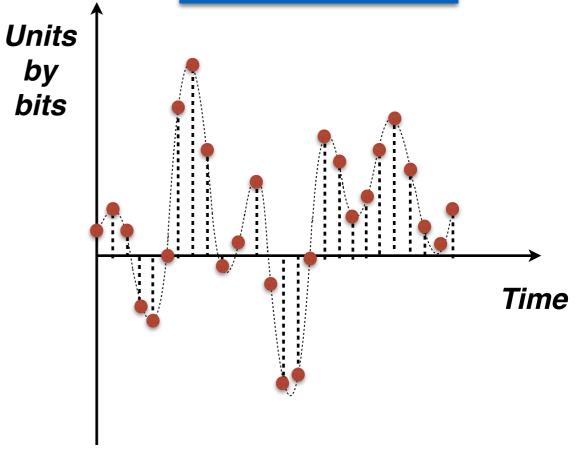
Big Data 資料分析流程: PySpark, SparkR, Spark/Scala

Data Analytics Project Life Cycle (5 stages)





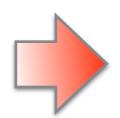




Data Acquisition (數據擷取)

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

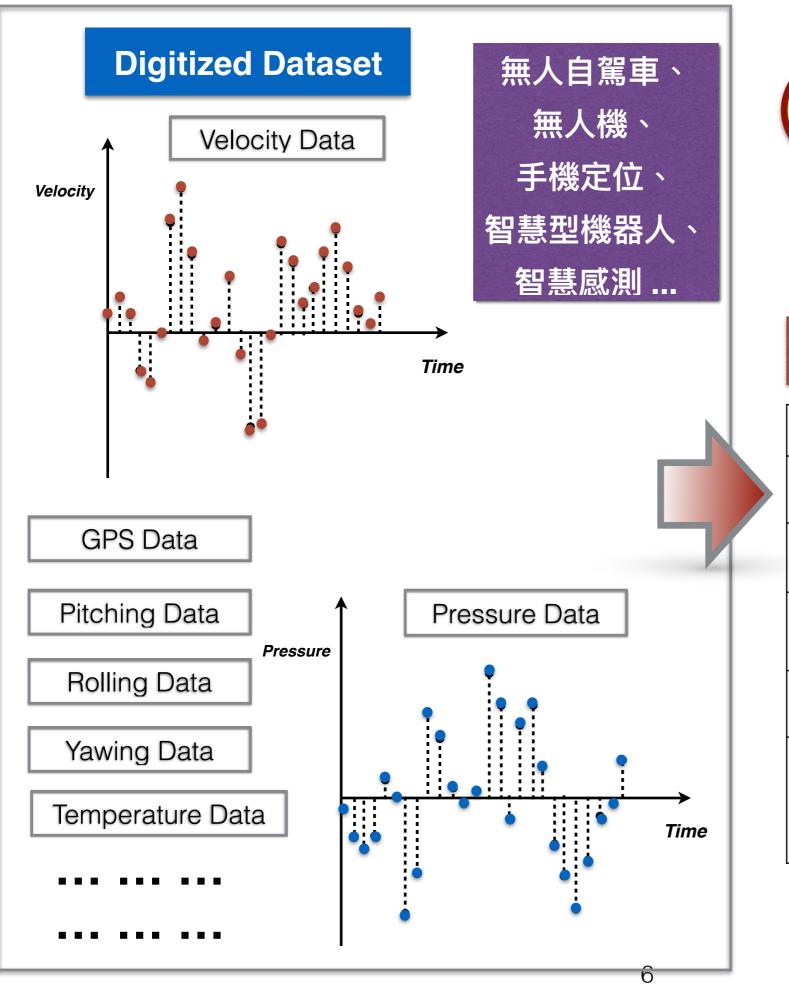
Analog Signal (類比訊號)



A/D C



Digital Signal (數位訊號)





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			•••	•••
	•••	•••		•••



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	•••		•••	
		•••		

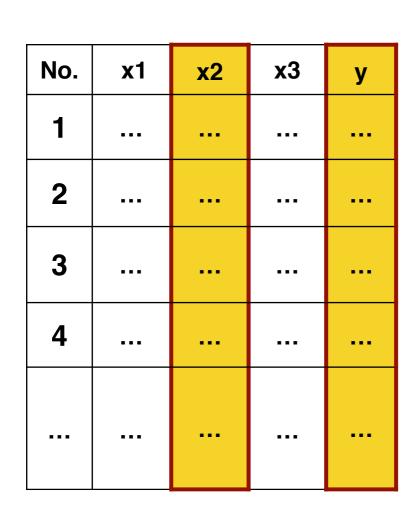
[Another Example]: AAPL 股票

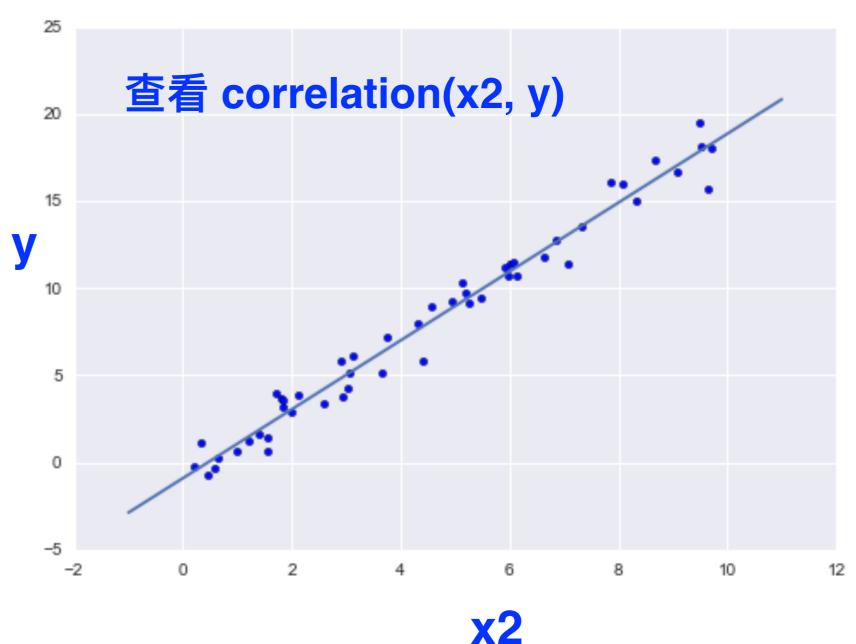
Date	Open	High	Low	Closed	
d ₁	•••	•••			
d ₂					
d ₃					
•••	•••		•••		
	•••	•••	•••		



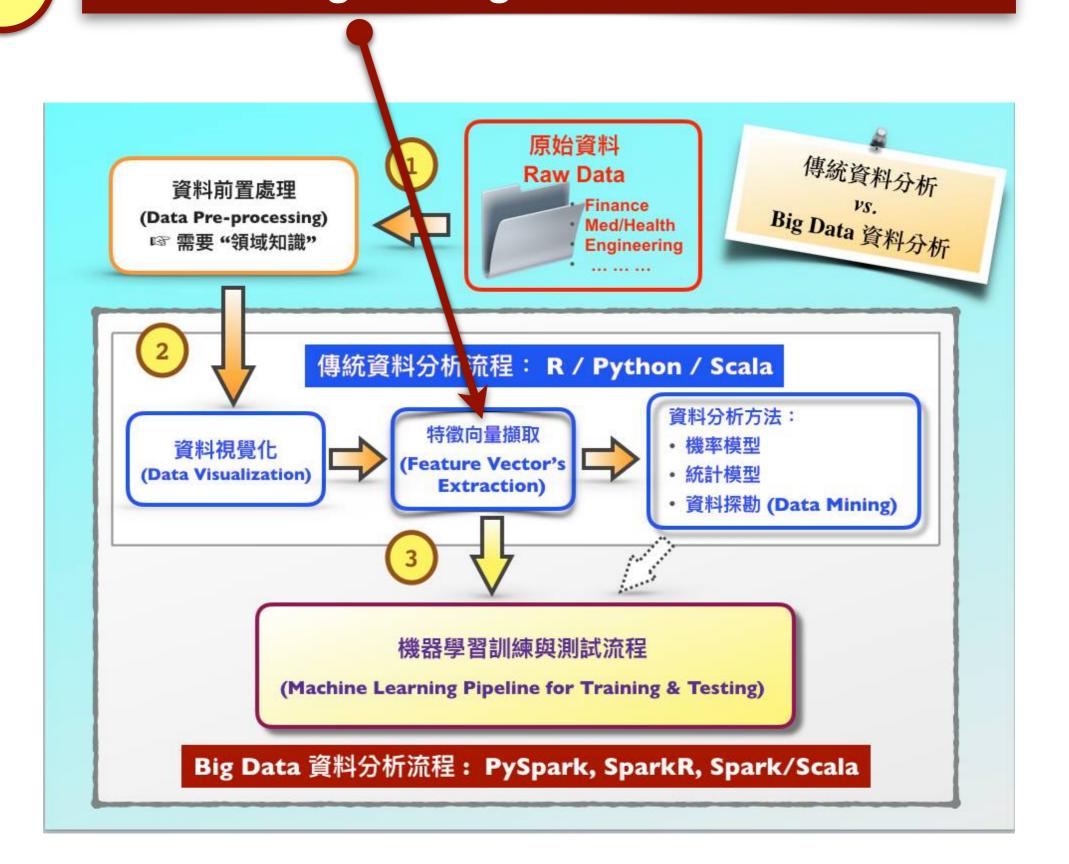
Data Visualization:從 Data Tables 繪圖

[需要用到的 Python 技術]: Matplotlib



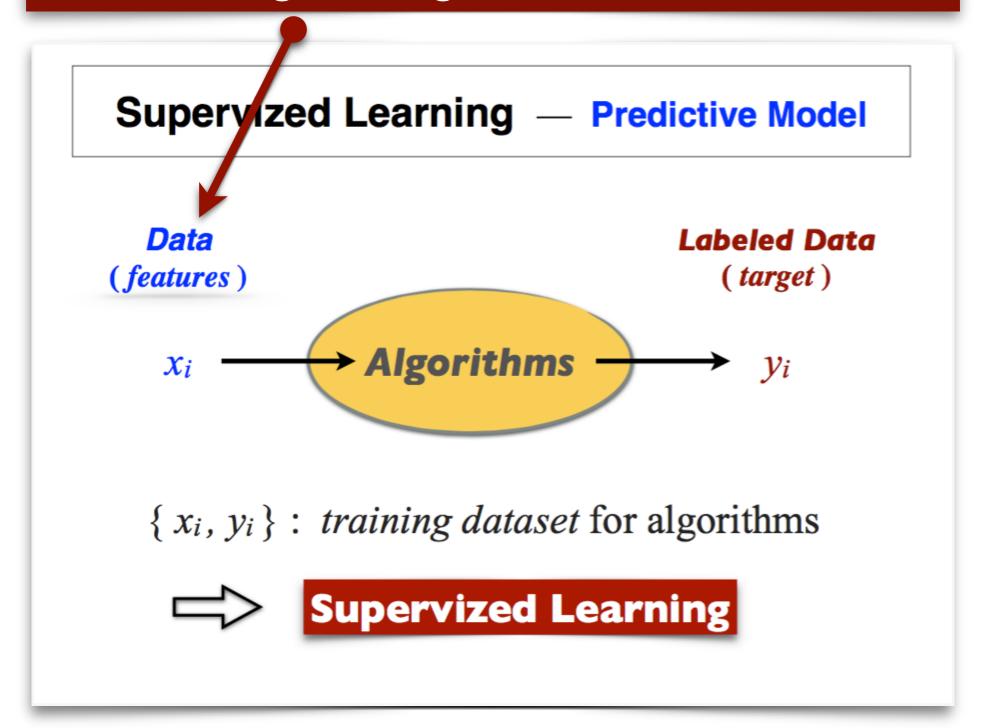


Feature Engineering: 決定 Feature Variables





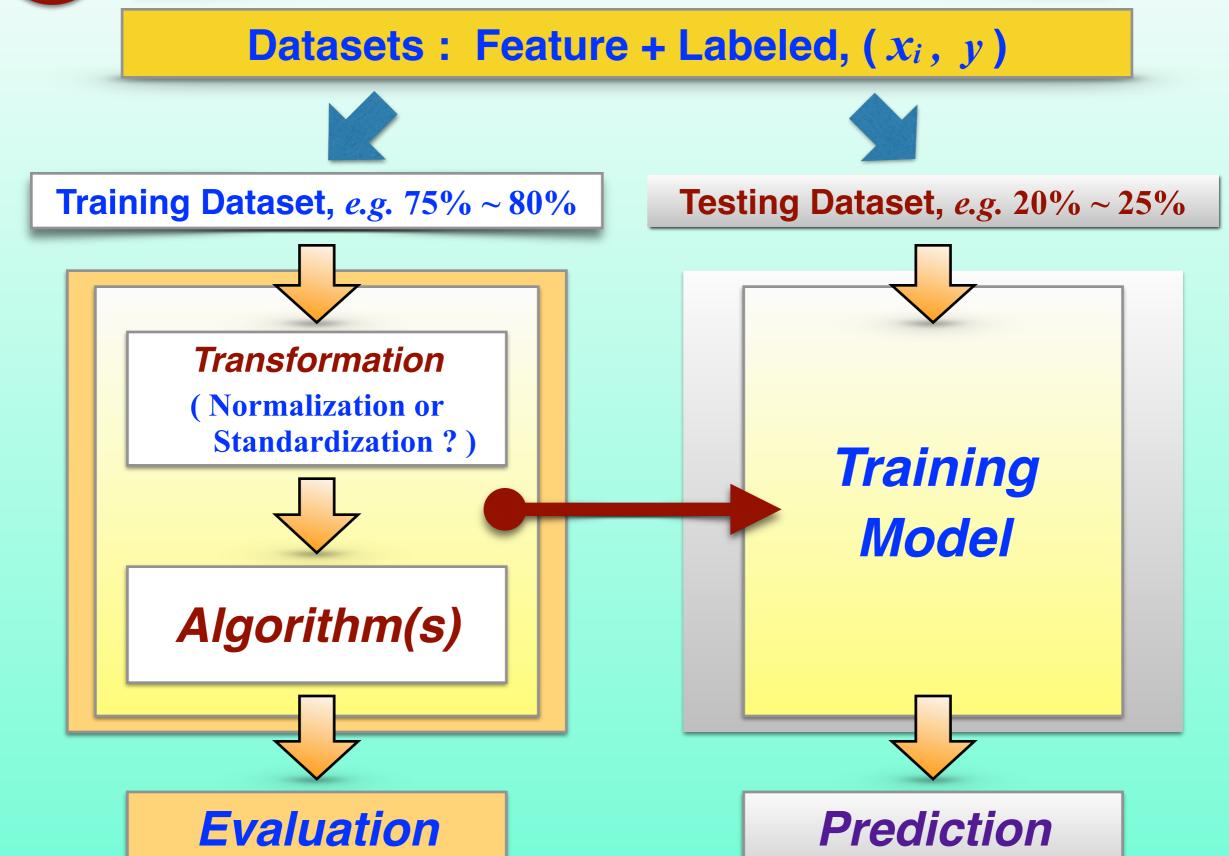
Feature Engineering: 決定 Feature Variables



=> Machine Learning Training & Testing Pipeline



Machine Learning Training & Testing Pipeline



11



Machine Learning Training & Testing Pipeline

[需要用到的 Python 技術]: Scikit-Learn

Q: 為什麼 Feature Dataset, xi 要拆解成 75% (or 80%): 25% (or 20%)?

可不可以拆解成其他比例呢? 例如:60%:40%

Q:如何決定該選擇 Normalization 或 Standardization?

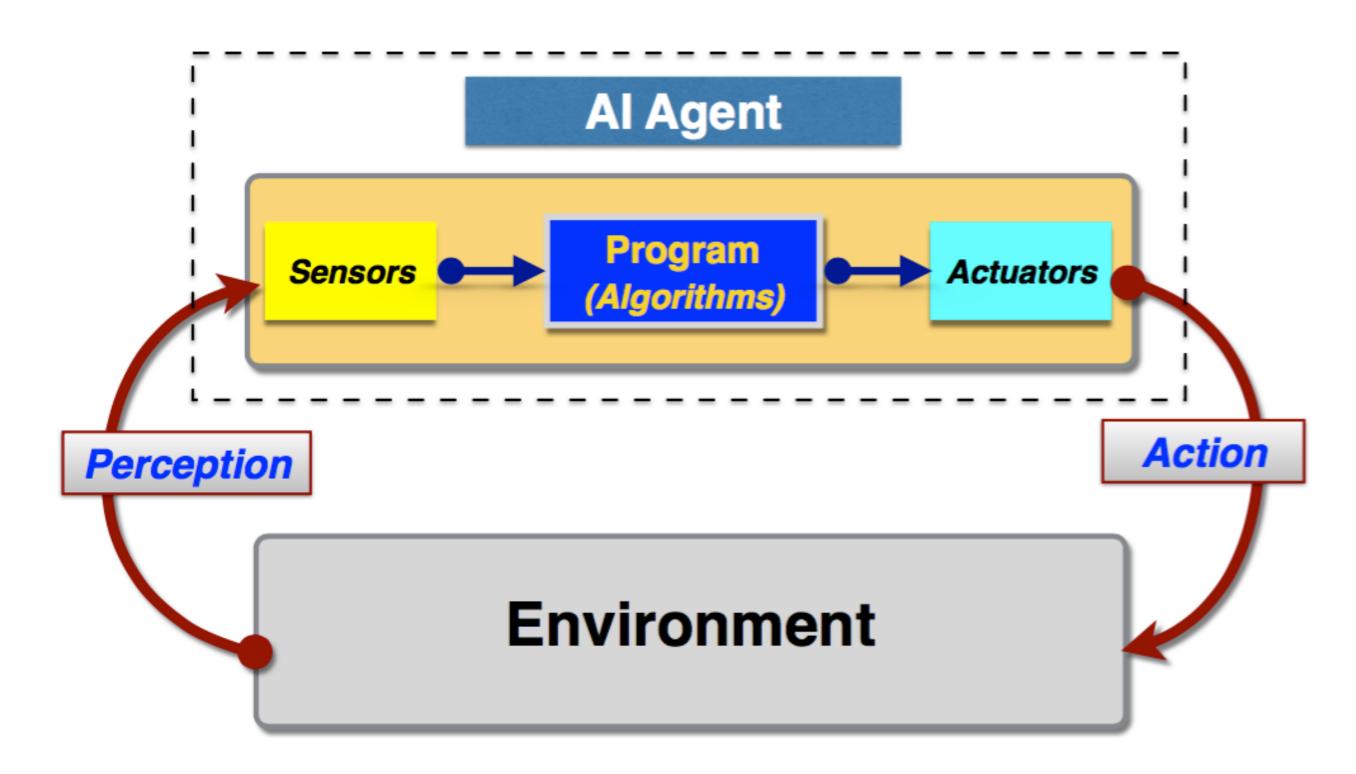
Q:如何在Training階段選擇Algorithms?

Q:如何在Training階段,進行Evaluation (評估)? [labeled data]

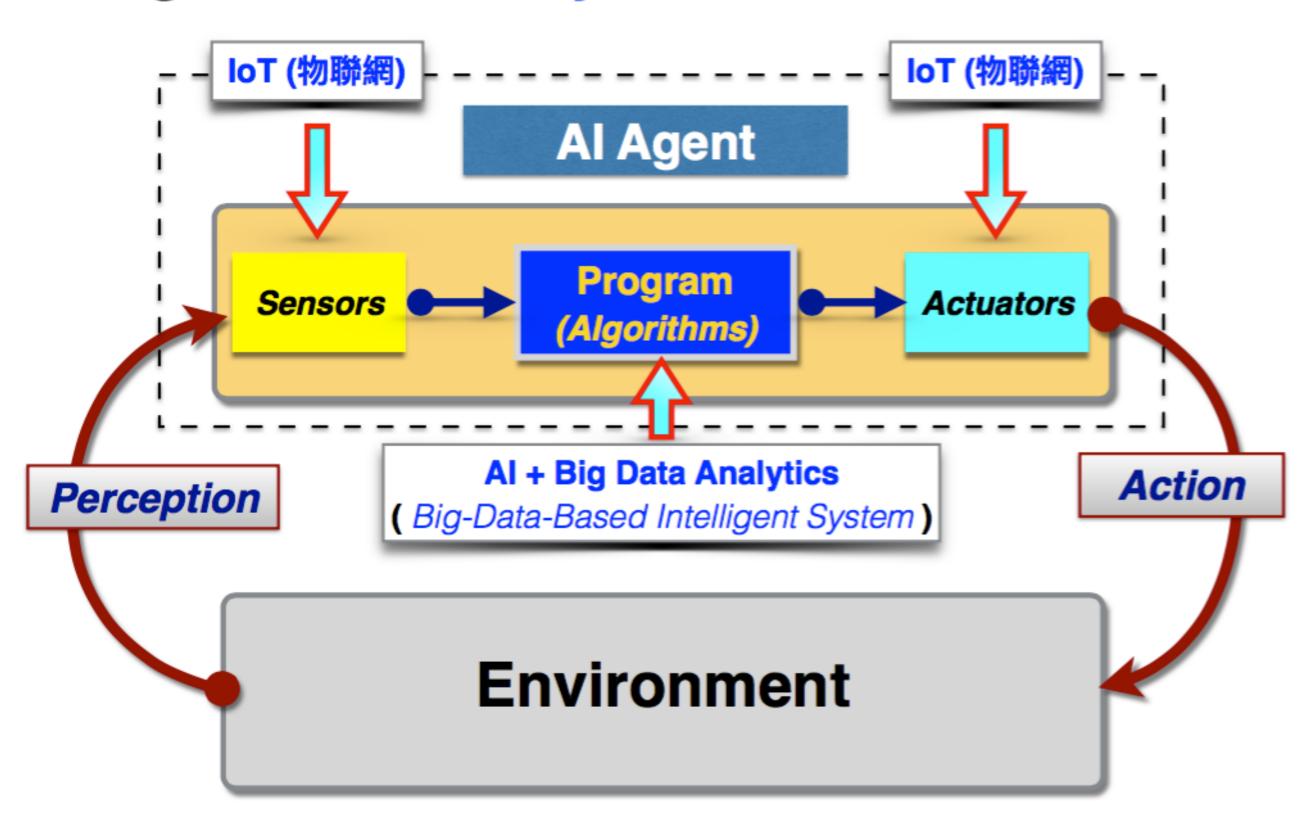
Q:如果 Training Model 在 Testing 階段,預測結果表現"很好"或"不好",該如何處理呢?

Risk Analysis: Test Result => Baseline => Strategy
Risk Management => Execution

Al Agent



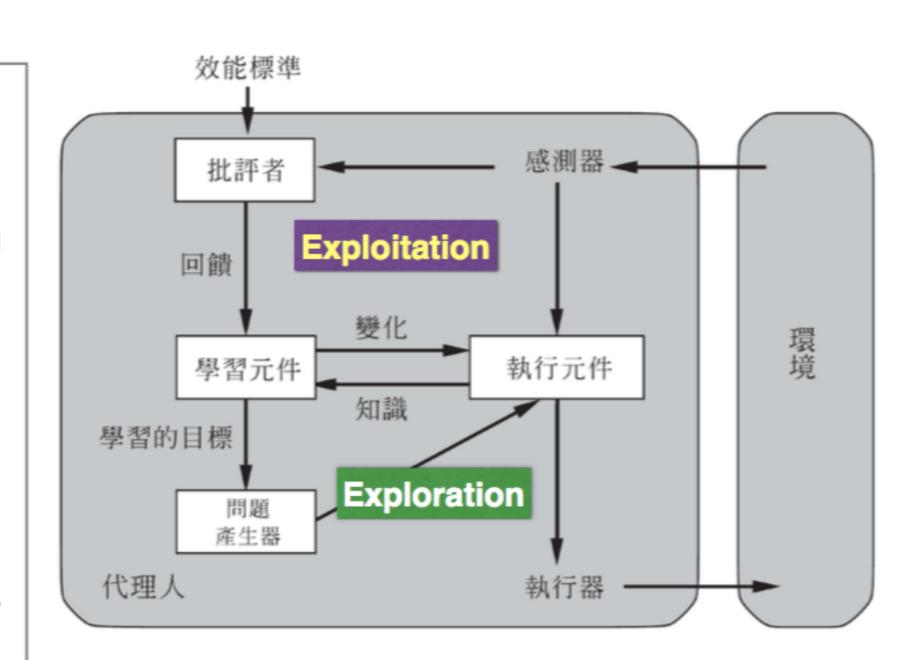
Al Agent for Industry 4.0 Solutions



一般的學習型代理人 Exploitation vs. Exploration

4 個概念上的元件:

- 學習元件 負責做出改進
- 執行元件 負責選擇外部動作 (相 當於先前提過的代理人
- 批評者 評價代理人做得如何之後,學習元件利用來自批評者的回饋來決定應該如何修改執行元件, 使得在未來能夠做得更好
- 問題產生器 負責提出探索活動, 收集新資訊及經驗



一般的學習型代理人

REFERENCE

Jake VanderPlas, "*Python Data Science Handbook*", 2016, O'Reily.

https://jakevdp.github.io/PythonDataScienceHandbook/

Code from Github:

https://github.com/jakevdp/PythonDataScienceHandbook