

機器學習簡介

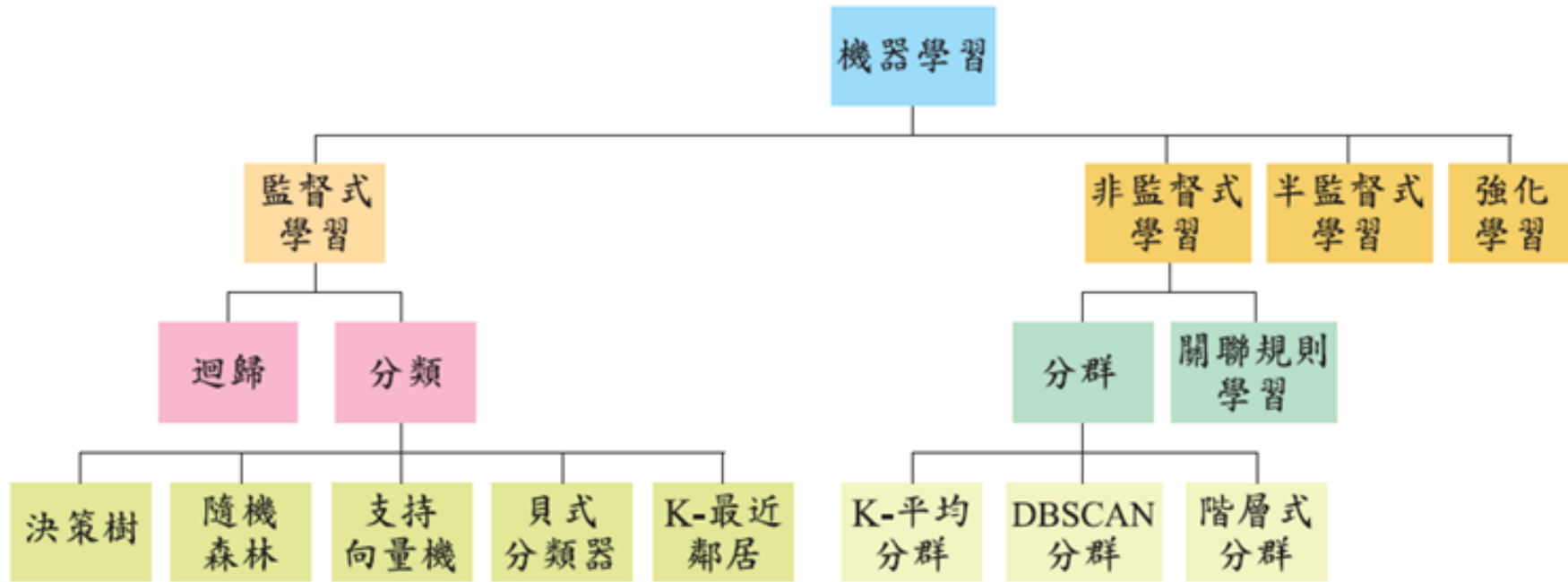
- 機器學習(Machine Learning)是一種數據分析技術，它教導計算機模仿人類從經驗中學習。
- Clustering - 聚類
- Classification - 分類
- Regression - 回歸



什麼是機器學習？

機器學習簡介

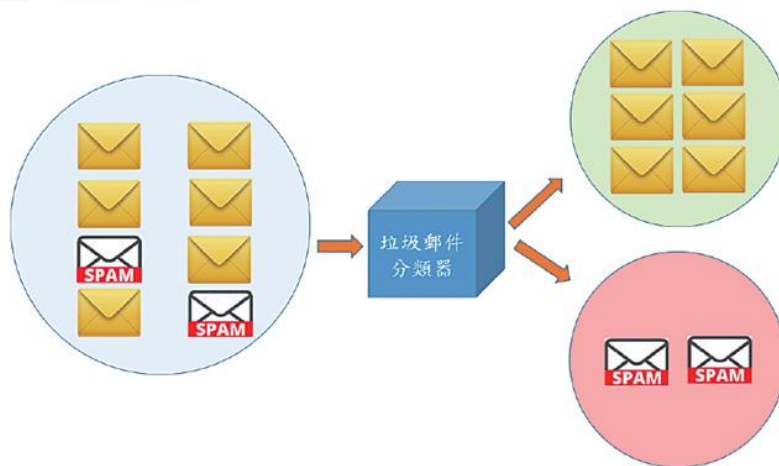
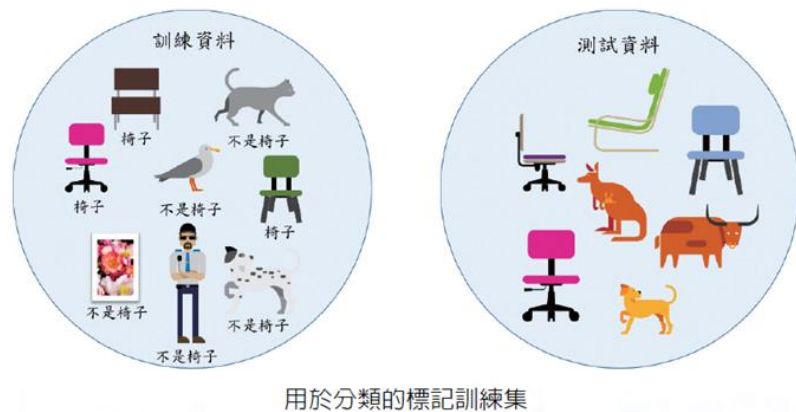
- 機器學習的分類



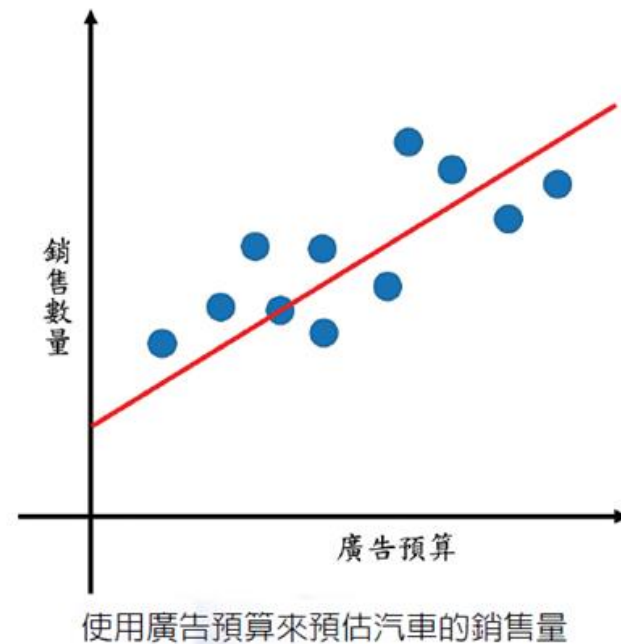
機器學習的分類

監督式學習

- 監督式學習 (Supervised learning)



用於監督學習的垃圾郵件分類



非監督式學習

- 非監督式學習 (Un-Supervised learning)

Machine 1	machine B	Machine 2
machine A	Machine C	machine 3

(a)原始資料

Machine 1	machine A	Machine 2
machine B	Machine C	machine 3

(b)依據顏色分群

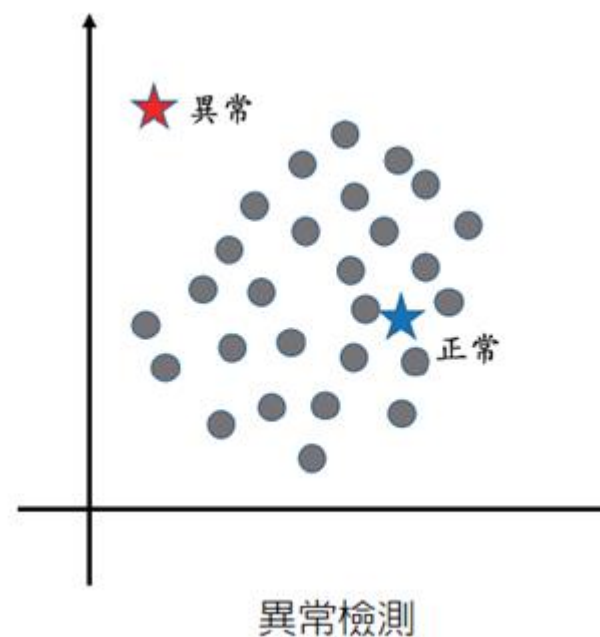
Machine 1	Machine C	Machine 2
machine A	machine B	machine 3

(c)依據字首的大小寫分群

machine A	machine B	Machine C
Machine 1	Machine 2	machine 3

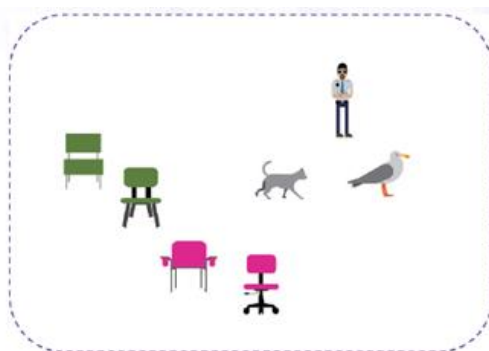
(d)依據字中是否包含數字分群

無監督學習的分群

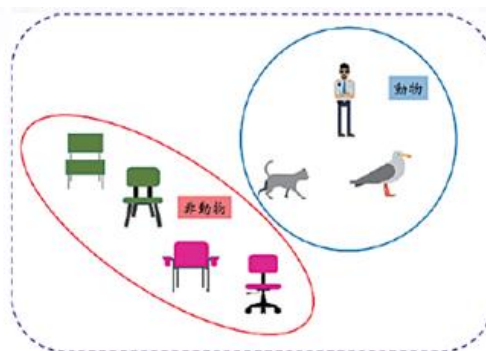


非監督式學習

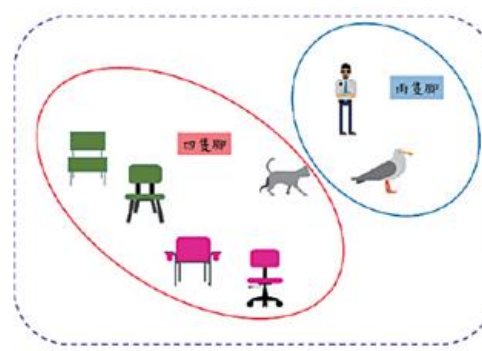
- 非監督式學習 (Un-Supervised learning)



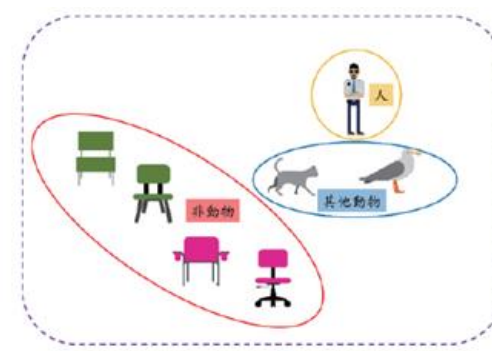
(a)原始資料集



(b)依據是否為動物來分成二群



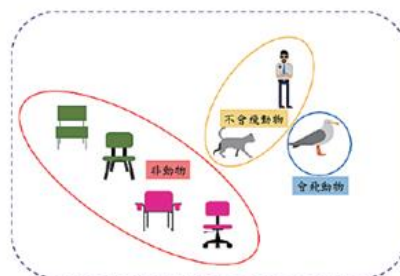
(c)依據腳的數量來分成兩隻腳和四隻腳兩群



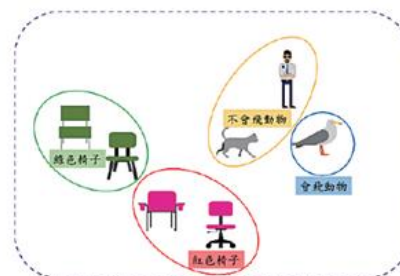
(d)依據人、其他動物和非動物分成三群

不同的分群結果

不同的分群結果 (續)



(e)依據會飛動物、不會飛動物和非動物分成三群



(f)依據會飛動物、不會飛動物、綠色椅子和紅色椅子分成四群

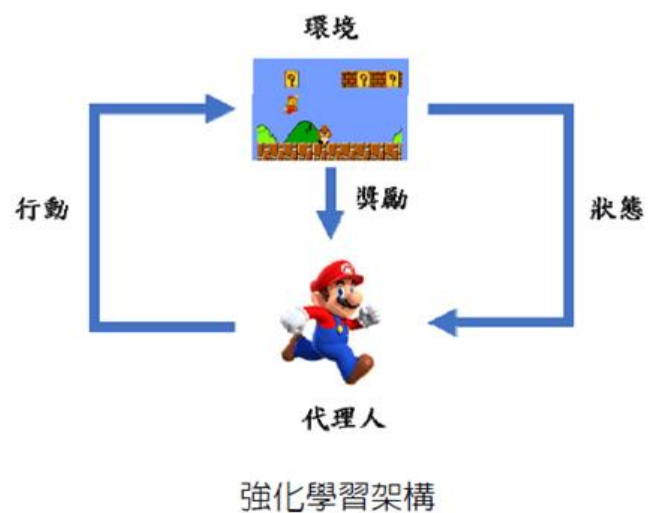
不同的分群結果

半監督式學習

- 半監督式學習 (Semi-supervised learning)
 - 半監督式學習顧名思義就是結合監督式學習與非監督式學習的一種方法。
 - 在實際的應用中，有些領域人工標示的樣本成本很高，但是沒有標籤的樣本數量卻非常龐大。
 - 所以半監督式學習就是利用大量的無標籤樣本和少量的標籤樣本來解決標籤樣本不足的問題。

強化學習

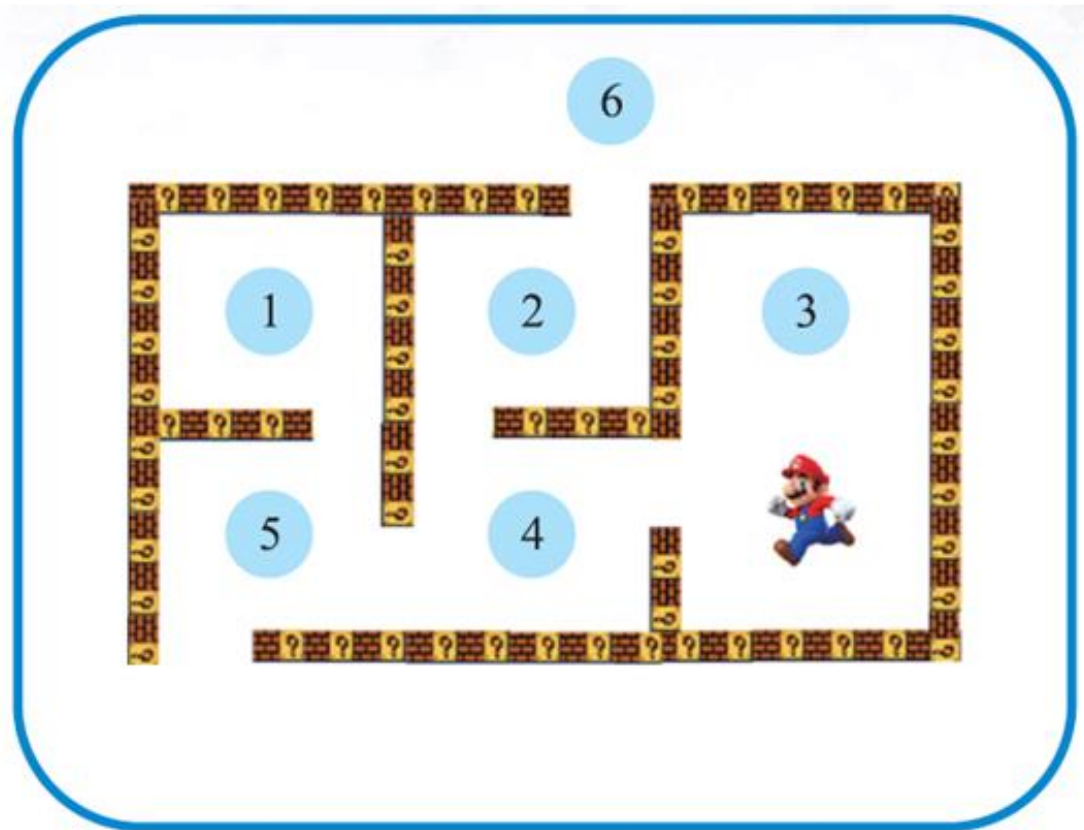
- 強化學習 (Reinforcement learning)



AlphaGo

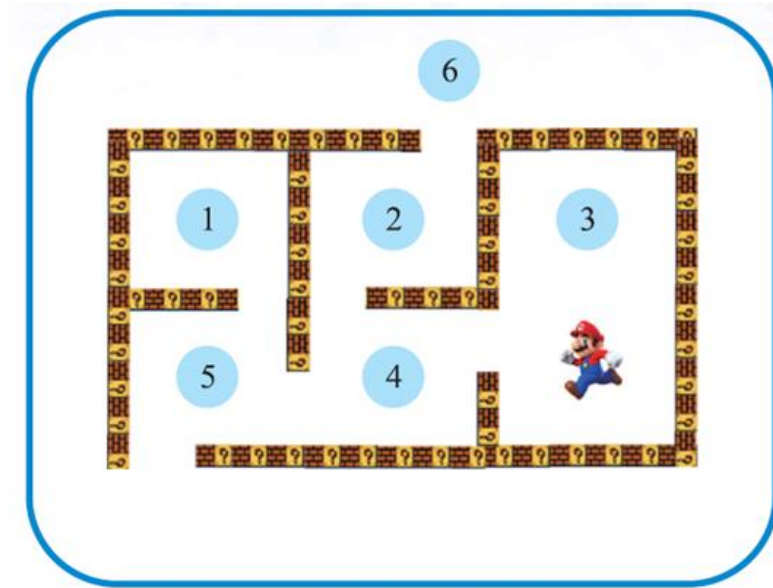
強化學習 (Q-learning)

- 強化學習之 Q-learning

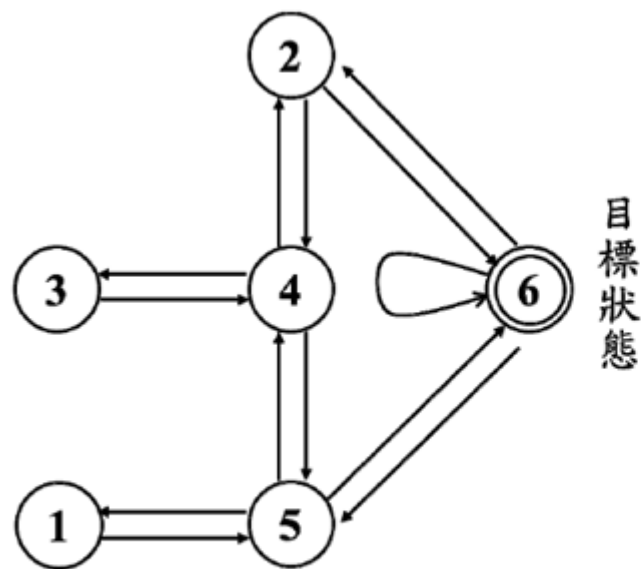


迷宮的地圖

強化學習 (Q-learning)

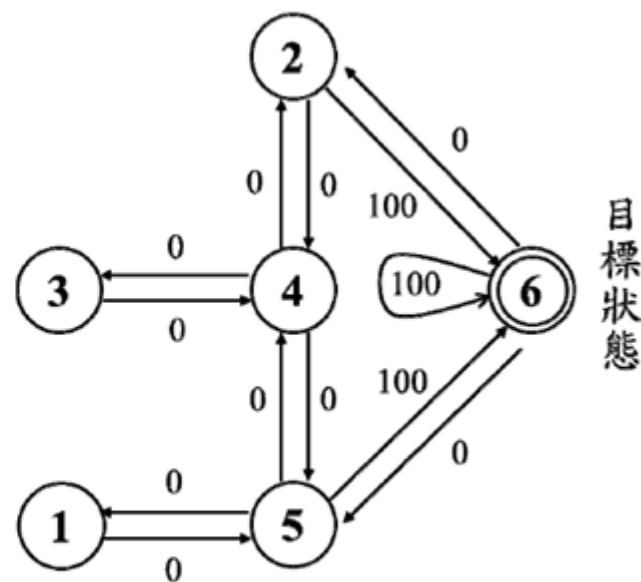


迷宮的地圖



目標狀態

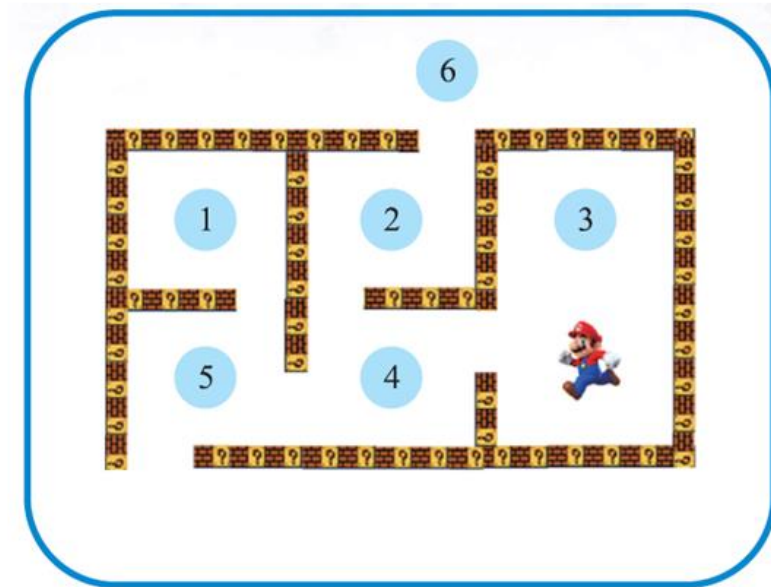
以圖來表示迷宮地圖



目標狀態

標上獎勵的圖

強化學習 (Q-learning)



迷宮的地圖

$$Q(\text{狀態}, \text{行動}) = R(\text{狀態}, \text{行動}) + \text{折扣因子} * \text{Max}[Q(\text{下個狀態}, \text{所有行動})]$$

行: 行動
列: 狀態

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

強化學習 (Q-learning)

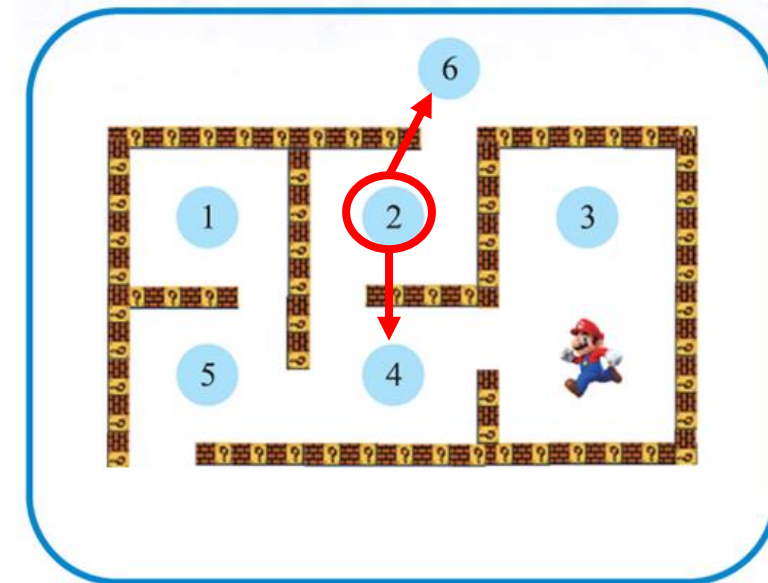
Step1.

折扣因子=0.8

初始狀態=2

Q矩陣=0

假設選到行動=6



迷宮的地圖

$$Q(\text{狀態}, \text{行動}) = R(\text{狀態}, \text{行動}) + \text{折扣因子} * \text{Max}[Q(\text{下個狀態}, \text{所有行動})]$$

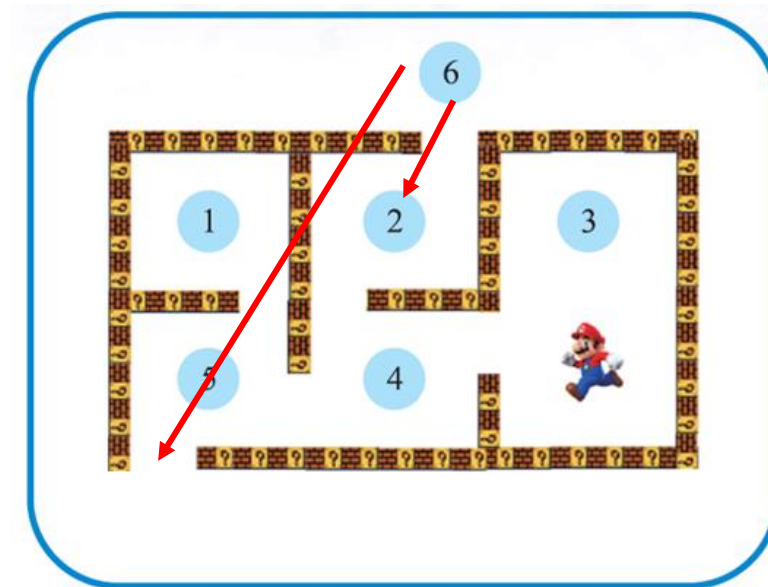
$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

強化學習 (Q-learning)

Step2.

計算 $Q(2, 6)$



迷宮的地圖

$$\begin{aligned} Q(2, 6) &= R(2, 6) + 0.8 * \text{Max}[Q(6, 2), Q(6, 5), Q(6, 6)] \\ &= 100 + 0.8 * 0 \\ &= 100 \end{aligned}$$

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

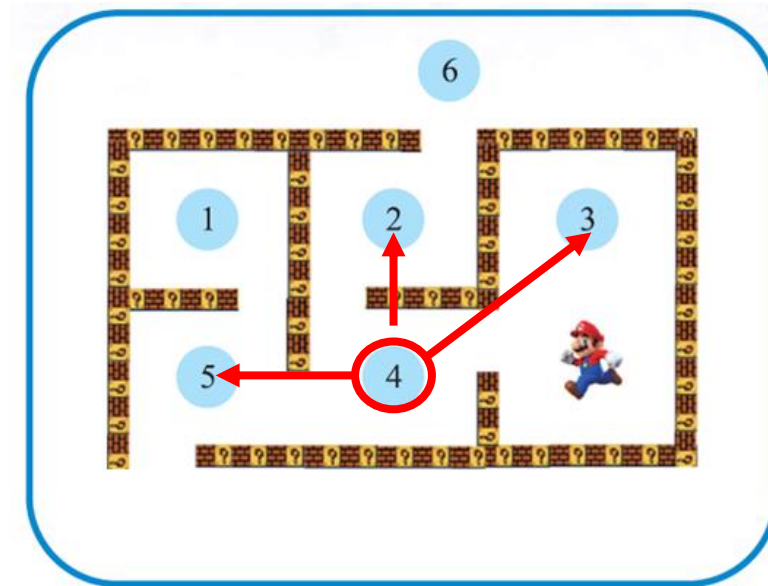
強化學習 (Q-learning)

Step1.

折扣因子=0.8

初始狀態=4

假設選到行動=2



迷宮的地圖

$$Q(\text{狀態}, \text{行動}) = R(\text{狀態}, \text{行動}) + \text{折扣因子} * \text{Max}[Q(\text{下個狀態}, \text{所有行動})]$$

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

強化學習 (Q-learning)

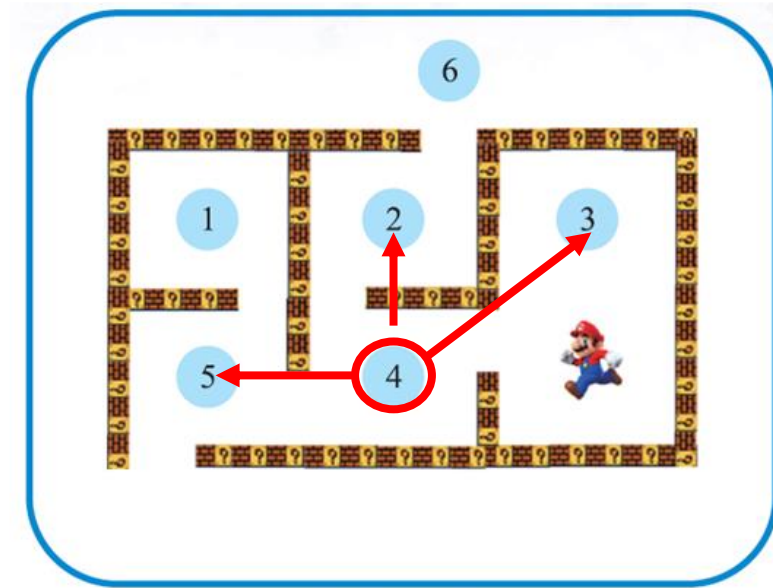
Step2.

計算 $Q(4, 2)$

$$\begin{aligned} Q(4, 2) &= R(4, 2) + 0.8 * \text{Max}[Q(2, 4), Q(2, 6)] \\ &= 0 + 0.8 * 100 \\ &= 80 \end{aligned}$$

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$



迷宮的地圖

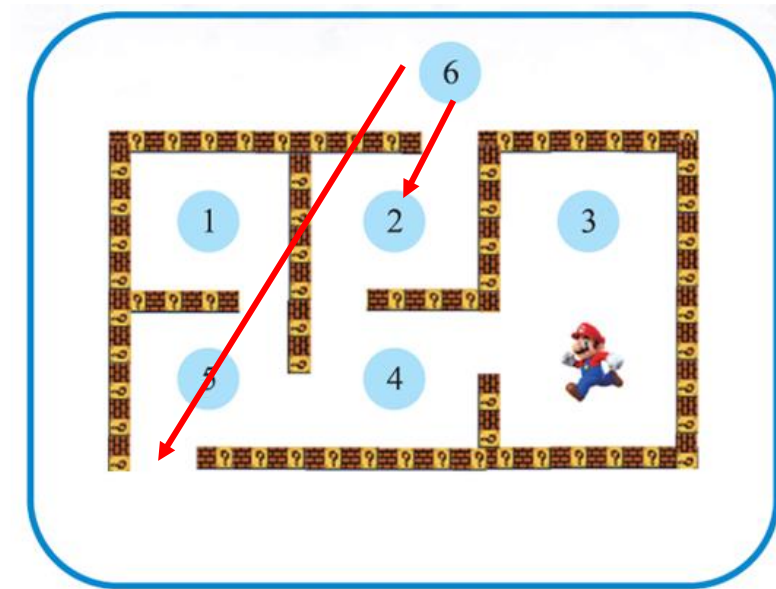
強化學習 (Q-learning)

Step3.

折扣因子=0.8

狀態=2

假設選到行動=6



迷宮的地圖

$$Q(\text{狀態}, \text{行動}) = R(\text{狀態}, \text{行動}) + \text{折扣因子} * \text{Max}[Q(\text{下個狀態}, \text{所有行動})]$$

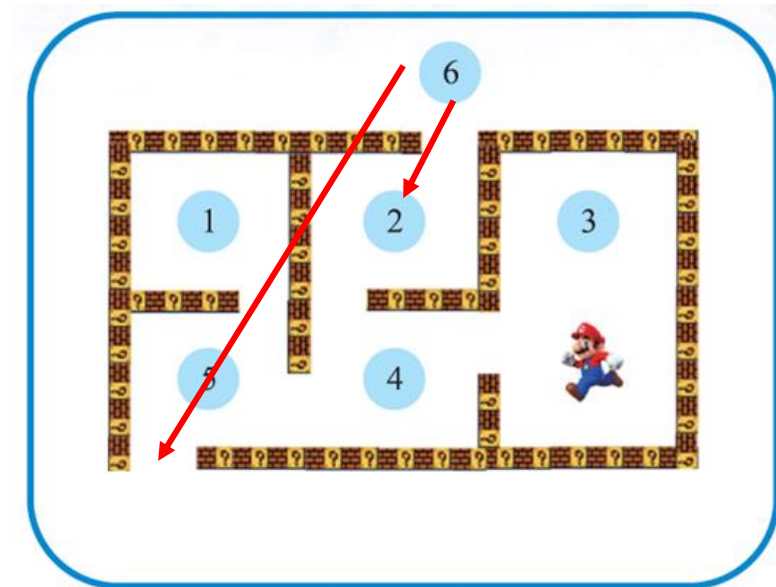
$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

強化學習 (Q-learning)

Step4.

計算 $Q(2, 6)$



迷宮的地圖

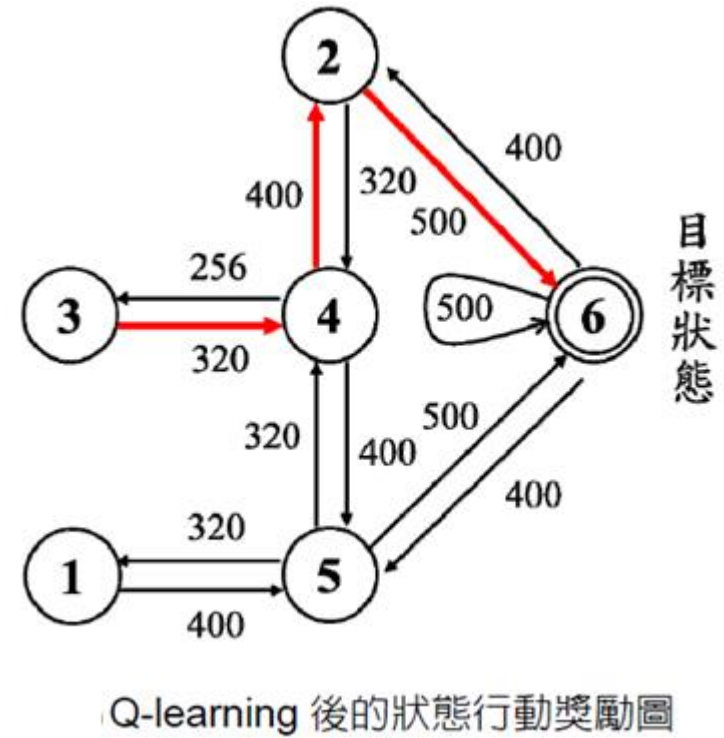
$$\begin{aligned} Q(2, 6) &= R(2, 6) + 0.8 * \text{Max}[Q(6, 2), Q(6, 5), Q(6, 6)] \\ &= 100 + 0.8 * 0 \\ &= 100 \end{aligned}$$

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$R = \begin{bmatrix} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{bmatrix}$$

強化學習 (Q-learning)

$$Q = \begin{bmatrix} 0 & 0 & 0 & 0 & 400 & 0 \\ 0 & 0 & 0 & 320 & 0 & 500 \\ 0 & 0 & 0 & 320 & 0 & 0 \\ 0 & 400 & 256 & 0 & 400 & 0 \\ 320 & 0 & 0 & 320 & 0 & 500 \\ 0 & 400 & 0 & 0 & 400 & 500 \end{bmatrix}$$

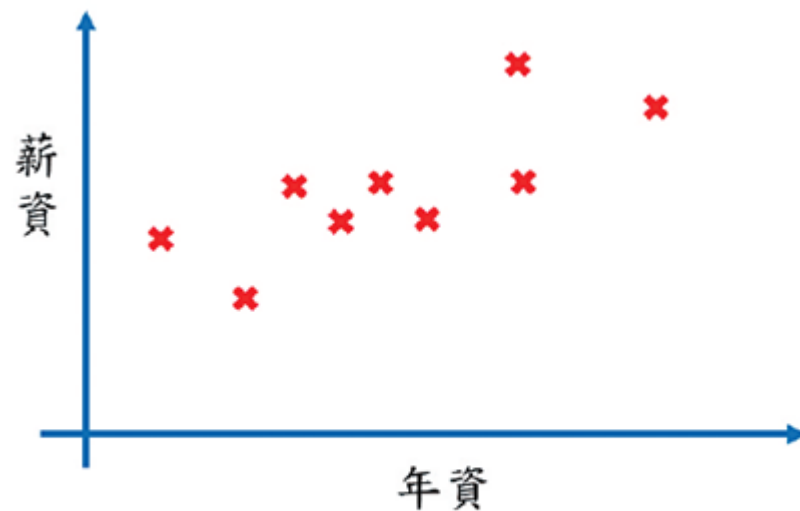


機器學習演算法

- 迴歸 (Regression)

年資和薪資的資料集

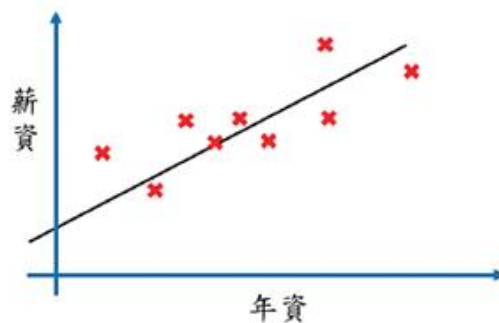
年資	薪資
1	30000
2	25000
3	35000
3.5	32000
4	35000
4.5	32000
6	35000
6	50000
7	40000



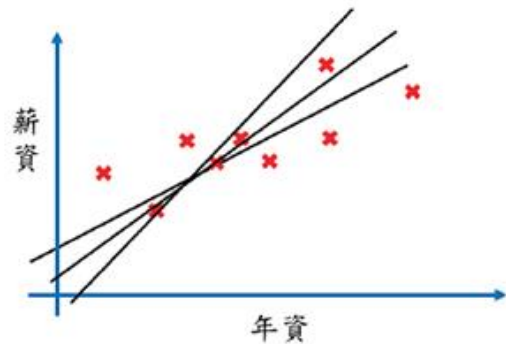
年資和薪資的二維分佈圖

迴歸 (Regression)

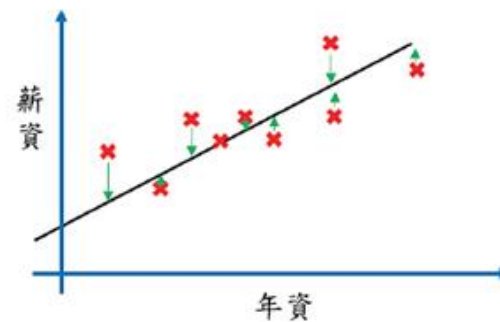
- 迴歸預測



(a) 迴歸直線方程式



(b) 無數條的直線方程式

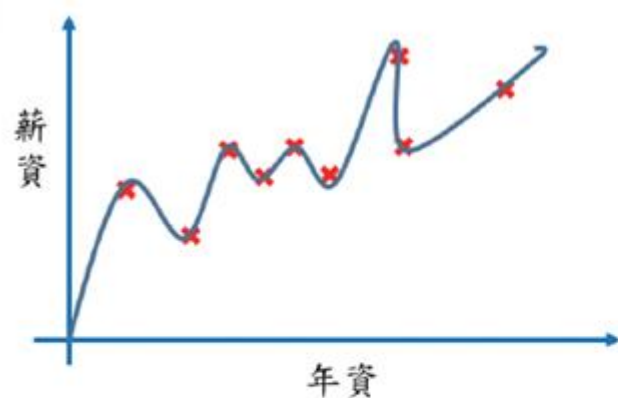


(c) 誤差平方和為最小的直線

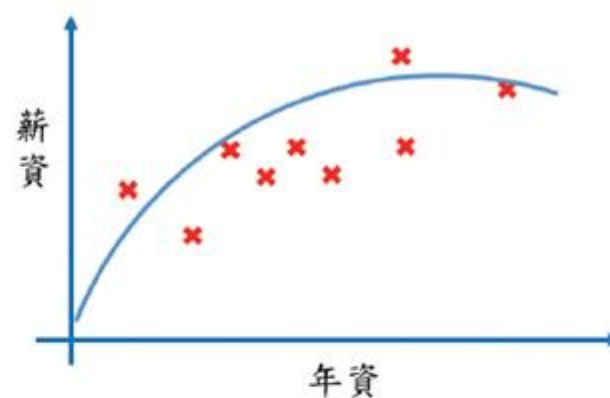
迴歸預測

迴歸 (Regression)

- 非線性函數



(a)二次多項式的函數線

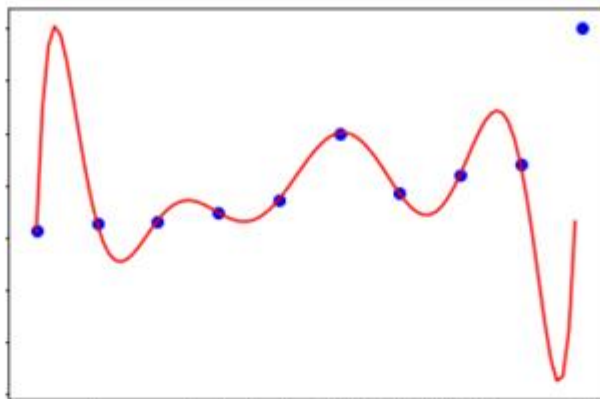


(b)高次多項式的函數線

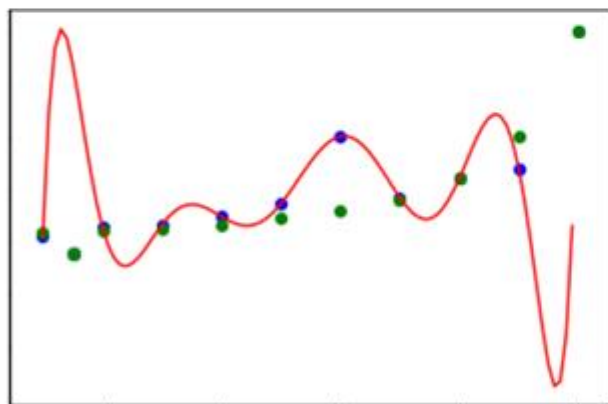
非線性函數

迴歸 (Regression)

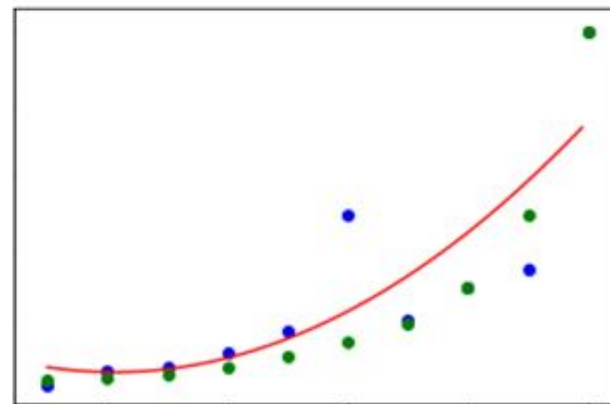
- 預測模型



(a)十次方多項式的模型



(b)預測值在十次方多項式模型的趨勢



(c)預測值在二次方多項式模型的趨勢

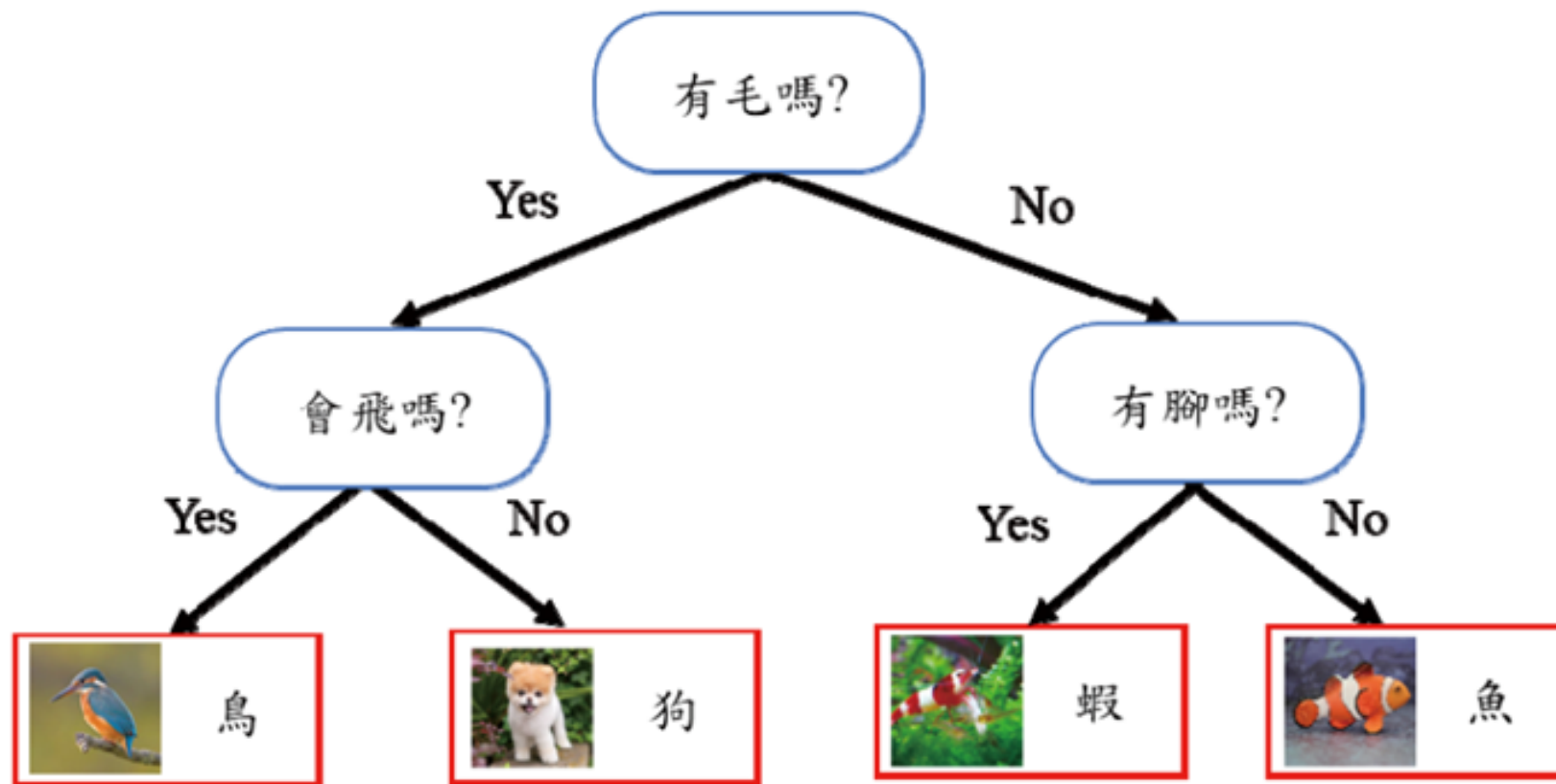
預測模型

迴歸 (Regression)

以下的範例是Boston的房價資料集。

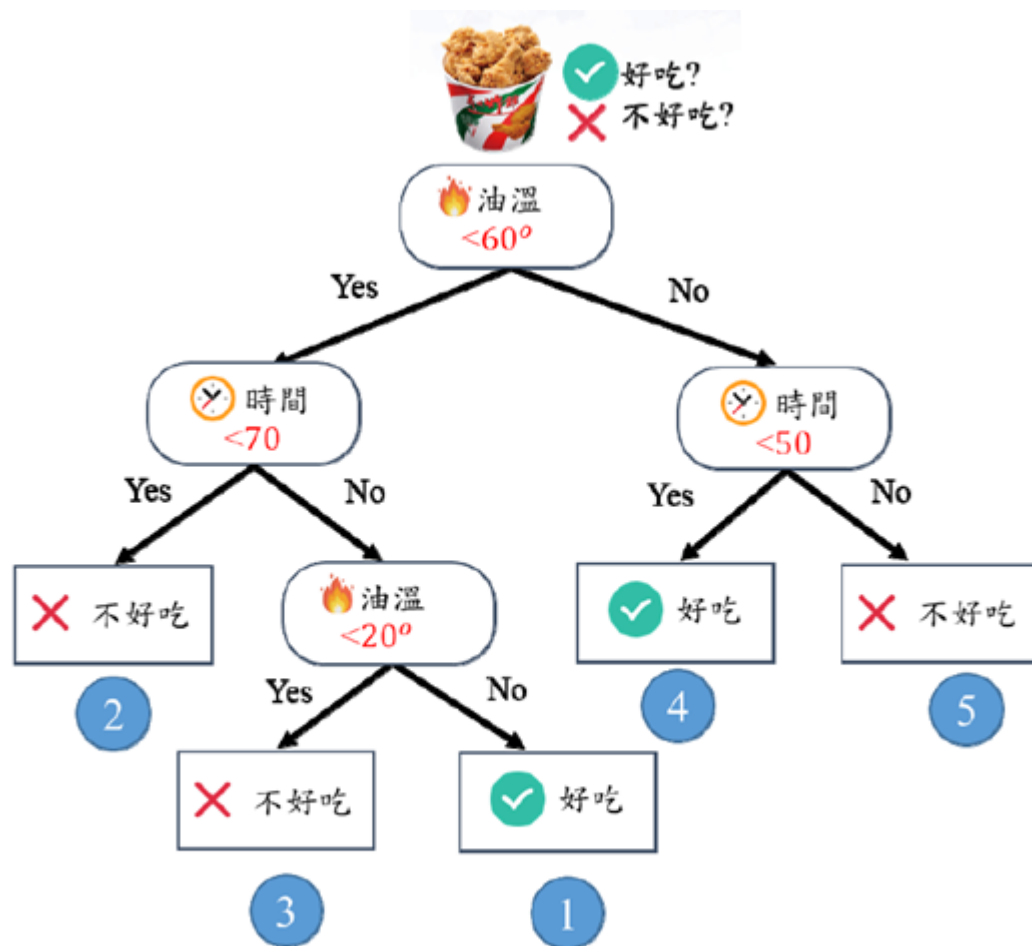
- 共有14個特徵欄位，506筆資料。
- CRIM：人均犯罪率。
- ZN：住宅用地超過25,000平方英尺的比例。
- INDUS：城鎮非零售商用土地的比例。
- CHAS：是否鄰近查爾斯河(1：是、0：否)。
- NOX：一氧化氮濃度。
- RM：住宅的平均房間數。
- AGE：1940年之前建造的自用房屋比例。
- DIS：到波士頓五個中心區域的加權距離。
- RAD：到達高速公路的方便性指標。
- TAX：每萬元的全價值房屋稅。
- PTRATIO：城鎮的師生比例。
- B： $1000 * (B_k - 0.63) * 2$ ，(B_k ：城鎮中黑人的比例)。
- LSTAT：低收入人口的比例。
- MEDV：自住房屋的平均房價(單位：千元美金)。

決策樹 (Decision Tree)



分辨何種動物的決策樹

決策樹 (Decision Tree)



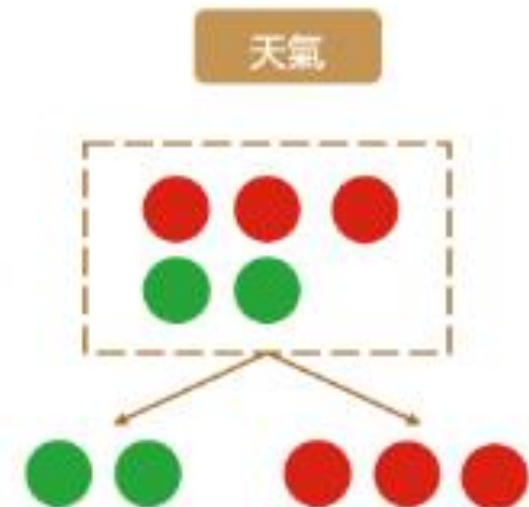
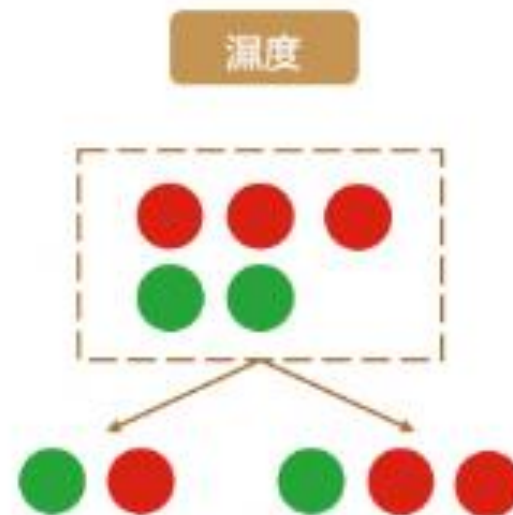
判斷炸雞是否好吃的決策樹

炸雞的數據

編號	油溫	油炸時間	好不好吃
1	50	80	好
2	45	60	不好
3	19	100	不好
4	100	30	好
5	100	70	不好
6	70	30	?

決策樹 (Decision Tree)

- 正常舉行
- 取消舉行



決策樹 (Decision Tree)

- 分割原則-資訊增益(Information gain, 簡稱IG)
 - 熵(Entropy)
 - Gini不純度(Gini Impurity)

資訊增益

獲得的資訊量 原本的資訊量 經由分割後的資訊量

$$IG(D_p, f) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N_p} I(D_j)$$

二元分類資訊增益

獲得的資訊量 原本的資訊量 分割後左邊資訊量 分割後右邊資訊量

$$IG(D_p, f) = I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right})$$

決策樹 (Decision Tree)

- 分割原則-資訊增益(Information gain, 簡稱IG)
- 評估分割資訊量

$$Entropy = - \sum_j p_j \log_2 p_j$$

$$Gini = 1 - \sum_j p_j^2$$

決策樹 (Decision Tree)

- 分割原則-資訊增益(Information gain, 簡稱IG)
- 熵 (Entropy)
 - 計算Information Gain 的一種方法

$$\text{Entropy} = -\sum p_j \log_2 p_j$$

$$\text{Information Gain} = -p * \log_2 p - q * \log_2 q$$

p : 是的機率 q : 否的機率



$$\text{Info}(6, 0) = -\frac{6}{6} \log_2 \left(\frac{6}{6}\right) - \frac{0}{6} \log_2 \left(\frac{0}{6}\right) = 0$$



$$\text{Info}(3, 3) = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1$$

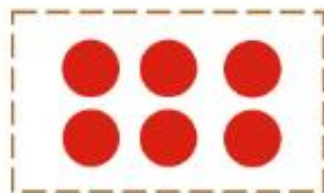
決策樹 (Decision Tree)

- 分割原則-資訊增益(Information gain, 簡稱IG)
- Gini 不純度 (Gini Impurity)
 - 另外一種計算Information Gain 的方法

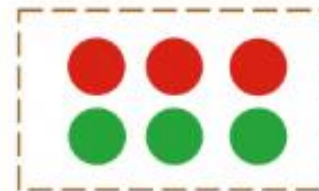
$$\text{Gini} = 1 - \sum p_j^2$$

$$\text{Gini Impurity} = 1 - (p^2 + q^2)$$

p : 是的機率 q : 否的機率



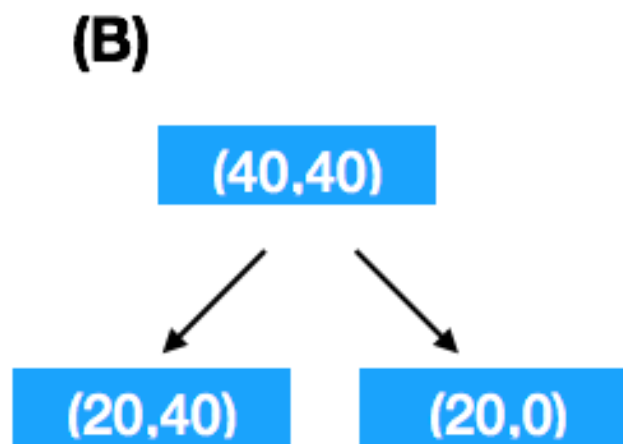
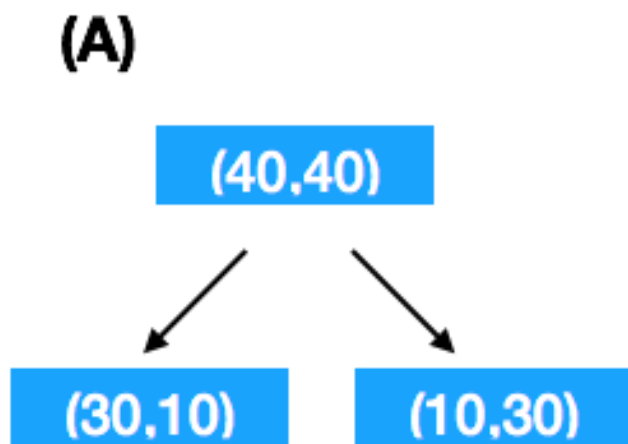
$$\text{Info}(6, 0) = 1 - (1^2 + 0^2) = 0$$



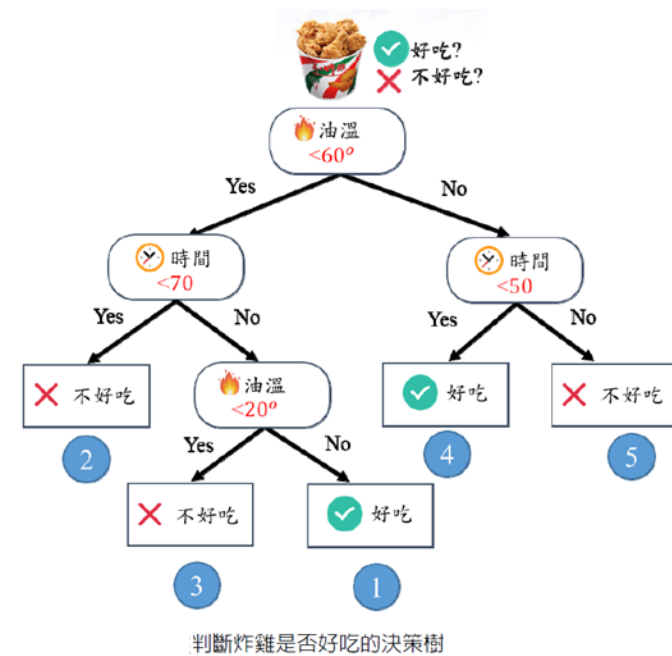
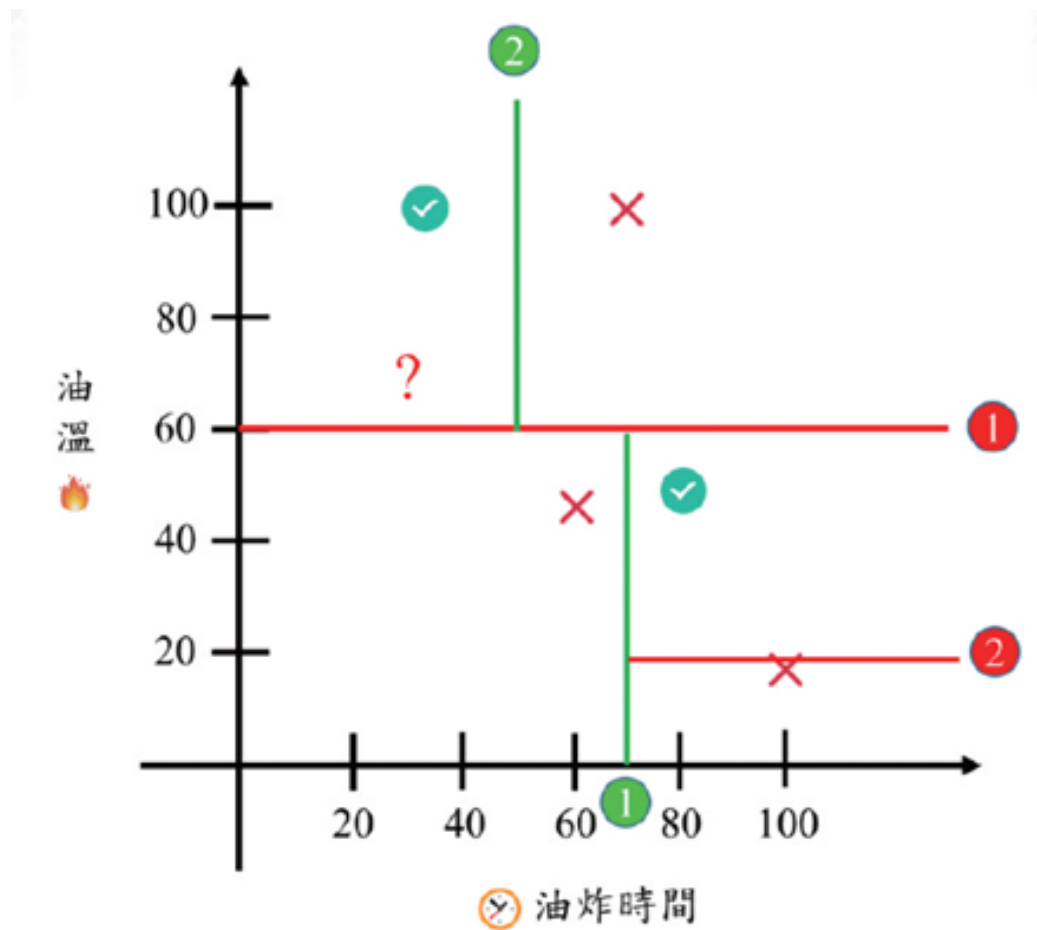
$$\text{Info}(3, 3) = 1 - (0.5^2 + 0.5^2) = 0.5$$

決策樹 (Decision Tree)

- 分割原則-資訊增益(Information gain, 簡稱IG)
- 有80筆資料，有40是1類別、40筆是2類別。使用兩種不同的切割方法A與B



決策樹 (Decision Tree)



決策樹 (Decision Tree)

優點

- 簡單且高度可解釋性
- 低計算時間複雜度
- 每個決策階段都相當的明確清楚
- 幾乎沒有要調整的超參數

缺點

- 模型容易過度擬合
- 當標籤類別種類多時樹會很複雜

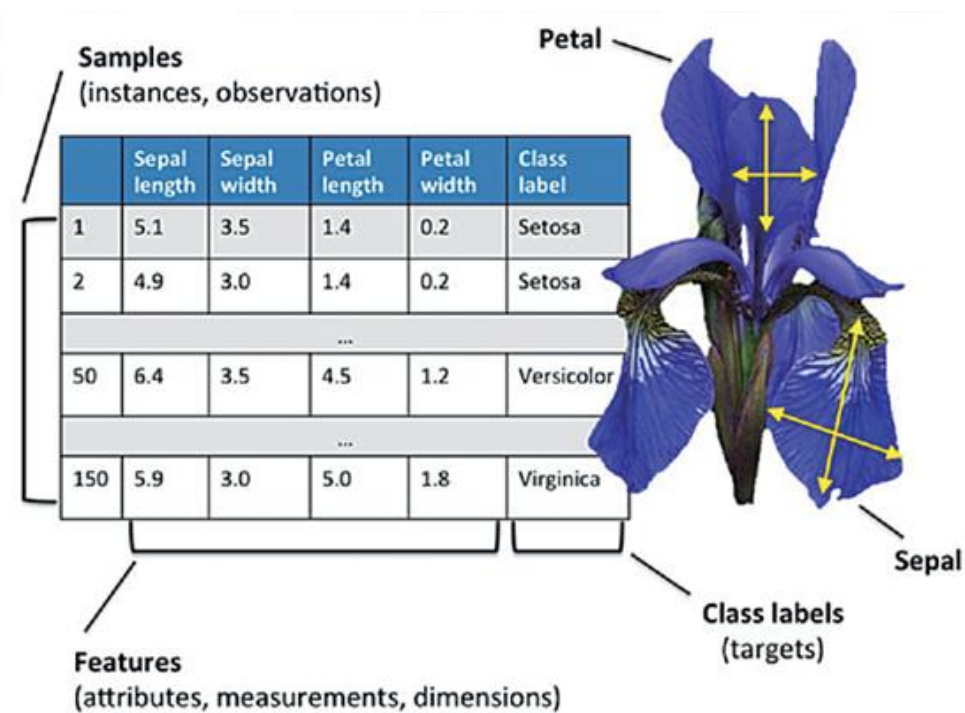
決策樹 (Decision Tree)

範例練習

- 以下的範例是鳶(ㄣㄅ)尾花的分類資料集。
- 共有4個欄位，150筆資料。
 - sepal length (cm)：花萼的長度。
 - sepal width (cm)：花萼的寬度。
 - petal length (cm)：花瓣的長度。
 - petal width (cm)：花瓣的寬度。

決策樹 (Decision Tree)

- Setosa
- Versicolor
- Virginica



鳶尾花的花瓣長寬和花萼長寬

- Entropy: 使用每種特徵分類後的資訊量，資訊量越多就越優先作為決策條件
- Gini Impurity: 指的是分類器的分錯機率，越小代表錯誤機率越小故選為決策的條件。

資料集怎麼來？

Kaggle Dataset scikit-learn

Q Search

Datasets

Explore, analyze, and share quality data. Learn more about data types, creating, and collaborating.

+ New Dataset

Your Work

Q Search datasets

Filters

All datasets

Computer Science

Education

Classification

Computer Vision

NLP

Data Visualization

Pre-Trained Model

Trending Datasets

See All

Porter Delivery Time Estimation

Ranit Sarkar · Updated 3 days ago

Usability 10.0 · 5 MB

1 File (CSV)

15

World Happiness Report 2005-Present

Usama Buttar · Updated 3 days ago

Usability 9.4 · 126 kB

1 File (CSV)

9

Deloitte Media Consumption Survey

Caroline(Yuanmo) Zhu · Updated 11 days ...

Usability 10.0 · 368 kB

2 Files (CSV)

16

Ukraine/Russia Conflict Dataset

Kyle Graupe · Updated 14 hours ago

Usability 9.4 · 5 MB

2 Files (CSV)

3

sklearn.datasets : Datasets ¶

The `sklearn.datasets` module includes utilities to load datasets, including methods to load and fetch popular reference datasets. It also features some artificial data generators.

User guide: See the [Dataset loading utilities](#) section for further details.

Loaders

<code>datasets.clear_data_home</code> <code>((data_home))</code>	Delete all the content of the data home cache.
<code>datasets.dump_svmlight_file</code> <code>(X, y, f[, ...])</code>	Dump the dataset in svmlight / libsvm file format.
<code>datasets.fetch_20newsgroups</code> <code>((data_home, ...))</code>	Load the filenames and data from the 20 newsgroups dataset.
<code>datasets.fetch_20newsgroups_vectorized</code> <code>((...))</code>	Load the 20 newsgroups dataset and transform it into tf-idf vectors.
<code>datasets.fetch_california_housing</code> <code>((...))</code>	Loader for the California housing dataset from StatLib.
<code>datasets.fetch_covtype</code> <code>((data_home, ...))</code>	Load the covtype dataset, downloading it if necessary.
<code>datasets.fetch_kddcup99</code> <code>((subset, data_home, ...))</code>	Load and return the kddcup 99 dataset (classification).
<code>datasets.fetch_lfw_pairs</code> <code>((subset, ...))</code>	Loader for the Labeled Faces in the Wild (LFW) pairs dataset
<code>datasets.fetch_lfw_people</code> <code>((data_home, ...))</code>	Loader for the Labeled Faces in the Wild (LFW) people dataset
<code>datasets.fetch_mldata</code> <code>(dataname[, ...])</code>	Fetch an mldata.org data set
<code>datasets.fetch_olivetti_faces</code> <code>((data_home, ...))</code>	Loader for the Olivetti faces data-set from AT&T.
<code>datasets.fetch_rcv1</code> <code>((data_home, subset, ...))</code>	Load the RCV1 multilabel dataset, downloading it if necessary.
<code>datasets.fetch_species_distributions</code> <code>((...))</code>	Loader for species distribution dataset from Phillips et.
<code>datasets.get_data_home</code> <code>((data_home))</code>	Return the path of the scikit-learn data dir.
<code>datasets.load_boston</code> <code>((return_X_y))</code>	Load and return the boston house-prices dataset (regression).
<code>datasets.load_breast_cancer</code> <code>((return_X_y))</code>	Load sklearn.datasets.breast cancer wisconsin dataset (classification).
<code>datasets.load_diabetes</code> <code>((return_X_y))</code>	Load and return the diabetes dataset (regression).
<code>datasets.load_digits</code> <code>((n_class, return_X_y))</code>	Load and return the digits dataset (classification).
<code>datasets.load_files</code> <code>(container_path[, ...])</code>	Load text files with categories as subfolder names.
<code>datasets.load_iris</code> <code>((return_X_y))</code>	Load and return the iris dataset (classification).
<code>datasets.load_linnerud</code> <code>((return_X_y))</code>	Load and return the linnerud dataset (multivariate regression).
<code>datasets.load_mlcomp</code> <code>(*args, **kwargs)</code>	DEPRECATED: since the http://mlcomp.org/ website will shut down in March 2017, the load_mlcomp function was deprecated in version 0.19 and will be removed in 0.21.
<code>datasets.load_sample_image</code> <code>(image_name)</code>	Load the numpy array of a single sample image
<code>datasets.load_sample_images</code> <code>()</code>	Load sample images for image manipulation.

35

Iris dataset

鳶尾花資料集是非常著名的生物資訊資料集之一，取自美國加州大學歐文分校的機器學習資料庫。

<http://archive.ics.uci.edu/ml/datasets/Iris>

