



機器學習

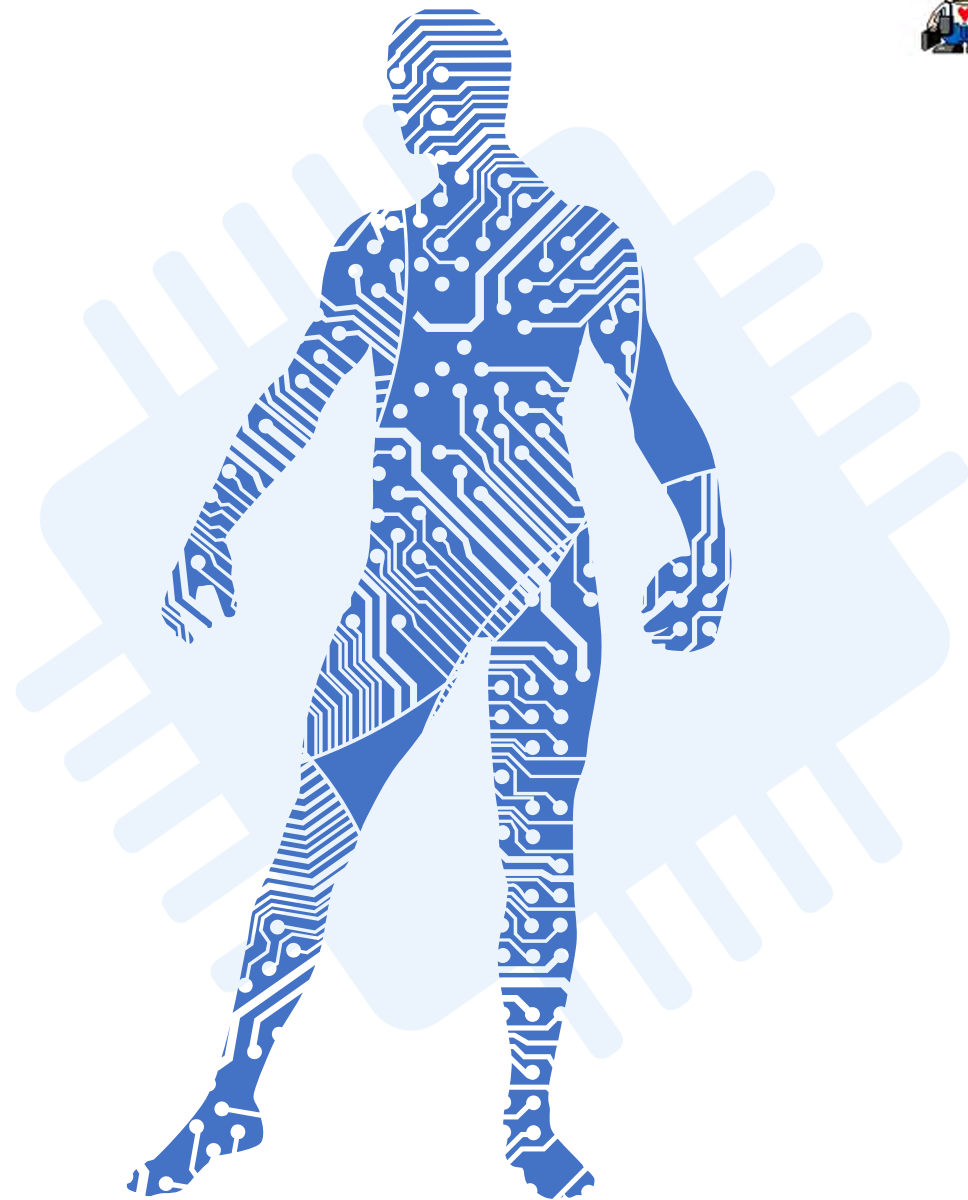
第 14 章 K 平均集群法 (K-Means Clustering)

講師：紀俊男



本章大綱

- 原理解說
- 資料前處理
- 實作 K 平均法
- 將結果視覺化
- 本章總結





AI

原理解説

何謂「集群演算法 (Clustering)」

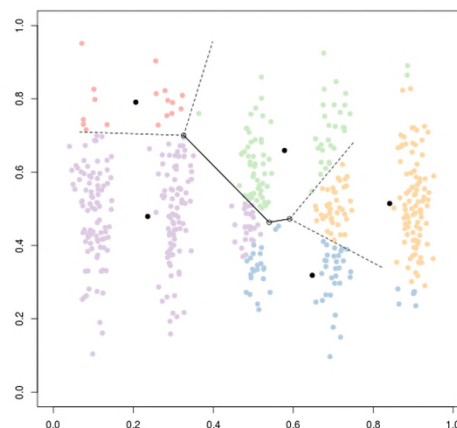


| Genre | Age | Income (k\$) | Spending |
|--------|-----|--------------|----------|
| Male | 19 | 15 | 39 |
| Male | 21 | 15 | 81 |
| Female | 20 | 16 | 6 |
| Female | 23 | 16 | 77 |
| Female | 31 | 17 | 40 |
| Female | 22 | 17 | 76 |
| Female | 35 | 18 | 6 |
| Female | 23 | 18 | 94 |
| Male | 64 | 19 | 3 |

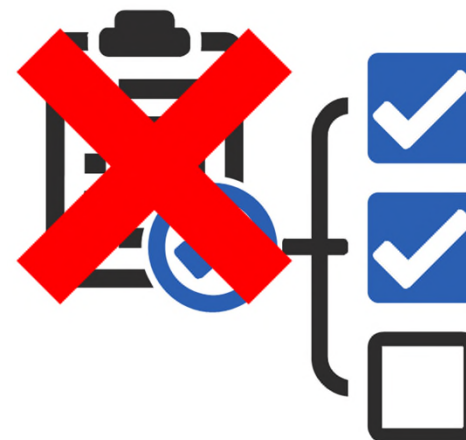
自變數 X

?

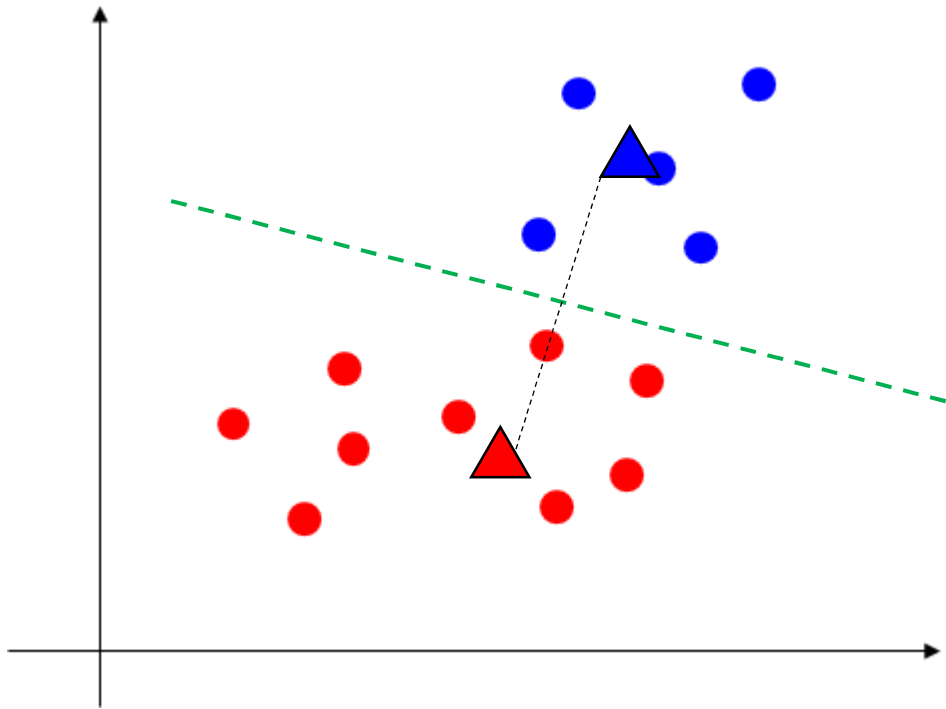
只有自變數 X，
沒有應變數 Y 的問題



利用相似度 (距離)
慢慢找到所屬群體



又稱為「無監督式學習」
(沒有老師、沒有答案)

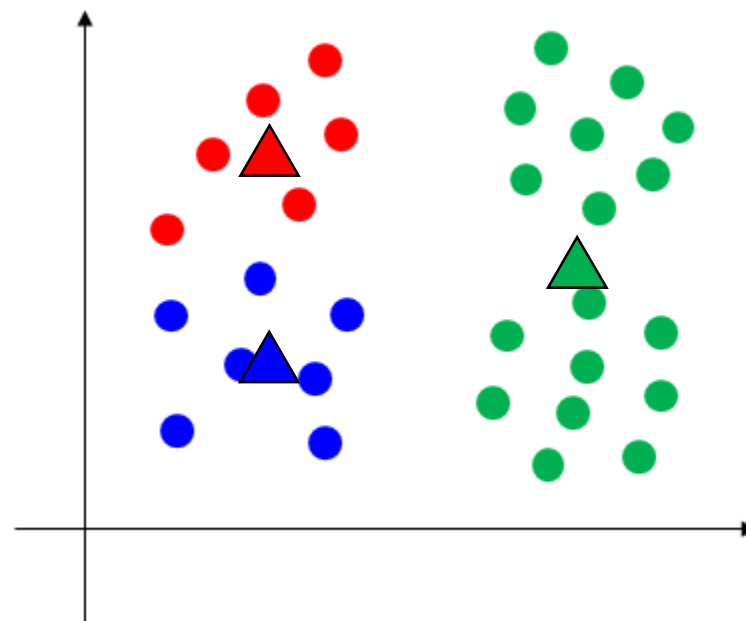
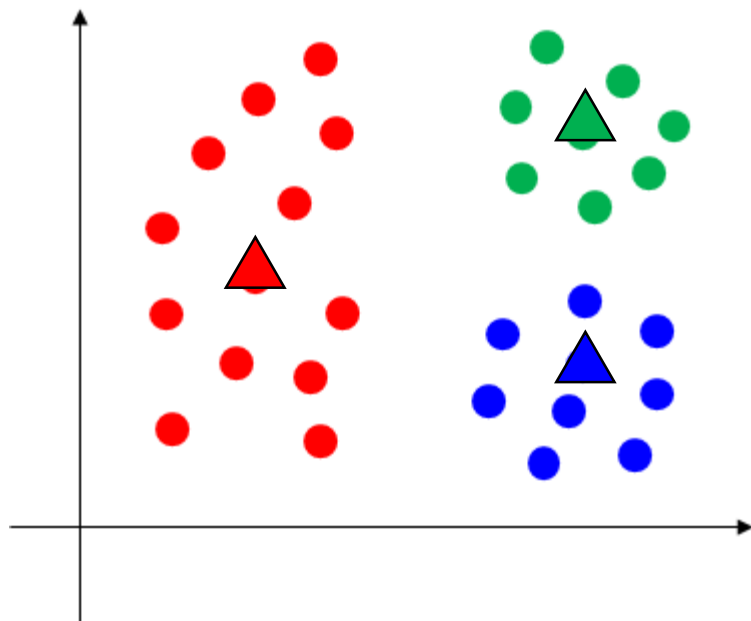


1. 假設要分 K 群，任選 K 個中心點。
2. 對所有樣本點，看距離哪個中心點近，分給該中心點所屬群。
3. 重新計算各群中心點。
4. 重複步驟 (2) ~ (3)。
5. 若樣本點所屬群沒有變化，流程結束。

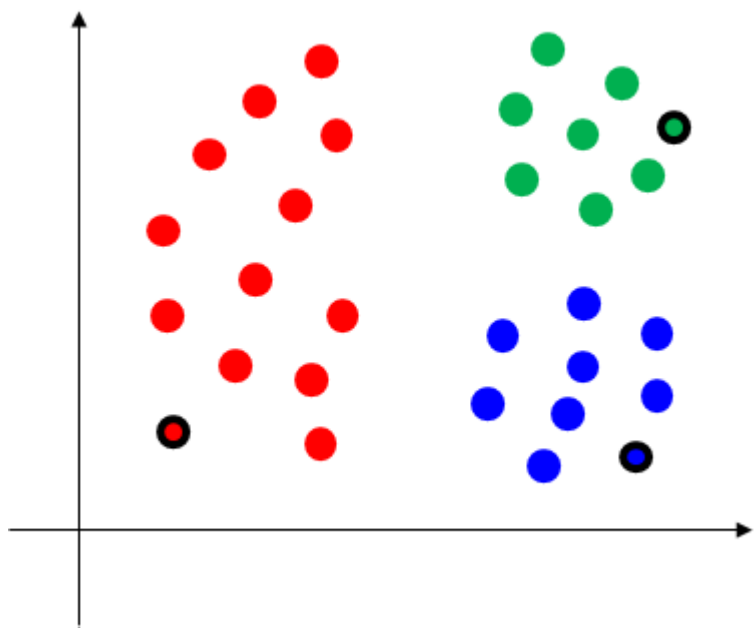
中心點隨機初始化陷阱



- 中心點選擇不同，有可能導致不同的集群結果

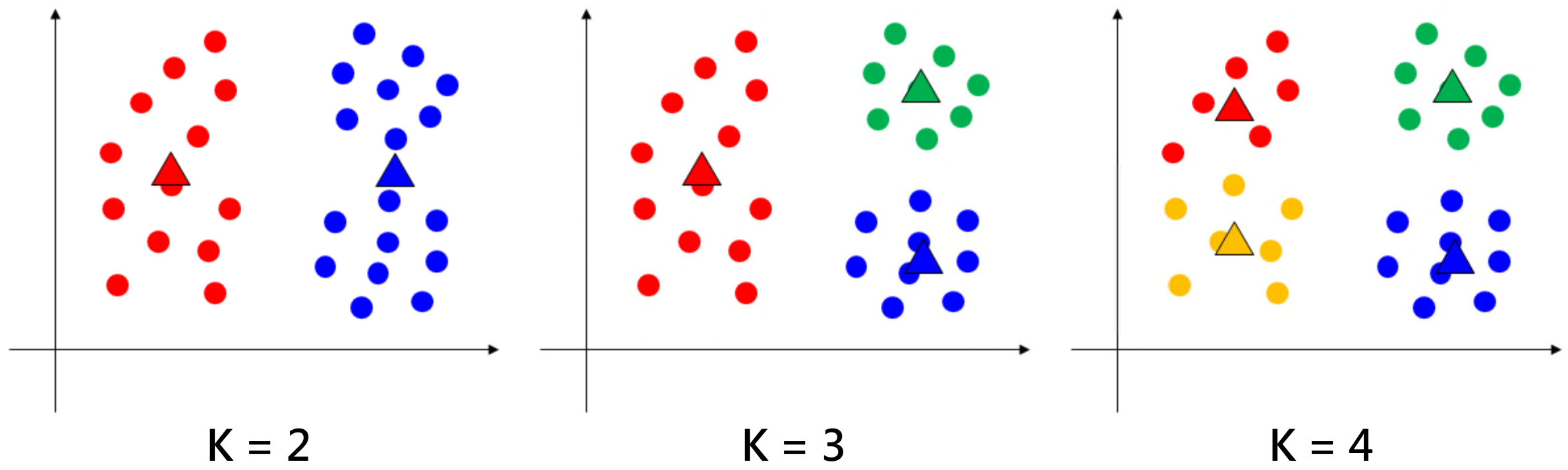


- 解決方法：**K-Means++ 演算法**
 - 挑選「彼此盡量遠離」的初始化中心點



1. 任意選擇 K 個樣本點當中心點。
2. 其它樣本點，看距離哪個**中心點**近，分給該中心點所屬群。
3. 重新計算各群中心點。找各群中，離**舊中心點****最遠**的樣本點，當**新中心點**。
4. 重複步驟（2）～（3）。
5. 若樣本點所屬群沒有變化，**開始一般的** K-Means 流程。

- 起始 K 值不同，最後答案會不一樣

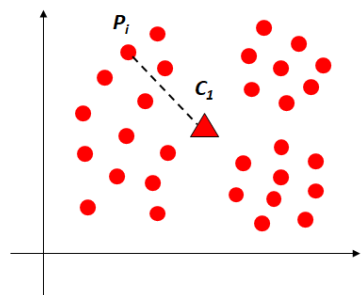


如何選擇正確的 K 值

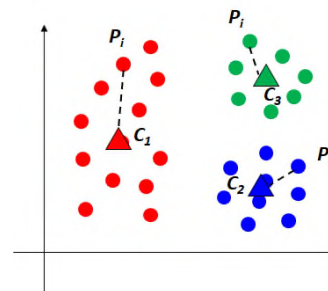


- 解決方法：使用「**群內平方和**」（Within Cluster Sum of Square）

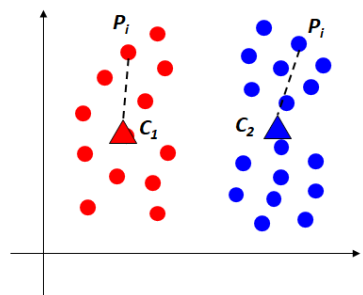
$$WCSS = \sum_{P_i \text{ in Cluster}_1} \text{distance}(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster}_2} \text{distance}(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster}_3} \text{distance}(P_i, C_3)^2 \dots$$



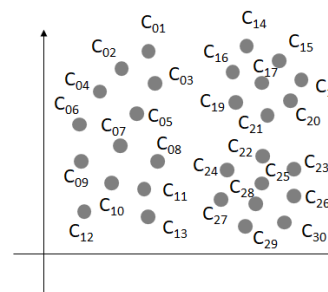
$$WCSS = \sum_{P_i \text{ in Cluster}_1} \text{distance}(P_i, C_1)^2$$



$$WCSS = \sum_{P_i \text{ in Cluster}_1} \text{distance}(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster}_2} \text{distance}(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster}_3} \text{distance}(P_i, C_3)^2$$

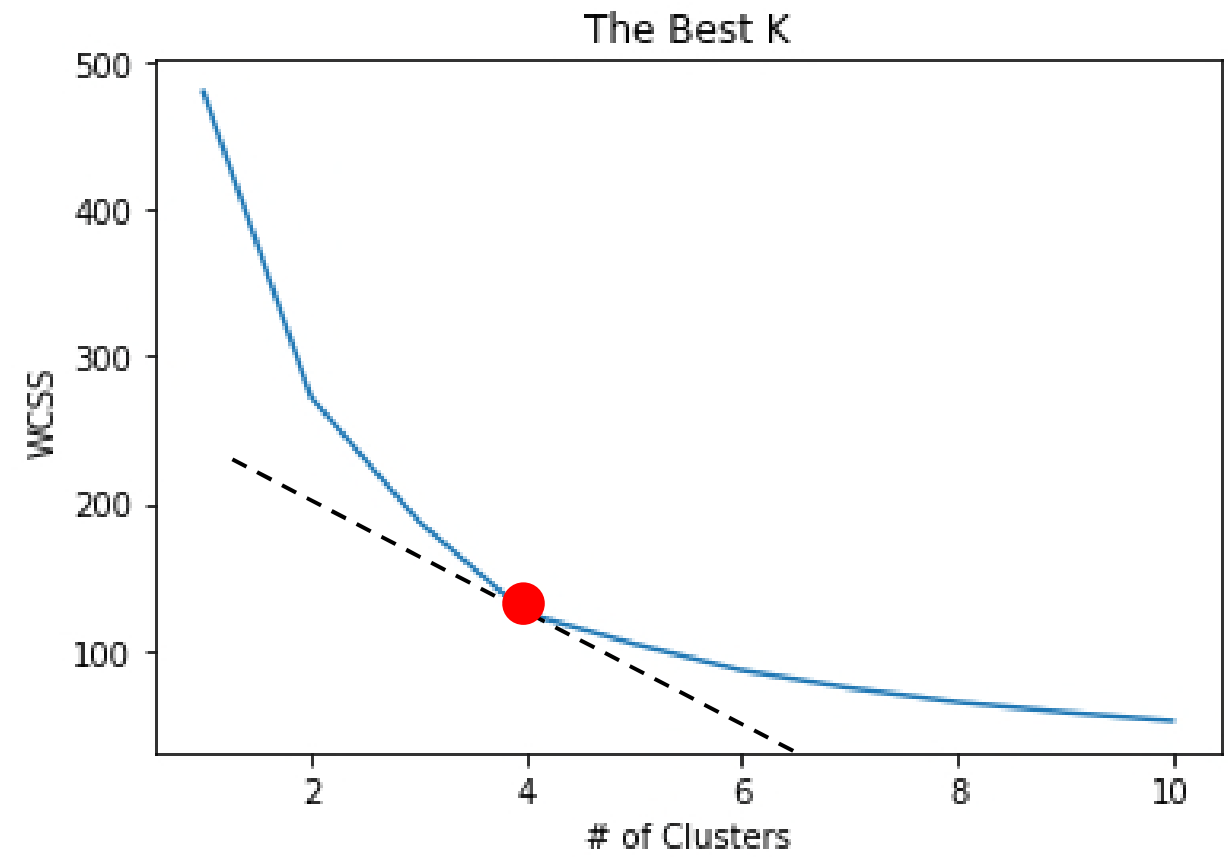


$$WCSS = \sum_{P_i \text{ in Cluster}_1} \text{distance}(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster}_2} \text{distance}(P_i, C_2)^2$$



$$WCSS = 0$$

- 最佳 K 值：下圖紅點處（依照「凸邊形優化 (Convex Optimization) 」定理）

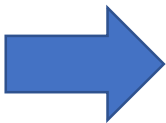




AI

資料前處理

- 依照講師指示，下載並瀏覽資料集



| | A | B | C | D | E |
|----|----|--------|-----|--------------|----------|
| 1 | ID | Gender | Age | Income (k\$) | Spending |
| 2 | 1 | Male | 19 | 15 | 39 |
| 3 | 2 | Male | 21 | 15 | 81 |
| 4 | 3 | Female | 20 | 16 | 6 |
| 5 | 4 | Female | 23 | 16 | 77 |
| 6 | 5 | Female | 31 | 17 | 40 |
| 7 | 6 | Female | 22 | 17 | 76 |
| 8 | 7 | Female | 35 | 18 | 6 |
| 9 | 8 | Female | 23 | 18 | 94 |
| 10 | 9 | Male | 64 | 19 | 3 |

- ID：客戶代號
- Gender：性別
- Age：年齡
- Income(K\$)：收入
- Spending：花費指數 (0~100%)

目的：利用「年齡、收入...」
--> 將客戶分成若干群



撰寫程式碼

```
1 import HappyML.preprocessor as pp
2
3 # Load Dataset
4 dataset = pp.dataset(file="Mall_Customers.csv")
5
6 # Decomposition
7 X = pp.decomposition(dataset, x_columns=[1, 2, 3, 4])
8
9 # One-Hot Encoding
10 X = pp.onehot_encoder(ary=X, columns=[0], remove_trap=True)
11
12 # Feature Scaling (for PCA Feature Selection)
13 X = pp.feature_scaling(fit_ary=X, transform_arys=X)
14
15 # Feature Selection (PCA)
16 from HappyML.preprocessor import PCASelector
17
18 selector = PCASelector()
19 X = selector.fit(x_ary=X, verbose=True, plot=True).transform(x_ary=X)
```

資料前處理流程：

1. 載入資料
2. 取得自變數
3. 處理缺失資料
(無缺失資料)
4. 類別資料數位化
(+ 去除虛擬變數陷阱)
5. 特徵縮放 (PCA 必須)
6. 特徵選擇 (使用 PCA)



- 請撰寫下列程式碼，並予以執行，完成「**資料前處理**」的步驟：

```
1 import HappyML.preprocessor as pp
2
3 # Load Dataset
4 dataset = pp.dataset(file="Mall_Customers.csv")
5
6 # Decomposition
7 X = pp.decomposition(dataset, x_columns=[1, 2, 3, 4])
8
9 # One-Hot Encoding
10 X = pp.onehot_encoder(ary=X, columns=[0], remove_trap=True)
11
12 # Feature Scaling (for PCA Feature Selection)
13 X = pp.feature_scaling(fit_ary=X, transform_arys=X)
14
15 # Feature Selection (PCA)
16 from HappyML.preprocessor import PCASelector
17
18 selector = PCASelector()
19 X = selector.fit(x_ary=X, verbose=True, plot=True).transform(x_ary=X)
```





AI

實作 K 平均法

程式碼

載入必要套件

```
1 from sklearn.cluster import KMeans
2 import time
3
4 kmeans = KMeans(n_clusters=4, init="k-means++", random_state=int(time.time()))
5 Y_pred = kmeans.fit_predict(X)
```

產生物件本身
訓練 & 預測

集群個數 (K 值)

初始中心點演算法 ("random" 或 "k-means++")

執行結果

| | |
|---|---|
| | 0 |
| 0 | 1 |
| 1 | 1 |
| 2 | 2 |
| 3 | 1 |
| 4 | 2 |
| 5 | 1 |
| 6 | 2 |
| 7 | 1 |
| 8 | 2 |
| 9 | 1 |



- 請輸入下列程式碼，並執行看看，是否能取得客戶的**集群結果**：

```
1 from sklearn.cluster import KMeans
2 import time
3
4 kmeans = KMeans(n_clusters=4, init="k-means++", random_state=int(time.time()))
5 Y_pred = kmeans.fit_predict(X)
```



以「標準函式庫」找出最佳的 K 值



• 程式碼

```
1 from sklearn.cluster import KMeans
2 import time
3
4 # Find Best K
5 wcss = []
6 for i in range(1, 11):
7     kmeans = KMeans(n_clusters=i, init="k-means++", random_state=int(time.time()))
8     kmeans.fit(X)
9     wcss.append(kmeans.inertia_)
10
11 # Draw WCSS for each K
12 import matplotlib.pyplot as plt
13
14 plt.plot(range(1, 11), wcss)
15 plt.title("The Best K")
16 plt.xlabel("# of Clusters")
17 plt.ylabel("WCSS")
18 plt.show()
```

儲存 1~10 的 WCSS →

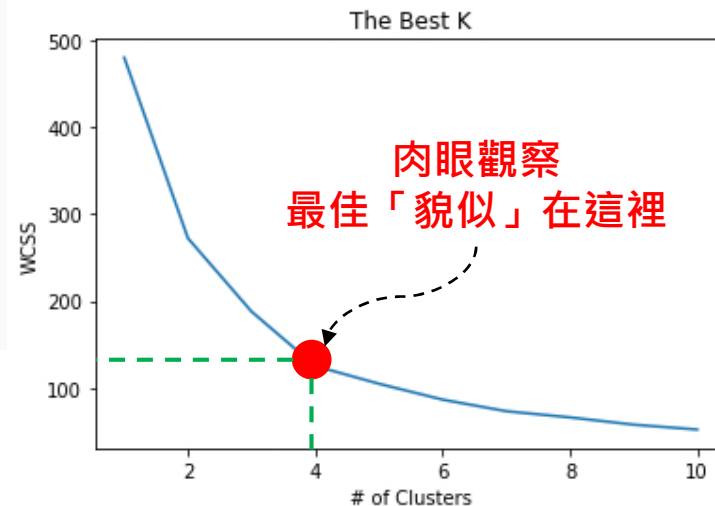
擬合這次 K=i 的結果 {

將這次的 WCSS 存起來 →

計算分成 1~10 群的 WCSS 值

可取得這次的 WCSS 值

繪製所有的 WCSS 值



A 隨堂練習：以「標準函式庫」找出最佳的 K 值



- 請輸入下列程式碼，並執行看看，然後以「**肉眼**」選擇最佳的 K 值：

```
1  from sklearn.cluster import KMeans
2  import time
3
4  # Find Best K
5  wcss = []
6  for i in range(1, 11):
7      kmeans = KMeans(n_clusters=i, init="k-means++", random_state=int(time.time()))
8      kmeans.fit(X)
9      wcss.append(kmeans.inertia_)
10
11 # Draw WCSS for each K
12 import matplotlib.pyplot as plt
13
14 plt.plot(range(1, 11), wcss)
15 plt.title("The Best K")
16 plt.xlabel("# of Clusters")
17 plt.ylabel("WCSS")
18 plt.show()
```



• 程式碼解說（1）：[/HappyML/clustering.py](#)

引入必要套件

類別的成員變數

建構函數

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import time
5
6 from sklearn.cluster import KMeans
7 from sklearn.preprocessing import MinMaxScaler
8
9
10 class KMeansCluster:
11     __cluster = None
12     __best_k = None
13     __max_k = None
14     __strategy = None
15     __random_state = None
16     __centroids = None
17
18     def __init__(self, best_k="auto", max_k=10, random_state=int(time.time())):
19         if type(best_k) is int:
20             self.__strategy = "fixed"
21             self.best_k = best_k
22         else:
23             self.__strategy = "auto"
24             self.best_k = 8
25
26         self.__max_k = max_k
27         self.__random_state = random_state
28
29         self.__cluster = KMeans(n_clusters=self.best_k, init="k-means++", random_state=self.__random_state)
```

- `best_k` = “auto” 或 整數
 - “auto”：交由本類別決定
 - 整數：自行指定
- `max_k` = 整數
 - 最佳 K 值尋找範圍的上限
 - 最佳 K 值 = 1 ~ `max_k`

• 程式碼解說 (2) :

cluster 的
setter & getter

best_k 的
setter & getter

centroids (中心點)
的 getter

/HappyML/clustering.py

```
31 @property
32 def cluster(self):
33     return self.__cluster
34
35 @cluster.setter
36 def cluster(self, cluster):
37     self.__cluster = cluster
38
39 @property
40 def best_k(self):
41     return self.__best_k
42
43 @best_k.setter
44 def best_k(self, best_k):
45     if (type(best_k) is int) and (best_k >= 1):
46         self.__best_k = best_k
47     else:
48         self.__best_k = 1
49
50 @property
51 def centroids(self):
52     return self.__centroids
```


• 程式碼解說 (3) :

/HappyML/clustering.py

計算 1~max_k 的所有 WCSS 值

找尋最佳的 K 值
(X 前進一格、Y 前進不足一格時停止)

verbose = True
印出找到的最佳 K 值

把所有的 WCSS 值繪製出來，
並把最佳 K 值也點描出來

使用最佳 K 值去擬合，
並保存所有中心點資訊

預測各樣本點的所屬群 →

```
54 def fit(self, x_ary, verbose=False, plot=False):
55     if self.__strategy == "auto":
56         wcss = []
57         for i in range(1, self.__max_k+1):
58             kmeans = KMeans(n_clusters=i, init="k-means++", random_state=self.__random_state)
59             kmeans.fit(x_ary)
60             wcss.append(kmeans.inertia_)
61
62         scaler = MinMaxScaler(feature_range=(0, len(wcss)-1))
63         wcss_scaled = scaler.fit_transform(np.array(wcss).reshape(-1, 1)).ravel()
64         for i in range(1, wcss_scaled.shape[0]):
65             if (wcss_scaled[i-1]-wcss_scaled[i]) < 1:
66                 break
67         self.best_k = i
68
69         if verbose:
70             print("The best clusters = {}".format(self.best_k))
71
72         if plot:
73             plt.plot(range(1, len(wcss)+1), wcss, color="blue")
74             plt.scatter(x=self.best_k, y=wcss[self.best_k], color="red")
75             plt.title("The Best Cluster")
76             plt.xlabel("# of Clusters")
77             plt.ylabel("WCSS")
78             plt.show()
79
80         # Fit the Model
81         self.cluster = KMeans(n_clusters=self.best_k, random_state=self.__random_state)
82         self.cluster.fit(x_ary)
83         self.__centroids = self.cluster.cluster_centers_
84
85     return self
86
87 def predict(self, x_ary, y_column="Result"):
88     return pd.DataFrame(self.cluster.predict(x_ary), index=x_ary.index, columns=[y_column])
```


使用「快樂版函式庫」實作



• 呼叫範例

1
2
3
4
5
6
7
8

```
from HappyML.clustering import KMeansCluster  
cluster = KMeansCluster()  
Y_pred = cluster.fit(x_ary=X, verbose=True, plot=True).predict(x_ary=X, y_column="Customer Type")  
# Optional, Attach the Y_pred to Dataset & Save as .CSV file  
dataset = pp.combine(dataset, Y_pred)  
dataset.to_csv("Mall_Customers_answers.csv", index=False)
```

訓練 & 預測
將 Y_pred 結果，附加到 dataset 尾部
將 dataset 輸出成 .CSV 檔

• 執行結果

The best clusters = 4



Y_pred

| Index | Customer Type |
|-------|---------------|
| 0 | 3 |
| 1 | 3 |
| 2 | 1 |
| 3 | 3 |
| 4 | 1 |
| 5 | 3 |
| 6 | 1 |
| 7 | 3 |
| 8 | 1 |
| 9 | 3 |

| ID | Gender | Age | Income (k\$) | Spending (1-100) | Customer Type |
|----|----------|-----|--------------|------------------|---------------|
| 0 | 1 Male | 19 | 15 | 39 | 0 |
| 1 | 2 Male | 21 | 15 | 81 | 0 |
| 2 | 3 Female | 20 | 16 | 6 | 1 |
| 3 | 4 Female | 23 | 16 | 77 | 0 |
| 4 | 5 Female | 31 | 17 | 40 | 1 |
| 5 | 6 Female | 22 | 17 | 76 | 0 |
| 6 | 7 Female | 35 | 18 | 6 | 0 |
| 7 | 8 Female | 23 | 18 | 94 | 0 |

| ID | Gender | Age | Income (k\$) | Spending (1-100) | Customer Type |
|----|--------|-----|--------------|------------------|---------------|
| 1 | Male | 19 | 15 | 39 | 0 |
| 2 | Male | 21 | 15 | 81 | 0 |
| 3 | Female | 20 | 16 | 6 | 3 |
| 4 | Female | 23 | 16 | 77 | 0 |
| 5 | Female | 31 | 17 | 40 | 0 |
| 6 | Female | 22 | 17 | 76 | 0 |
| 7 | Female | 35 | 18 | 6 | 3 |



- 請先將先前使用 **標準函式庫** 製作的 K 平均法**註解**掉。
- 輸入下列程式，並執行看看：

```
1 from HappyML.clustering import KMeansCluster
2
3 cluster = KMeansCluster()
4 Y_pred = cluster.fit(x_ary=X, verbose=True, plot=True).predict(x_ary=X, y_column="Customer Type")
5
6 # Optional, Attach the Y_pred to Dataset & Save as .CSV file
7 dataset = pp.combine(dataset, Y_pred)
8 dataset.to_csv("Mall_Customers_answers.csv", index=False)
```

- 觀察下列輸出：
 - 最佳 **K 值**、**Y_pred**、**dataset** 尾部、輸出的 **.CSV** 檔

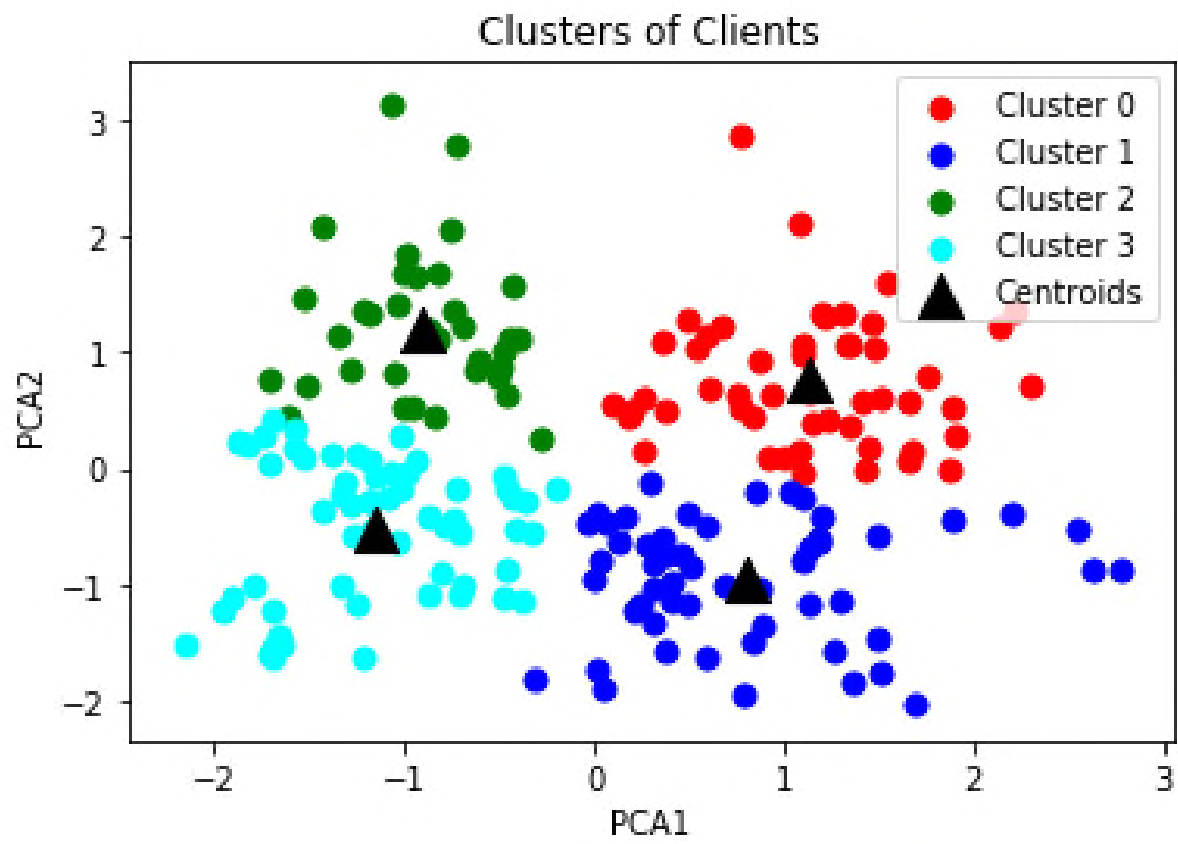




AI

將結果視覺化

- 先看執行結果



• 程式碼

引入必要套件

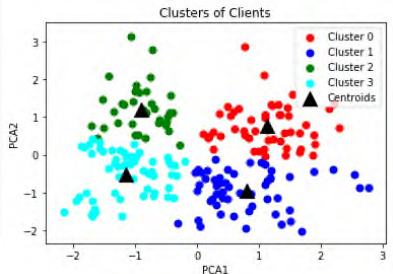
將 DataFrame 轉成 NDAarray

繪製四個群
樣本點所在位置

繪製四個群
的中心點

繪製標題、軸線標籤
與圖示

```
1 import matplotlib.pyplot as plt
2
3 # Draw Samples
4 Y_array = Y_pred.values.ravel()
5 plt.scatter(X.iloc[Y_array==0, 0], X.iloc[Y_array==0, 1], s=50, c="red", label="Cluster 0")
6 plt.scatter(X.iloc[Y_array==1, 0], X.iloc[Y_array==1, 1], s=50, c="blue", label="Cluster 1")
7 plt.scatter(X.iloc[Y_array==2, 0], X.iloc[Y_array==2, 1], s=50, c="green", label="Cluster 2")
8 plt.scatter(X.iloc[Y_array==3, 0], X.iloc[Y_array==3, 1], s=50, c="cyan", label="Cluster 3")
9
10 # Draw Centroids
11 centroids = cluster.centroids
12 plt.scatter(centroids[:, 0], centroids[:, 1], s=200, c="black", marker="^", label="Centroids")
13
14 # Labels & Legends
15 plt.title("Clusters of Clients")
16 plt.xlabel("PCA1")
17 plt.ylabel("PCA2")
18 plt.legend(loc="best")
19 plt.show()
```





- 請輸入下列程式碼，並執行看看，是否能將集群結果**視覺化**：

```
1 import matplotlib.pyplot as plt
2
3 # Draw Samples
4 Y_array = Y_pred.values.ravel()
5 plt.scatter(X.iloc[Y_array==0, 0], X.iloc[Y_array==0, 1], s=50, c="red", label="Cluster 0")
6 plt.scatter(X.iloc[Y_array==1, 0], X.iloc[Y_array==1, 1], s=50, c="blue", label="Cluster 1")
7 plt.scatter(X.iloc[Y_array==2, 0], X.iloc[Y_array==2, 1], s=50, c="green", label="Cluster 2")
8 plt.scatter(X.iloc[Y_array==3, 0], X.iloc[Y_array==3, 1], s=50, c="cyan", label="Cluster 3")
9
10 # Draw Centroids
11 centroids = cluster.centroids
12 plt.scatter(centroids[:, 0], centroids[:, 1], s=200, c="black", marker="^", label="Centroids")
13
14 # Labels & Legends
15 plt.title("Clusters of Clients")
16 plt.xlabel("PCA1")
17 plt.ylabel("PCA2")
18 plt.legend(loc="best")
19 plt.show()
```





- 程式碼解說 (1) :

/HappyML/model_drawer.py

引入調色盤 (Color Map) 物件 → `import matplotlib.cm as cm`

如果傳入的自變數，
不是只有兩個特徵，則退出

將 DataFrame 轉成 NDArray →

找出應變數 Y 究竟有幾種答案 →

```
1 import matplotlib.cm as cm
2
3 def cluster_drawer(x, y, centroids, title="", font='Arial Unicode MS'):
4     # Check for x has only two columns
5     if x.shape[1] != 2:
6         print("ERROR: x must have only two features to draw!!")
7         return None
8
9     # Change y from DataFrame to NDArray
10    y_ndarray = y.values.ravel()
11
12    # Get how many classes in y
13    y_unique = np.unique(y_ndarray)
```




- 程式碼解說（2）：

/HappyML/model_drawer.py

| | | |
|-----------------------|----|--|
| 幾種 Y 就有幾種色 → | 15 | # Iterate all classes in y |
| 繪製樣本點 | 17 | colors = cm.rainbow(np.linspace(0, 1, len(y_unique))) |
| | 18 | for val, col in zip(y_unique, colors): |
| | 19 | plt.scatter(x.iloc[y_ndarray==val, 0], x.iloc[y_ndarray==val, 1], s=50, c=col, label="Cluster {}".format(val)) |
| 繪製中心點 | 20 | # Draw Centroids |
| | 21 | plt.scatter(centroids[:, 0], centroids[:, 1], s=200, c="black", marker="^", label="Centroids") |
| 設定中文字形 | 22 | |
| | 23 | # Labels & Legends |
| | 24 | # for showing Chinese characters |
| | 25 | plt.rcParams['font.sans-serif']=[font] |
| | 26 | plt.rcParams['axes.unicode_minus'] = False |
| 繪製標題、 軸線標籤、 與圖示 | 27 | |
| | 28 | plt.title(title) |
| | 29 | plt.xlabel(x.columns[0]) |
| | 30 | plt.ylabel(x.columns[1]) |
| | 31 | plt.legend(loc="best") |
| | 32 | plt.show() |

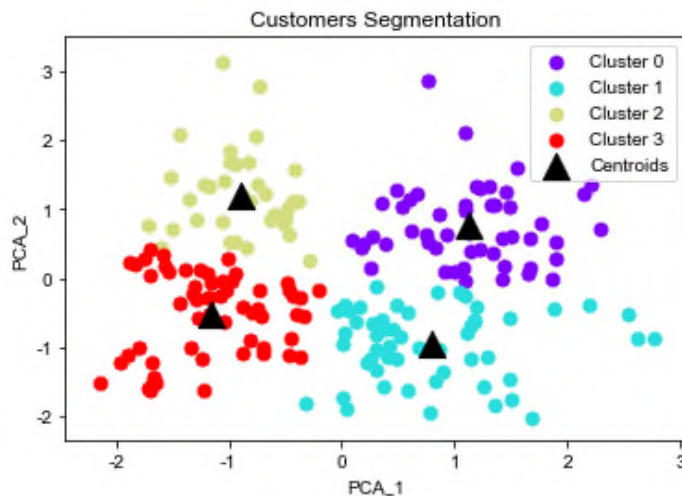
使用「快樂版函式庫」視覺化



• 呼叫範例

```
1 import HappyML.model_drawer as md
2
3 md.cluster_drawer(x=X, y=Y_pred, centroids=cluster.centroids, title="Customers Segmentation")
```

• 執行結果





- 請先將先前使用 **標準函式庫** 視覺化的程式碼**註解**掉。
- 輸入下列程式，並執行看看：

```
1 import HappyML.model_drawer as md
2
3 md.cluster_drawer(x=X, y=Y_pred, centroids=cluster.centroids, title="Customers Segmentation")
```






AI

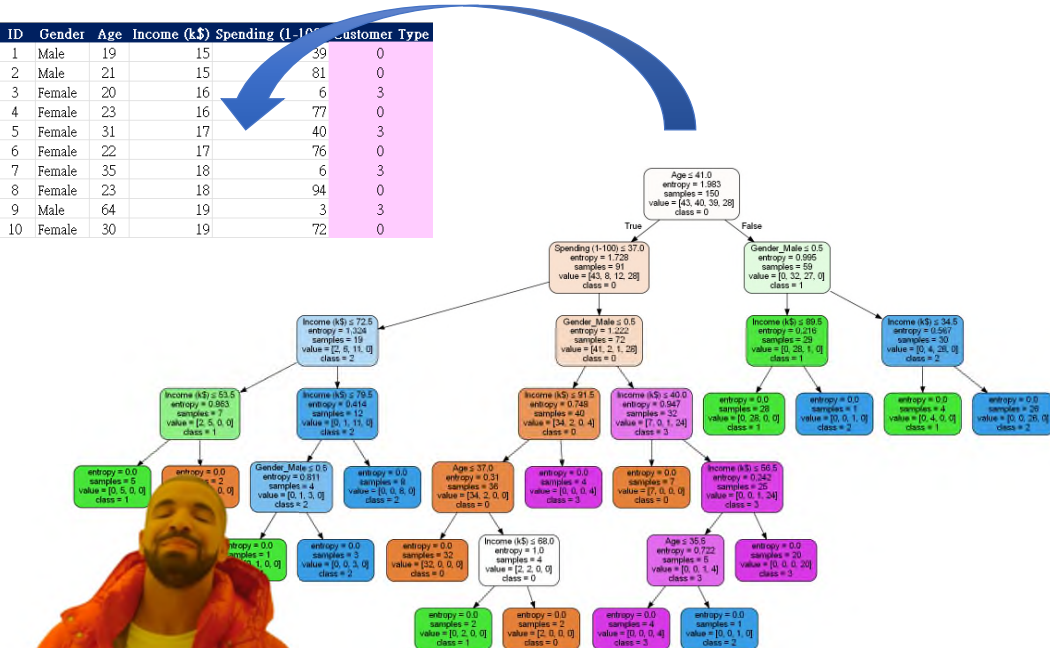

分析集群結果

- 用「肉眼」分析集群原因
- 利用「決策樹」分析集群原因



| ID | Gender | Age | Income (k\$) | Spending (1-100) | Customer Type |
|----|--------|-----|--------------|------------------|---------------|
| 1 | Male | 19 | 15 | 39 | 0 |
| 2 | Male | 21 | 15 | 81 | 0 |
| 3 | Female | 20 | 16 | 6 | 3 |
| 4 | Female | 23 | 16 | 77 | 0 |
| 5 | Female | 31 | 17 | 40 | 3 |
| 6 | Female | 22 | 17 | 76 | 0 |
| 7 | Female | 35 | 18 | 6 | 3 |
| 8 | Female | 23 | 18 | 94 | 0 |
| 9 | Male | 64 | 19 | 3 | 3 |
| 10 | Female | 30 | 19 | 72 | 0 |

| ID | Gender | Age | Income (k\$) | Spending (1-100) | Customer Type |
|----|--------|-----|--------------|------------------|---------------|
| 1 | Male | 19 | 15 | 39 | 0 |
| 2 | Male | 21 | 15 | 81 | 0 |
| 3 | Female | 20 | 16 | 6 | 3 |
| 4 | Female | 23 | 16 | 77 | 0 |
| 5 | Female | 31 | 17 | 40 | 3 |
| 6 | Female | 22 | 17 | 76 | 0 |
| 7 | Female | 35 | 18 | 6 | 3 |
| 8 | Female | 23 | 18 | 94 | 0 |
| 9 | Male | 64 | 19 | 3 | 3 |
| 10 | Female | 30 | 19 | 72 | 0 |



```
graph TD
    Root["Age ≤ 41.0  
entropy = 1.883  
samples = 150  
value = [43, 40, 30, 28]  
class = 0"]
    Root -- True --> Node1["Spending (1-100) ≤ 37.0  
entropy = 1.728  
samples = 91  
value = [43, 8, 12, 28]  
class = 0"]
    Root -- False --> Node2["Gender_Male ≤ 0.5  
entropy = 0.995  
samples = 59  
value = [0, 32, 27, 0]  
class = 1"]
    Node1 --> Node1L["Income (k$) ≤ 72.5  
entropy = 1.324  
samples = 19  
value = [2, 6, 11, 0]  
class = 2"]
    Node1 --> Node1R["Gender_Male ≤ 0.5  
entropy = 1.222  
samples = 72  
value = [41, 2, 1, 29]  
class = 0"]
    Node2 --> Node2L["Income (k$) ≤ 69.5  
entropy = 0.216  
samples = 29  
value = [0, 26, 0, 0]  
class = 1"]
    Node2 --> Node2R["Income (k$) ≤ 34.5  
entropy = 0.567  
samples = 30  
value = [0, 4, 28, 0]  
class = 2"]
    Node1L --> Node1LL["Income (k$) ≤ 53.5  
entropy = 0.853  
samples = 7  
value = [2, 5, 0, 0]  
class = 1"]
    Node1L --> Node1LR["Income (k$) ≤ 79.5  
entropy = 0.811  
samples = 12  
value = [0, 1, 11, 0]  
class = 2"]
    Node1R --> Node1RL["Income (k$) ≤ 91.5  
entropy = 0.745  
samples = 40  
value = [34, 2, 0, 4]  
class = 0"]
    Node1R --> Node1RR["Income (k$) ≤ 40.0  
entropy = 0.947  
samples = 32  
value = [0, 0, 1, 34]  
class = 3"]
    Node2L --> Node2LL["Age ≤ 37.0  
entropy = 1.0  
samples = 4  
value = [0, 0, 0, 0]  
class = 0"]
    Node2L --> Node2LR["Income (k$) ≤ 56.5  
entropy = 0.722  
samples = 25  
value = [0, 0, 1, 34]  
class = 3"]
    Node2R --> Node2RL["Income (k$) ≤ 56.5  
entropy = 0.722  
samples = 25  
value = [0, 0, 1, 34]  
class = 3"]
    Node2R --> Node2RR["Income (k$) ≤ 34.5  
entropy = 0.567  
samples = 30  
value = [0, 4, 28, 0]  
class = 2"]
    Node1LL --> Node1LLL["entropy = 0.0  
samples = 5  
value = [0, 5, 0, 0]  
class = 1"]
    Node1LL --> Node1LLR["entropy = 0.0  
samples = 2  
value = [0, 0, 0, 0]  
class = 0"]
    Node1LR --> Node1LRL["Gender_Male ≤ 0.5  
entropy = 0.811  
samples = 4  
value = [0, 1, 3, 0]  
class = 2"]
    Node1LR --> Node1LRR["entropy = 0.0  
samples = 8  
value = [0, 0, 6, 0]  
class = 2"]
    Node1RL --> Node1RLR["Age ≤ 37.0  
entropy = 1.0  
samples = 36  
value = [34, 2, 0, 0]  
class = 0"]
    Node1RL --> Node1RLR2["Income (k$) ≤ 68.0  
entropy = 1.0  
samples = 4  
value = [2, 2, 0, 0]  
class = 0"]
    Node1RR --> Node1RRR["Age ≤ 35.5  
entropy = 0.722  
samples = 7  
value = [0, 0, 1, 4]  
class = 3"]
    Node1RR --> Node1RRR2["Income (k$) ≤ 56.5  
entropy = 0.722  
samples = 25  
value = [0, 0, 1, 34]  
class = 3"]
    Node2LL --> Node2LLL["entropy = 0.0  
samples = 4  
value = [0, 0, 0, 0]  
class = 0"]
    Node2LL --> Node2LLR["entropy = 0.0  
samples = 2  
value = [0, 2, 0, 0]  
class = 1"]
    Node2LR --> Node2LRL["Income (k$) ≤ 56.5  
entropy = 0.722  
samples = 25  
value = [0, 0, 1, 34]  
class = 3"]
    Node2LR --> Node2LRR["Income (k$) ≤ 34.5  
entropy = 0.567  
samples = 30  
value = [0, 4, 28, 0]  
class = 2"]
    Node2RL --> Node2RLR["Income (k$) ≤ 56.5  
entropy = 0.722  
samples = 25  
value = [0, 0, 1, 34]  
class = 3"]
    Node2RL --> Node2RLR2["Income (k$) ≤ 34.5  
entropy = 0.567  
samples = 30  
value = [0, 4, 28, 0]  
class = 2"]
    Node2RR --> Node2RRR["Income (k$) ≤ 56.5  
entropy = 0.722  
samples = 25  
value = [0, 0, 1, 34]  
class = 3"]
    Node2RR --> Node2RRR2["Income (k$) ≤ 34.5  
entropy = 0.567  
samples = 30  
value = [0, 4, 28, 0]  
class = 2"]
```



決策樹的
前處理程式碼

```
1 # In[] Preprocessing
2 import HappyML.preprocessor as pp
3
4 # Load Dataset
5 dataset = pp.dataset(file="Mall_Customers_answers.csv")
6
7 # Decomposition
8 X, Y = pp.decomposition(dataset, x_columns=[1, 2, 3, 4], y_columns=[5])
9
10 # One-Hot Encoding
11 X = pp.onehot_encoder(ary=X, columns=[0], remove_trap=True)
12
13 # Feature Selection
14 from HappyML.preprocessor import KBestSelector
15 selector = KBestSelector(best_k="auto")
16 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
17
18 # Split Training / Testing Set
19 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
20
21 # In[] Decision Tree
22 from HappyML.classification import DecisionTree
23
24 classifier = DecisionTree()
25 Y_pred = classifier.fit(X_train, Y_train).predict(X_test)
```

決策樹模型
訓練 & 預測



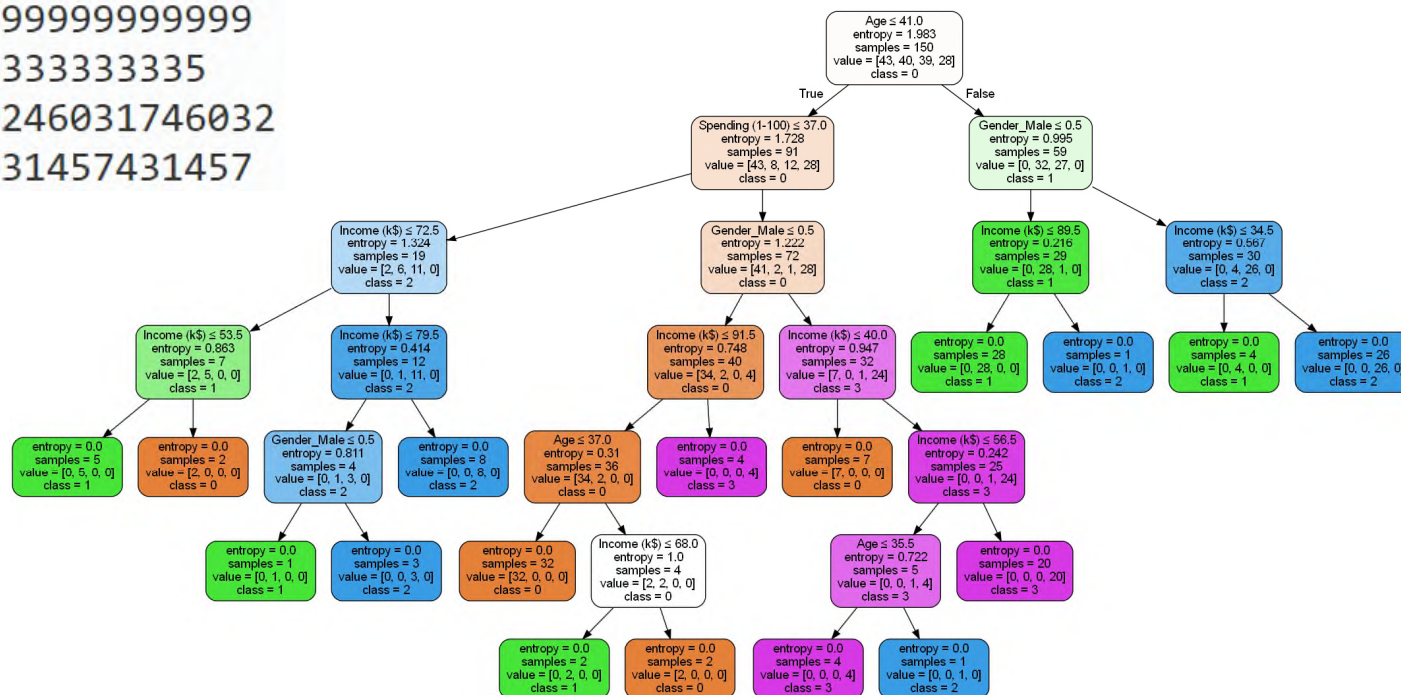
決策樹效能
(對資料集之解釋程度)

決策樹繪製

```
21 # Performance
22 from HappyML.performance import KFoldClassificationPerformance
23
24 K = 10
25 kfp = KFoldClassificationPerformance(x_ary=X, y_ary=Y,
26                                     classifier=classifier.classifier, k_fold=K)
27
28 print("----- Decision Tree Classification -----")
29 print("{} Folds Mean Accuracy: {}".format(K, kfp.accuracy()))
30 print("{} Folds Mean Recall: {}".format(K, kfp.recall()))
31 print("{} Folds Mean Precision: {}".format(K, kfp.precision()))
32 print("{} Folds Mean F1-Score: {}".format(K, kfp.f_score()))
33
34 # Visualization
35 GRAPHVIZ_INSTALL = "C:/Program Files/Graphviz/bin"
36 import HappyML.model_drawer as md
37 from IPython.display import Image, display
38
39 graph = md.tree_drawer(classifier=classifier.classifier,
40                        feature_names=X_test.columns, target_names="0123",
41                        graphviz_bin=GRAPHVIZ_INSTALL)
42 display(Image(graph.create_png()))
```


● 執行結果

```
1 Number of Features Selected: 4
2 ----- Decision Tree Classification -----
3 10 Folds Mean Accuracy: 0.8899999999999999
4 10 Folds Mean Recall: 0.8920833333333333
5 10 Folds Mean Precision: 0.9119246031746032
6 10 Folds Mean F1-Score: 0.8869931457431457
```





- 請撰寫下列程式碼，繪製出決策樹，並且分析集群原因（1）：

```
1 # Decomposition
2 X, Y = pp.decomposition(dataset, x_columns=[1, 2, 3, 4], y_columns=[5])
3
4 # One-Hot Encoding
5 X = pp.onehot_encoder(ary=X, columns=[0], remove_trap=True)
6
7 # Feature Selection
8 from HappyML.preprocessor import KBestSelector
9 selector = KBestSelector(best_k="auto")
10 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
11
12 # Split Training / Testing Set
13 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
14
15 # Decision Tree
16 from HappyML.classification import DecisionTree
17
18 classifier = DecisionTree()
19 Y_pred = classifier.fit(X_train, Y_train).predict(X_test)
```



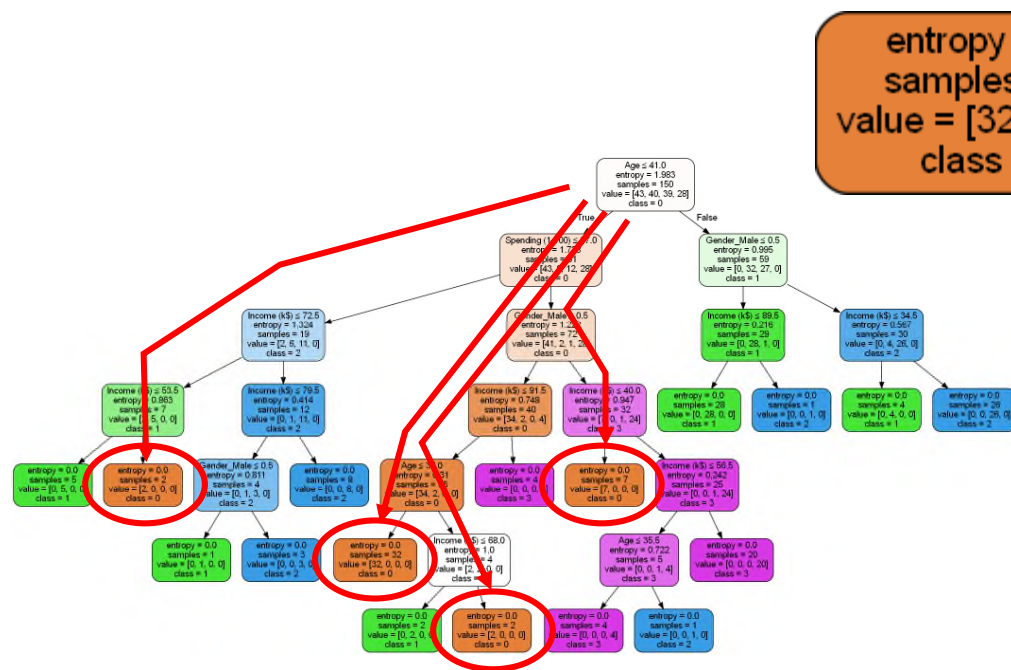


- 請撰寫下列程式碼，繪製出決策樹，並且分析集群原因（2）：

```
21 # Performance
22 from HappyML.performance import KFoldClassificationPerformance
23
24 K = 10
25 kfp = KFoldClassificationPerformance(x_ary=X, y_ary=Y,
26                                     classifier=classifier.classifier, k_fold=K)
27
28 print("----- Decision Tree Classification -----")
29 print("{} Folds Mean Accuracy: {}".format(K, kfp.accuracy()))
30 print("{} Folds Mean Recall: {}".format(K, kfp.recall()))
31 print("{} Folds Mean Precision: {}".format(K, kfp.precision()))
32 print("{} Folds Mean F1-Score: {}".format(K, kfp.f_score()))
33
34 # Visualization
35 GRAPHVIZ_INSTALL = "C:/Program Files/Graphviz/bin"
36 import HappyML.model_drawer as md
37 from IPython.display import Image, display
38
39 graph = md.tree_drawer(classifier=classifier.classifier,
40                        feature_names=X_test.columns, target_names="0123",
41                        graphviz_bin=GRAPHVIZ_INSTALL)
42 display(Image(graph.create_png()))
```



以「決策樹」分析集群結果



entropy = 0.0
samples = 32
value = [32, 0, 0, 0]
class = 0

| Types | Gender | Age | Income (k\$) | Spending (1-100) | Sample |
|-------|--------|------|--------------|------------------|--------|
| 0 | | ~ 41 | 53.5 ~ 72.5 | ~ 37 | 2 |
| 0 | F | ~ 37 | ~ 91.5 | 37 ~ | 32 |
| 0 | F | ~ 41 | 68 ~ 91.5 | 37 ~ | 2 |
| 0 | M | ~ 41 | ~ 40 | 37 ~ | 7 |

Class = 0

女性、年輕、高薪、常客



以「決策樹」分析集群結果



完整分析結果



| Types | Gender | Age | Income (k\$) | Spending (1-100) | Sample | |
|-------|--------|-----------|--------------|------------------|--------|-----------------|
| 0 | | ~41 | 53.5 ~ 72.5 | ~37 | 2 | |
| 0 | F | ~37 | ~91.5 | 37 ~ | 32 | 高薪年輕女性、常客 |
| 0 | F | ~41 | 68 ~ 91.5 | 37 ~ | 2 | |
| 0 | M | ~41 | ~40 | 37 ~ | 7 | |
| 1 | | ~41 | ~53.5 | ~37 | 5 | |
| 1 | F | ~41 | 72.5 ~ 79.5 | ~37 | 1 | |
| 1 | F | ~41 | ~68 | 37 ~ | 2 | |
| 1 | F | 41 ~ | ~89.5 | ~37 | 28 | 本次目標 高薪熟女、稀客 |
| 1 | M | 41 ~ | ~34.5 | | 4 | |
| 2 | M | ~41 | 72.5 ~ 79.5 | ~37 | 3 | |
| 2 | | ~41 | 79.5 ~ | ~37 | 8 | |
| 2 | M | 35.5 ~ 41 | 40 ~ 56 | 37 ~ | 1 | |
| 2 | F | 41 ~ | 89.5 ~ | | 1 | |
| 2 | M | 41 ~ | 34.5 ~ | 37 ~ | 26 | 小資熟男、常客 |
| 3 | F | ~41 | 91.5 ~ | 37 ~ | 4 | |
| 3 | M | ~41 | 56 ~ | 37 ~ | 20 | 中產年輕男性、常客 |

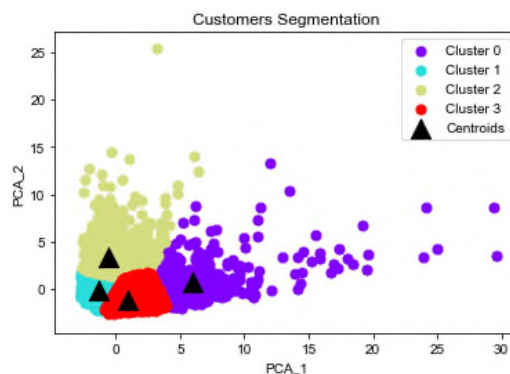
- 資料集說明
 - 請下載資料集 **CreditCards.csv** (資料原始出處：<https://is.gd/sJd4hy>)
 - 欄位名稱解說如下：
 - CUST_ID**：客戶編號
 - BALANCE**：帳戶餘額
 - BALANCE_FREQUENCY**：餘額更新頻率 (0~100%)
 - PURCHASES**：購物金額
 - ONEOFF_PURCHASES**：單筆最大購物金額
 - ...
- 題目要求
 - 請利用 **K-Means 演算法**，將客戶集群。
 - 請將客戶集群結果，以**圖表**繪製出來。
 - 將客戶集群結果，**附加**到原始 dataset 後方。
 - 對附加了集群結果的 dataset，重新做**資料前處理**。
 - 前處理時，使用 **KBestSelector** 挑選「顯著性高」的特徵。
 - 前處理完畢後，使用「**決策樹**」執行「**監督式學習**」，並用 **10-Folds 交叉驗證**，說明你的模型與真實結果有多貼近。
 - 最後繪製出決策樹。說明這樣的集群結果，最重要的**決定因素**是哪個？



課後作業：信用卡客戶集群分析



- 提示
 - 本資料集有「**缺失資料**」，請記得先補齊，再執行運算。
- 輸出結果



```
----- Decision Tree Classification -----  
10 Folds Mean Accuracy: 0.9191101183921022  
10 Folds Mean Recall: 0.8945332130323873  
10 Folds Mean Precision: 0.9055827011793317  
10 Folds Mean F1-Score: 0.8990512580449794
```



- 何謂「**集群演算法 (Clustering)**」
 - 只有**自變數 X**，沒有**應變數 Y**的問題。
 - 利用**相似度 (距離)**，慢慢找到所屬群體。
 - 又稱為「**無監督式學習**」。
- **K-Means 理論**
 - K 平均法**集群流程**。
 - 中心點隨機初始化陷阱解決方法：**K-Means++ 演算法**。
 - 選擇正確的 K 值的方法：取最佳的 **WCSS 值**。
- **K-Means 相關套件**
 - 演算法本身：`sklearn.cluster.KMeans`

