

AI



機器學習

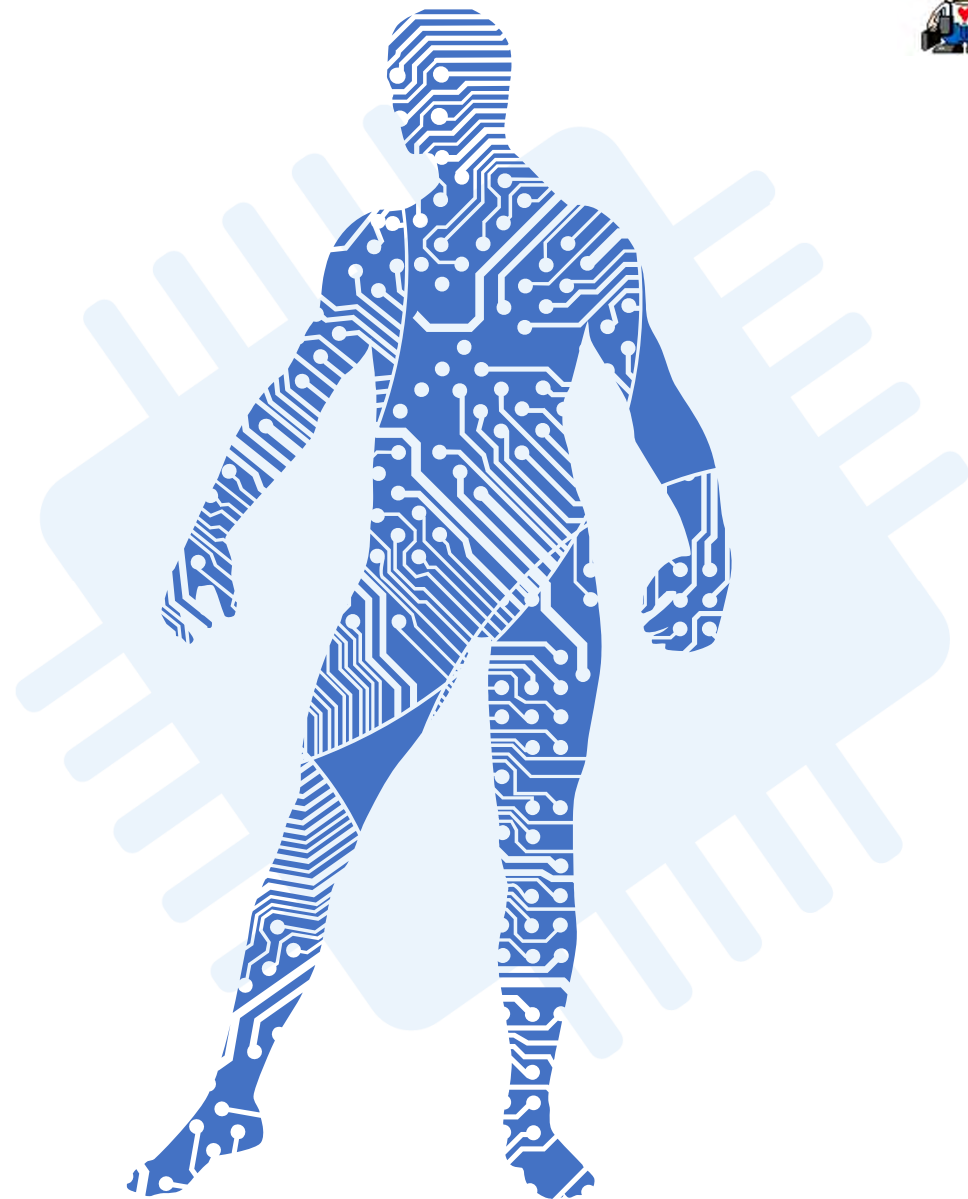
第 13 章 隨機森林 (Random Forest)

講師：紀俊男



本章大綱

- 原理解說
- 資料前處理
- 實作隨機森林
- PCA 降維
- 本章總結





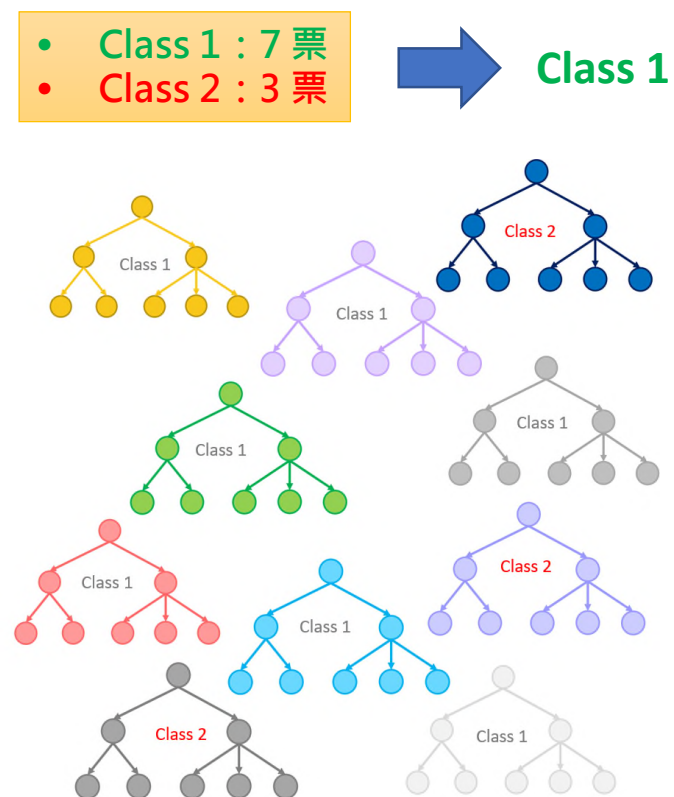
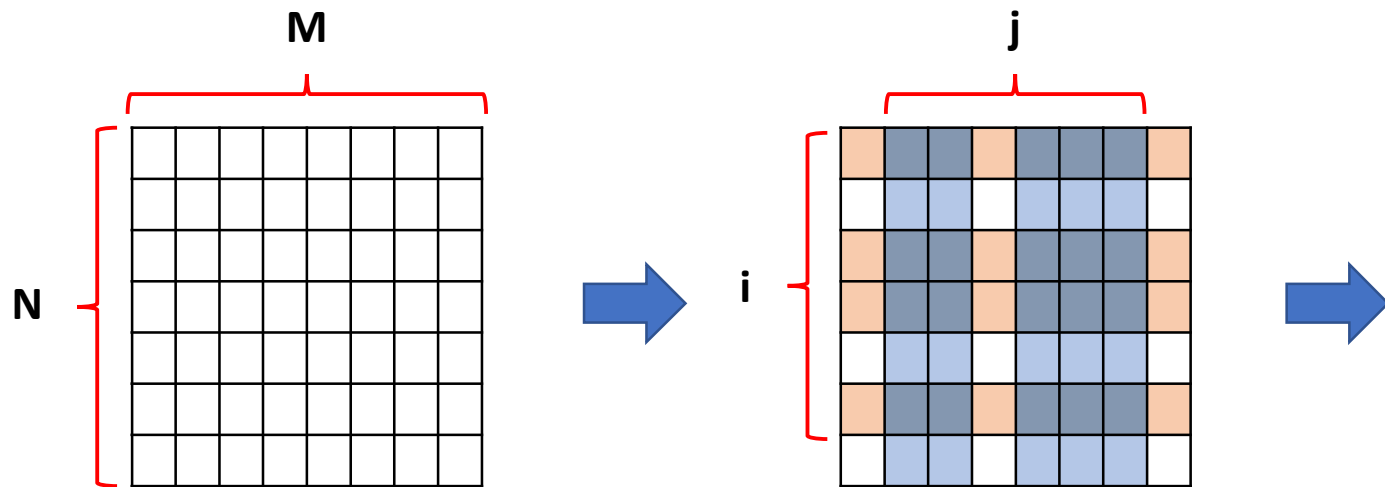
AI

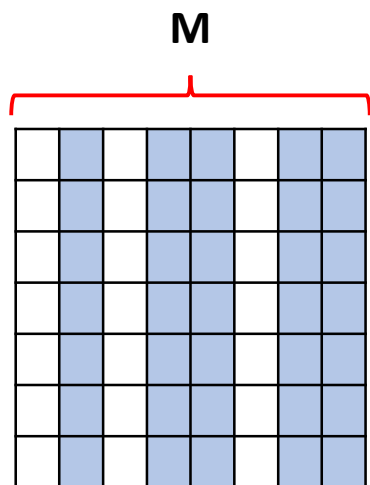
原理解説

使用多棵決策樹「投票」

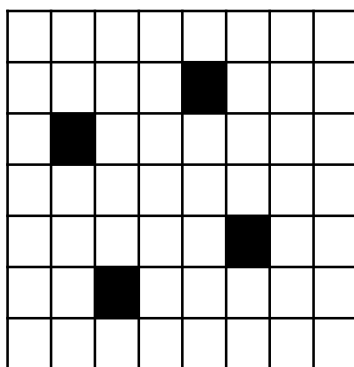


- N : 樣本點個數 M : 自變數個數
- 對於每棵樹
 - C_i^N : 以 i 個樣本點建立一棵決策樹 (取後放回)
 - C_j^M : 以 j 個自變數建立一棵決策樹 (取後放回)
- 對於預測結果, 投票決定





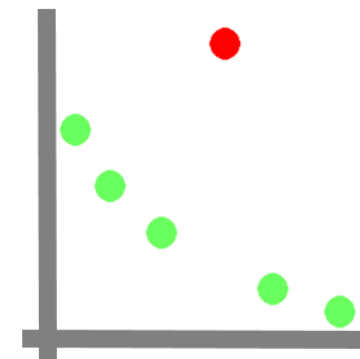
能處理大量、
無相關性的
自變數



對「**缺失資料**」
(**Missing Data**)
耐受度良好



能產生較為
「**不偏 (Un-bias)**」
的訓練集



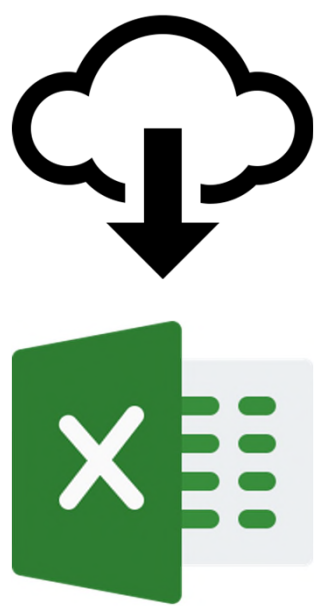
對於「**離群值**」
(**Outliers**)
抗性較佳



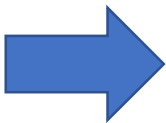
AI

資料前處理

- 依照講師指示，**下載**並**瀏覽**資料集



Wine.csv



| | K | L | M | N |
|----|------|-------|---------|------|
| 1 | Hue | OD280 | Proline | Type |
| 2 | 1.04 | 3.92 | 1065 | 1 |
| 3 | 1.05 | 3.4 | 1050 | 1 |
| 4 | 1.03 | 3.17 | 1185 | 1 |
| 5 | 0.86 | 3.45 | 1480 | 1 |
| 6 | 1.04 | 2.93 | 735 | 1 |
| 7 | 1.05 | 2.85 | 1450 | 1 |
| 8 | 1.02 | 3.58 | 1290 | 1 |
| 9 | 1.06 | 3.58 | 1295 | 1 |
| 10 | 1.08 | 2.85 | 1045 | 1 |

目的：利用「酒精濃度、蘋果酸...」
--> 推算「紅酒等級」

- **Alcohol**：酒精濃度
- **Malic_Acid**：蘋果酸濃度
- **Ash**：粉煤灰濃度
- **Ash_Alcanity**：粉煤灰鹼度
- **Magnesium**：鎂濃度
- **Total_Phenols**：總酚濃度
- **Flavanoids**：黃烷類濃度
- **Nonflavanoid phenols**：非黃酮酚濃度
- **Proanthocyanins**：前花青素濃度
- **Color intensity**：透光度
- **Hue**：色澤
- **OD280/OD315 of diluted wines**：蛋白質吸光度...等指標
- **Proline**：脯氨酸
- **Type**：等級
- 原始來源：<https://is.gd/Kbmb9h>

撰寫程式碼

```
1 import HappyML.preprocessor as pp
2
3 # Load Data, also can be loaded by sklearn.datasets.load_wine()
4 dataset = pp.dataset(file="Wine.csv")
5
6 # Decomposition
7 X, Y = pp.decomposition(dataset, x_columns=[i for i in range(13)], y_columns=[13])
8
9 # By KBestSelector
10 from HappyML.preprocessor import KBestSelector
11 selector = KBestSelector(best_k="auto")
12 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
13
14 # Split Training / Testing Set
15 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
```

- 資料前處理流程：
1. 載入資料

2. 切分自變數、應變數

3. 處理缺失資料
(無缺失資料)

4. 類別資料數位化
(無類別資料)

5. 特徵選擇

6. 切分訓練集、測試集

7. 特徵縮放
(暫無需要)



- 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 Data, also can be loaded by sklearn.datasets.load_wine()
4 dataset = pp.dataset(file="Wine.csv")
5
6 # Decomposition
7 X, Y = pp.decomposition(dataset, x_columns=[i for i in range(13)], y_columns=[13])
8
9 # By KBestSelector
10 from HappyML.preprocessor import KBestSelector
11 selector = KBestSelector(best_k="auto")
12 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
13
14 # Split Training / Testing Set
15 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
```





AI

實作隨機森林

使用「標準函式庫」實作



程式碼

載入必要套件

```
1 from sklearn.ensemble import RandomForestClassifier
2 import time
3
```

決策樹個數
(可用網格搜尋推估) 可選 "entropy" (ID3)
或者 "gini" (CART)

產生物件本身
訓練
預測

```
4 → classifier = RandomForestClassifier(n_estimators=10, criterion="entropy", random_state=int(time.time()))
5 → classifier.fit(X_train, Y_train.values.ravel())
6 → Y_pred = classifier.predict(X_test)
```

執行結果

| Y_test | Index | Type | Y_pred |
|--------|-------|------|--------|
| | 15 | 1 | 0 |
| | 150 | 3 | 1 |
| | 9 | 1 | 3 |
| | 129 | 2 | 1 |
| | 36 | 1 | 2 |
| | 131 | 3 | 3 |
| | 134 | 3 | 1 |
| | 100 | 2 | 3 |
| | 62 | 2 | 2 |
| | 96 | 2 | 2 |

- 請撰寫下列程式碼，並執行之：

```
1 from sklearn.ensemble import RandomForestClassifier
2 import time
3
4 classifier = RandomForestClassifier(n_estimators=10, criterion="entropy", random_state=int(time.time()))
5 classifier.fit(X_train, Y_train.values.ravel())
6 Y_pred = classifier.predict(X_test)
```

- 執行完畢後，請比較 **Y_test (真實值)** 與 **Y_pred (預測值)** 的差異。



• 程式碼解說（1）：

/HappyML/classification.py

引入必要套件

```
1 { from sklearn.ensemble import RandomForestClassifier
2   import time
3
```

類別的成員變數

```
4 class RandomForest:
5     __classifier = None
6     __n_estimators = None
7     __criterion = None
8     __y_columns = None
9
```

幾棵決策樹

演算法

建構函數

```
10 def __init__(self, n_estimators=10, criterion="entropy"):
11     self.__n_estimators = n_estimators
12     self.__criterion = criterion
13     self.__classifier = RandomForestClassifier(n_estimators=self.__n_estimators,
14                                              criterion=self.__criterion, random_state=int(time.time()))
```

• 程式碼解說（2）：

/HappyML/classification.py

classifier 的
setter & getter

n_estimators
的 getter

訓練

預測

```
16  @property
17  def classifier(self):
18      return self.__classifier
19
20  @classifier.setter
21  def classifier(self, classifier):
22      self.__classifier = classifier
23
24  @property
25  def n_estimators(self):
26      return self.__n_estimators
27
28  def fit(self, x_train, y_train):
29      self.classifier.fit(x_train, y_train.values.ravel())
30      self.__y_columns = y_train.columns ← 保存 Y_train 的欄位名稱
31      return self
32
33  def predict(self, x_test):
34      return pd.DataFrame(self.classifier.predict(x_test), index=x_test.index, columns=self.__y_columns)
```

↓ 把預測出來的 Y_pred 重新包裝回 DataFrame 再傳回

• 呼叫範例

```
1 from HappyML.classification import RandomForest
2
3 classifier = RandomForest(n_estimators=10, criterion="entropy")
4 Y_pred = classifier.fit(X_train, Y_train).predict(X_test)
```

• 執行結果

Y_test

| Index | Type |
|-------|------|
| 15 | 1 |
| 150 | 3 |
| 9 | 1 |
| 129 | 2 |
| 36 | 1 |
| 131 | 3 |
| 134 | 3 |
| 100 | 2 |
| 62 | 2 |
| 96 | 2 |

↔

Y_pred

| Index | Type |
|-------|------|
| 15 | 1 |
| 150 | 3 |
| 9 | 1 |
| 129 | 2 |
| 36 | 1 |
| 131 | 3 |
| 134 | 2 |
| 100 | 2 |
| 62 | 2 |
| 96 | 2 |



隨堂練習：使用「快樂版函式庫」實作



- 請撰寫下列程式碼，並執行之：

```
1 from HappyML.classification import RandomForest
2
3 classifier = RandomForest(n_estimators=10, criterion="entropy")
4 Y_pred = classifier.fit(X_train, Y_train).predict(X_test)
```

- 執行完畢後，請比較 Y_{test} (真實值) 與 Y_{pred} (預測值) 的差異。



- 程式碼解說

```
1 from HappyML.performance import KFoldClassificationPerformance
2
3 K = 10
4 kfp = KFoldClassificationPerformance(x_ary=X, y_ary=Y, classifier=classifier.classifier, k_fold=K)
5
6 print("----- Random Forest Classification -----")
7 print("{} Folds Mean Accuracy: {}".format(K, kfp.accuracy()))
8 print("{} Folds Mean Recall: {}".format(K, kfp.recall()))
9 print("{} Folds Mean Precision: {}".format(K, kfp.precision()))
10 print("{} Folds Mean F1-Score: {}".format(K, kfp.f_score()))
```

- 執行結果

```
----- Random Forest Classification -----
10 Folds Mean Accuracy: 0.9502923976608187
10 Folds Mean Recall: 0.950436507936508
10 Folds Mean Precision: 0.9595899470899472
10 Folds Mean F1-Score: 0.9497772597772597
```



- 請撰寫下列程式碼，計算隨機森林的效能：

```
1 from HappyML.performance import KFoldClassificationPerformance
2
3 K = 10
4 kfp = KFoldClassificationPerformance(x_ary=X, y_ary=Y, classifier=classifier.classifier, k_fold=K)
5
6 print("----- Random Forest Classification -----")
7 print("{} Folds Mean Accuracy: {}".format(K, kfp.accuracy()))
8 print("{} Folds Mean Recall: {}".format(K, kfp.recall()))
9 print("{} Folds Mean Precision: {}".format(K, kfp.precision()))
10 print("{} Folds Mean F1-Score: {}".format(K, kfp.f_score()))
```

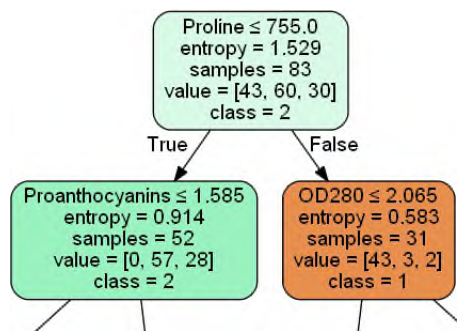


• 程式碼解說

```
1 GRAPHVIZ_INSTALL = "C:/Program Files/Graphviz/bin"
2
3 import HappyML.model_drawer as md
4 from IPython.display import Image, display
5
6 clfr = classifier.classifier.estimators_[0]
7 graph = md.tree_drawer(classifier=clfr, feature_names=X_test.columns,
8                        target_names="123", graphviz_bin=GRAPHVIZ_INSTALL)
9 display(Image(graph.create_png()))
```

範例：只顯示一棵樹
(可以用迴圈顯示所有樹)

• 執行結果





隨堂練習：隨機森林視覺化



- 請撰寫下列程式碼，將隨機森林視覺化：

```
1 GRAPHVIZ_INSTALL = "C:/Program Files/Graphviz/bin"
2
3 import HappyML.model_drawer as md
4 from IPython.display import Image, display
5
6 clfr = classifier.classifier.estimators_[0]
7 graph = md.tree_drawer(classifier=clfr, feature_names=X_test.columns,
8                        target_names="123", graphviz_bin=GRAPHVIZ_INSTALL)
9 display(Image(graph.create_png()))
```



隨機森林「決定邊界」視覺化



• 程式碼解說

前處理程式碼

```
1 import HappyML.preprocessor as pp
2
3 ...
4
5 # By KBestSelector
6 from HappyML.preprocessor import KBestSelector
7 selector = KBestSelector(best_k=2)
8 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
9
10 # Split Training / Testing Set
11 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
12
13 # Feature Scaling
14 X_train, X_test = pp.feature_scaling(fit_ary=X_train, transform_arys=(X_train, X_test))
```

best_k = 2
壓制在兩個特徵之下

一定要記得做「特徵縮放」
(維持各維度資料的比例尺一致)

圖形程式碼

```
1 import HappyML.model_drawer as md
2
3 md.classify_result(x=X_train, y=Y_train, classifier=classifier.classifier,
4                   fg_color=("orange", "blue", "white"), bg_color=("red", "green", "black"),
5                   title="訓練集 vs. 隨機森林模型", font="DFKai-sb")
6 md.classify_result(x=X_test, y=Y_test, classifier=classifier.classifier,
7                   fg_color=("orange", "blue", "white"), bg_color=("red", "green", "black"),
8                   title="測試集 vs. 隨機森林模型", font="DFKai-sb")
```

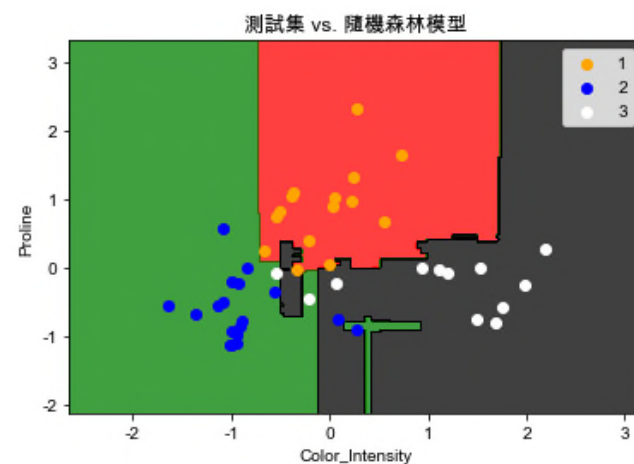
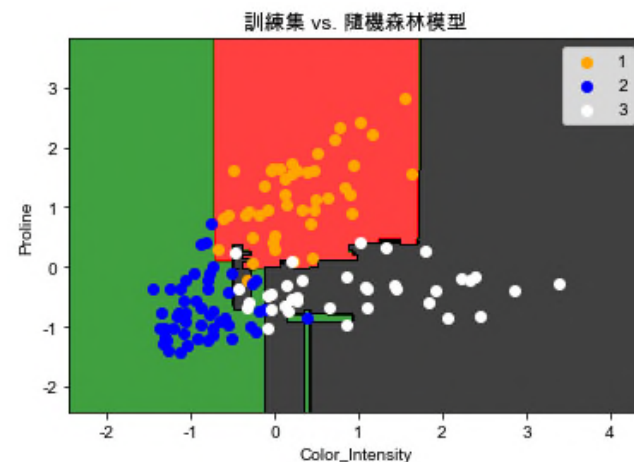
隨機森林「決定邊界」視覺化



- 執行結果

```
Number of Features Selected: 2
----- Random Forest Classification -----
10 Folds Mean Accuracy: 0.8783281733746129
10 Folds Mean Recall: 0.8668253968253969
10 Folds Mean Precision: 0.890489417989418
10 Folds Mean F1-Score: 0.8635534889946654
```

請各位記住這些數字，
以便稍後與PCA對比之用。



如「巧拼板」般的決定邊界



- 請先將 `best_k` 改為 `=2` :
 - `selector = KBestSelector(best_k=2)`
- 再將「特徵縮放」程式碼加上：
 - `X_train, X_test = pp.feature_scaling(fit_ary=X_train, transform_arys=(X_train, X_test))`
- 撰寫下列「決定邊界」視覺化的程式碼，並將整個程式執行一遍：

```
1 import HappyML.model_drawer as md
2
3 md.classify_result(x=X_train, y=Y_train, classifier=classifier.classifier,
4                   fg_color=("orange", "blue", "white"), bg_color=("red", "green", "black"),
5                   title="訓練集 vs. 隨機森林模型", font="DFKai-sb")
6 md.classify_result(x=X_test, y=Y_test, classifier=classifier.classifier,
7                   fg_color=("orange", "blue", "white"), bg_color=("red", "green", "black"),
8                   title="測試集 vs. 隨機森林模型", font="DFKai-sb")
```

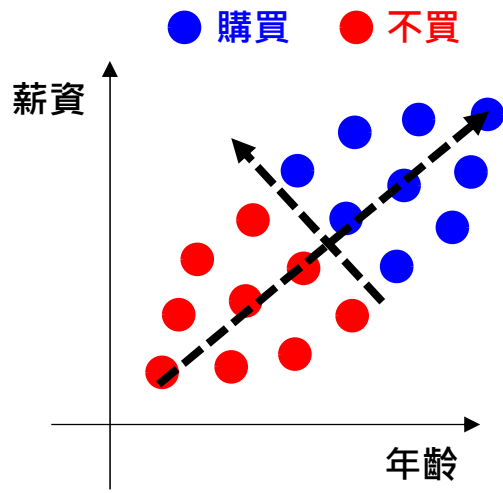
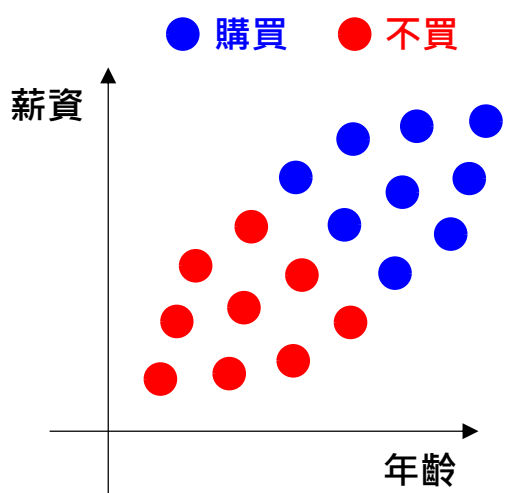




AI

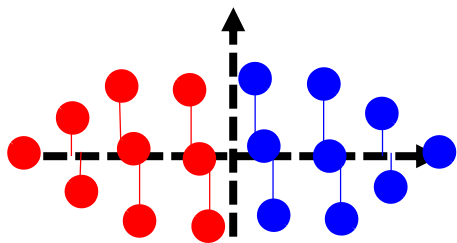
PCA 降維

何謂「PCA 降維」？



• PCA 降維法

- 以「**最大方差法**」找最適座標
(因此一定要先「**特徵縮放**」)
- 所有資料點往**新 X 座標**投影
(此乃線性代數之「**座標變換**」)
- 降維成功



往新 X 座標投影



僅剩一個維度，資訊量略損

使用「標準函式庫」實作



• 程式碼解說

1

```
1 # Feature Scaling
2 X = pp.feature_scaling(fit_ary=X, transform_arrys=X)
3
4 # PCA (with Python Libraries)
5 from sklearn.decomposition import PCA
6 import numpy as np
7 import matplotlib.pyplot as plt
8
```

2

```
9 pca = PCA(n_components=None)
```

3

```
10 pca.fit(X)
```

4

```
11 info_covered = pca.explained_variance_ratio_
```

5

```
12 cumulated_sum = np.cumsum(info_covered)
```

```
13 plt.plot(cumulated_sum, color="blue")
```

1. 「特徵縮放」

(PCA 靠最大方差法，一定得做！)

2. **n_components=None** (需要幾個特徵) 強制列出每一特徵所含資訊量

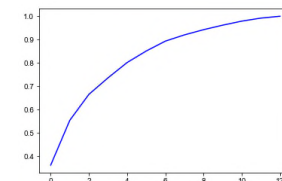
3. **info_covered** (NDArray) 列出每一特徵所含訊息量

```
array([0.36198848, 0.1920749 , 0.11123631, 0.0706903 , 0.06563294,
       0.04935823, 0.04238679, 0.02680749, 0.02222153, 0.01930019,
       0.01736836, 0.01298233, 0.00795215])
```

4. **cumulated_sum** (NDArray) 累進總和

```
array([0.36198848, 0.55406338, 0.66529969, 0.73598999, 0.80162293,
       0.85098116, 0.89336795, 0.92017544, 0.94239698, 0.96169717,
       0.97906553, 0.99204785, 1.          ])
```

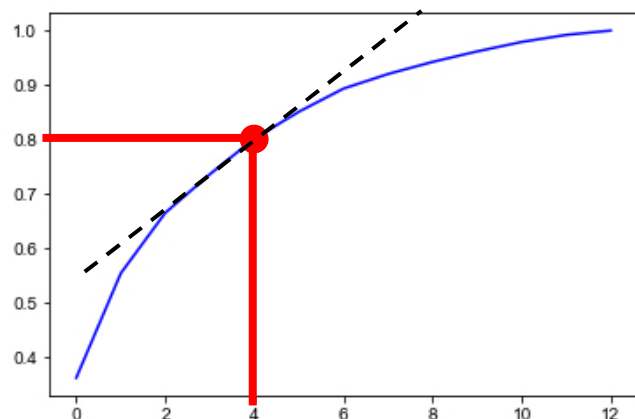
5. 累進總和作圖



使用「標準函式庫」實作



- 如何選取適合的 $n_components$ = ?



包含 1 2 3 4 5
array([[0.36198848, 0.55406338, 0.66529969, 0.73598999, 0.80162293,
0.85098116, 0.89336795, 0.92017544, 0.94239698, 0.96169717,
0.97906553, 0.99204785, 1.]])

個元素的「資訊量」

- 「凸邊形優化法」(Convex Optimization)
→ 現況：斜率最接近 1 者，有最佳解。
- 「資訊量選擇法」
→ 至少涵蓋 80% 資訊量的「特徵數」
- 「直接指定法」
→ 我想要「2 個特徵」

```
15 import pandas as pd
16
17 pca = PCA(n_components=2)
18 X_columns = ["PCA_{}".format(i+1) for i in range(2)]
19 X = pd.DataFrame(pca.fit_transform(X),
20                  index=X.index, columns=X_columns)
```

為了與 KBestSelector 比較，先用 2 個特徵

使用「標準函式庫」實作



- 程式從頭執行一遍的結果

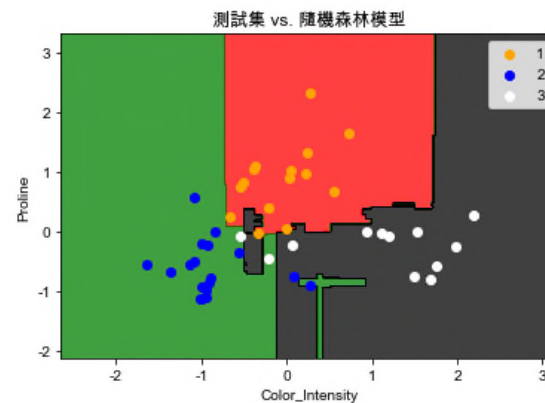
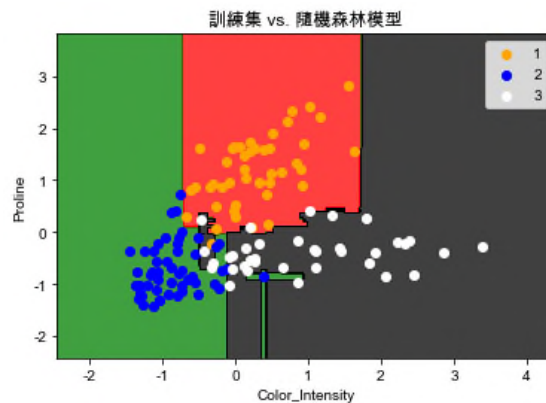
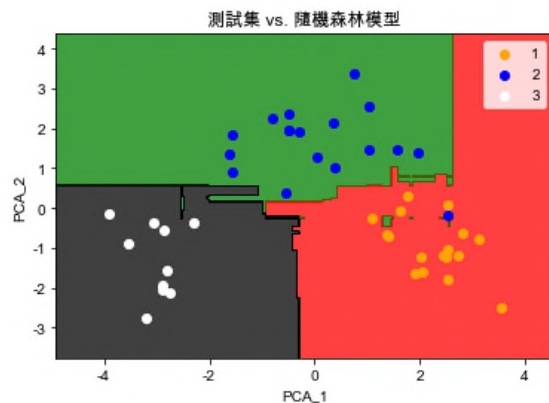
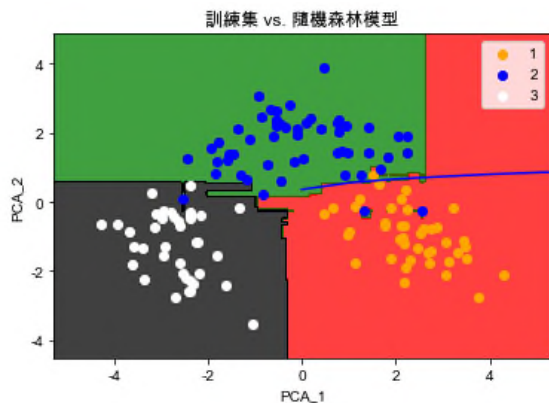
PCA Components = 2

```
----- Random Forest Classification -----  
10 Folds Mean Accuracy: 0.9329377364981081  
10 Folds Mean Recall: 0.9330952380952381  
10 Folds Mean Precision: 0.9451190476190476  
10 Folds Mean F1-Score: 0.9334974284974284
```

VS.

KBestSelector K=2

```
----- Random Forest Classification -----  
10 Folds Mean Accuracy: 0.8783281733746129  
10 Folds Mean Recall: 0.8668253968253969  
10 Folds Mean Precision: 0.890489417989418  
10 Folds Mean F1-Score: 0.8635534889946654
```



• 請撰寫下列程式碼，用 Python 實作 PCA 降維：

| | | |
|---------------|---|---|
| 特徵縮放 | { | 1 # Feature Scaling |
| | | 2 X = pp.feature_scaling(fit_ary=X, transform_arys=X) |
| | | 3 |
| 引入必要套件 | { | 4 # PCA (with Python Libraries) |
| | | 5 from sklearn.decomposition import PCA |
| | | 6 import numpy as np |
| | | 7 import matplotlib.pyplot as plt |
| | | 8 |
| 計算特徵值含金量 | { | 9 pca = PCA(n_components=None) |
| | | 10 pca.fit(X) |
| | | 11 info_covered = pca.explained_variance_ratio_ |
| 累進總和與製圖 | { | 12 cumulated_sum = np.cumsum(info_covered) |
| | | 13 plt.plot(cumulated_sum, color="blue") |
| | | 14 |
| 決定降剩兩維 | { | 15 import pandas as pd |
| | | 16 pca = PCA(n_components=2) |
| 實際降維 & 自製欄位名稱 | { | 17 X_columns = ["PCA_{}".format(i+1) for i in range(2)] |
| | | 18 X = pd.DataFrame(pca.fit_transform(X), |
| | | 19 index=X.index, columns=X_columns) |



• 原始碼講解 (1) :

引入必要套件

類別的成員變數

建構函數

/HappyML/preprocessor.py

```
1 from sklearn.decomposition import PCA
2 from sklearn.preprocessing import MinMaxScaler
3 import matplotlib.pyplot as plt
4
5 class PCASelector:
6     __selector = None
7     __best_k = None
8     __info_covered = None
9     __strategy = None
10
11     def __init__(self, best_k="auto"):
12         if type(best_k) is int:
13             self.__strategy = "fixed"
14             self.__best_k = best_k
15         elif type(best_k) is float:
16             self.__strategy = "percentage"
17             self.__info_covered = best_k
18             self.__best_k = None
19         else:
20             self.__strategy = "auto"
21             self.__best_k = None
22
23         self.__selector = PCA(n_components=self.__best_k)
```

- **best_k = "auto"**
由程式幫您選擇最佳的 K 值。
- **best_k = 0.80** (浮點數)
選擇至少涵蓋 80% 的 K 值。
- **best_k = 2** (整數)
直接選擇 2 個特徵。

• 原始碼講解 (2) :

/HappyML/preprocessor.py

selector 的
getter

best_k 的
getter

.fit()

25 {
26 def selector(self):
27 return self.__selector
28
29 {
30 def best_k(self):
31 return self.__best_k
32 }
33
34 {
35 def fit(self, x_ary, verbose=False, plot=False):
36 pca = PCA(n_components=None)
37 pca.fit(x_ary)
38 info_covered = pca.explained_variance_ratio_
39 cumulated_sum = np.cumsum(info_covered)
40 info_covered_dict = dict(zip([i+1 for i in range(info_covered.shape[0])], cumulated_sum))
41
42 if self.__strategy == "auto":
43 {
44 scaler = MinMaxScaler(feature_range=(0, info_covered.shape[0]-1))
45 scaled_info_covered = scaler.fit_transform(info_covered.reshape(-1, 1)).ravel()
46 for i in range(1, scaled_info_covered.shape[0]):
47 if (scaled_info_covered[i-1]-scaled_info_covered[i]) < 1:
48 break
49
50 找到最佳的 K 值 → self.__best_k = i

是否逐步顯示
↓

是否繪圖
↓

先計算每個特徵
所包含的資訊量

把 Y 軸份數，調整
成與 X 軸相同
當上升比例不足1，停止！

如果是「自動」

• 原始碼講解 (3) :

若是「涵蓋資訊」

執行 PCA 降維

顯示降維結果

繪製圖形

合成指定個數的特徵

```
47 elif self.__strategy == "percentage":
48     current_best = 1
49     cumulated_info = 0.0
50     for i in info_covered:
51         cumulated_info += i
52         if cumulated_info < self.__info_covered:
53             current_best += 1
54         else:
55             break
56     self.__best_k = current_best
57
58     { self.__selector = PCA(n_components=self.best_k)
59       self.selector.fit(x_ary)
60
61     { if verbose:
62       { print("Select {} components, covered information {:.2%}"
63             .format(self.best_k, info_covered_dict[self.best_k]))
64         Select 4 components, covered information 73.60%
65
66     { if plot:
67       { np.insert(cumulated_sum, 0, 0.0)
68         plt.plot(cumulated_sum, color="blue")
69         plt.scatter(x=self.best_k-1, y=cumulated_sum[self.best_k-1], color="red")
70         plt.title("Components vs. Information")
71         plt.xlabel("# of Components (0-based)")
72         plt.ylabel("Covered Information (%)")
73         plt.show()
74
75     return self
76
77     { def transform(self, x_ary):
78       { X_columns = ["PCA_{}".format(i) for i in range(1, self.best_k+1)]
79         return pd.DataFrame(self.selector.transform(x_ary), index=x_ary.index, columns=X_columns)
```

一直累加涵蓋資訊量，直到超過指定量為止



使用「快樂版函式庫」實作

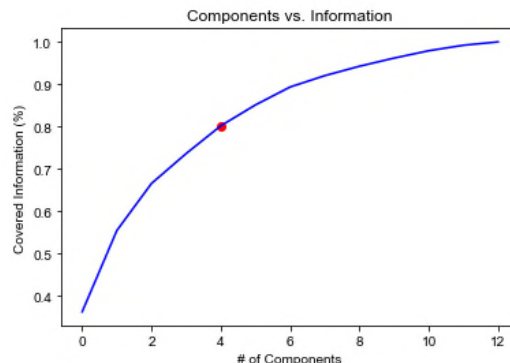


• 呼叫範例：

```
1 # Feature Scaling 不要漏掉「特徵縮放」
2 X = pp.feature_scaling(fit_ary=X, transform_arys=X)
3
4 # PCA with HappyML's Class
5 from HappyML.preprocessor import PCASelector
6
7 selector = PCASelector(best_k="auto") 開 "auto" 讓它自行選擇
8 X = selector.fit(x_ary=X, verbose=True, plot=True).transform(X)
```

• 執行結果：

Select 4 components, covered information 73.60%



----- Random Forest Classification -----

10 Folds Mean Accuracy: 0.9499656002751978

10 Folds Mean Recall: 0.9528571428571428

10 Folds Mean Precision: 0.9585978835978837

10 Folds Mean F1-Score: 0.952495130142189



隨堂練習：使用「快樂版函式庫」實作



- 請先註解掉「用 Python 製作 PCA」與「繪製決定邊界」程式碼：
- 撰寫下列程式碼，並重跑整個程式：

```
1 # # PCA without HappyML's Class
2 # from sklearn.decomposition import PCA
3 # import numpy as np
4 # import matplotlib.pyplot as plt
5 # import pandas as pd
6
7 # pca = PCA(n_components=None)
8 # pca.fit(X)
9 # info_covered = pca.explained_variance_ratio_
10 # cumulated_sum = np.cumsum(info_covered)
11 # plt.plot(cumulated_sum, color="blue")
12
13 # pca = PCA(n_components=2)
14 # X_columns = ["PCA_{}".format(i+1) for i in range(2)]
15 # X = pd.DataFrame(pca.fit_transform(X), index=X.index, columns=X_columns)
16
17 # import HappyML.model_drawer as md
18 #
19 # md.classify_result(x=X_train, y=Y_train, classifier=classifier.classifier,
20 #                   fg_color=("orange", "blue", "white"), bg_color=("red", "green", "black"),
21 #                   title="訓練集 vs. 隨機森林模型", font="DFKai-sb")
22 #
23 # md.classify_result(x=X_test, y=Y_test, classifier=classifier.classifier,
24 #                   fg_color=("orange", "blue", "white"), bg_color=("red", "green", "black"),
25 #                   title="測試集 vs. 隨機森林模型", font="DFKai-sb")
```

```
1 # Feature Scaling
2 X = pp.feature_scaling(fit_ary=X, transform_arrys=X)
3
4 # PCA with HappyML's Class
5 from HappyML.preprocessor import PCASelector
6
7 selector = PCASelector(best_k="auto")
8 X = selector.fit(x_ary=X, verbose=True, plot=True).transform(X)
```



課後作業：動物園分類



• 資料集說明

- 請下載資料集 [Zoo_Data.csv](#)、[Zoo_Class_Name.csv](#)。
- [Zoo_Data.csv](#) 包含動物如何被分類的資料
 - 自變數：hair (是否有毛髮)、feathers (是否有羽毛)、eggs (是否卵生) ...
 - 應變數：class_type (類別) = 1, 2, 3, 4, 5, 6, 7。
- [Zoo_Class_Name.csv](#) 包含「類別」vs.「類別名稱」的對應
 - 1=Mammal、2=Bird...

• 要求

- 請用「[隨機森林](#)」製作一款分類器，能根據生物特徵分類動物園裡的動物。
- 請自行選用「[PCA](#)」或「[KBest](#)」作為您的特徵選擇演算法。留下表現好的那一個。
- 請印出該分類器的「[確度](#)、[廣度](#)、[精度](#)、[F 值](#)」，來說明您的模型有多好。
- 請繪製出隨機森林其中一棵樹（使用 `HappyML.model_drawer.tree_drawer()`）。注意！決策樹「樹葉」要能夠清楚印出「[class=Mammal](#)」（哺乳類），而不是僅印出「[class=1](#)」。





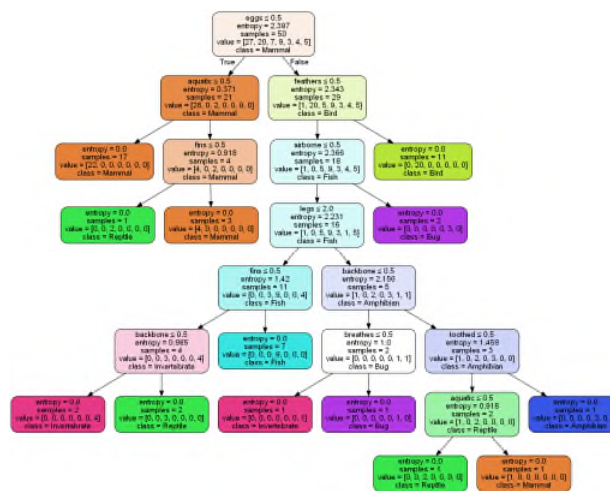
課後作業：動物園分類

- 提示

- 要能夠在決策樹的「樹葉」，印出「`class=Mammal`」，重點在於 `tree_drawer()` 的 `target_names=` 是否傳入正確的參數。
- 如果希望 1=Mammal, 2=Bird, 3=Reptile...，則 `target_names=["Mammal", "Bird", "Reptile", ...]`。
- 類別名稱可以用下列手法讀入：
 - `dataset_classname = pp.dataset("Zoo_Class_Name.csv")`
 - `class_names = [row['Class_Type'] for index, row in dataset_classname.iterrows()]`

● 輸出

```
----- Random Forest Classification -----
10 Folds Mean Accuracy: 0.9614646464646464
10 Folds Mean Recall: 0.9095238095238095
10 Folds Mean Precision: 0.9
10 Folds Mean F1-Score: 0.9038095238095238
```





本章總結



- 隨機森林的原理
 - 建立多棵「**決策樹**」，**投票**決定答案。
- 隨機森林的優點
 - 能處理**大量、無相關性**的自變數
 - 對「**缺失資料**」(Missing Data) 耐受度良好
 - 能產生較為「**不偏 (Un-bias)**」的訓練集
 - 對於「**離群值**」(Outliers) 抗性較佳
- 相關套件
 - 隨機森林：sklearn.ensemble.**RandomForestClassifier**
 - PCA：sklearn.decomposition.**PCA**

