



# 機器學習

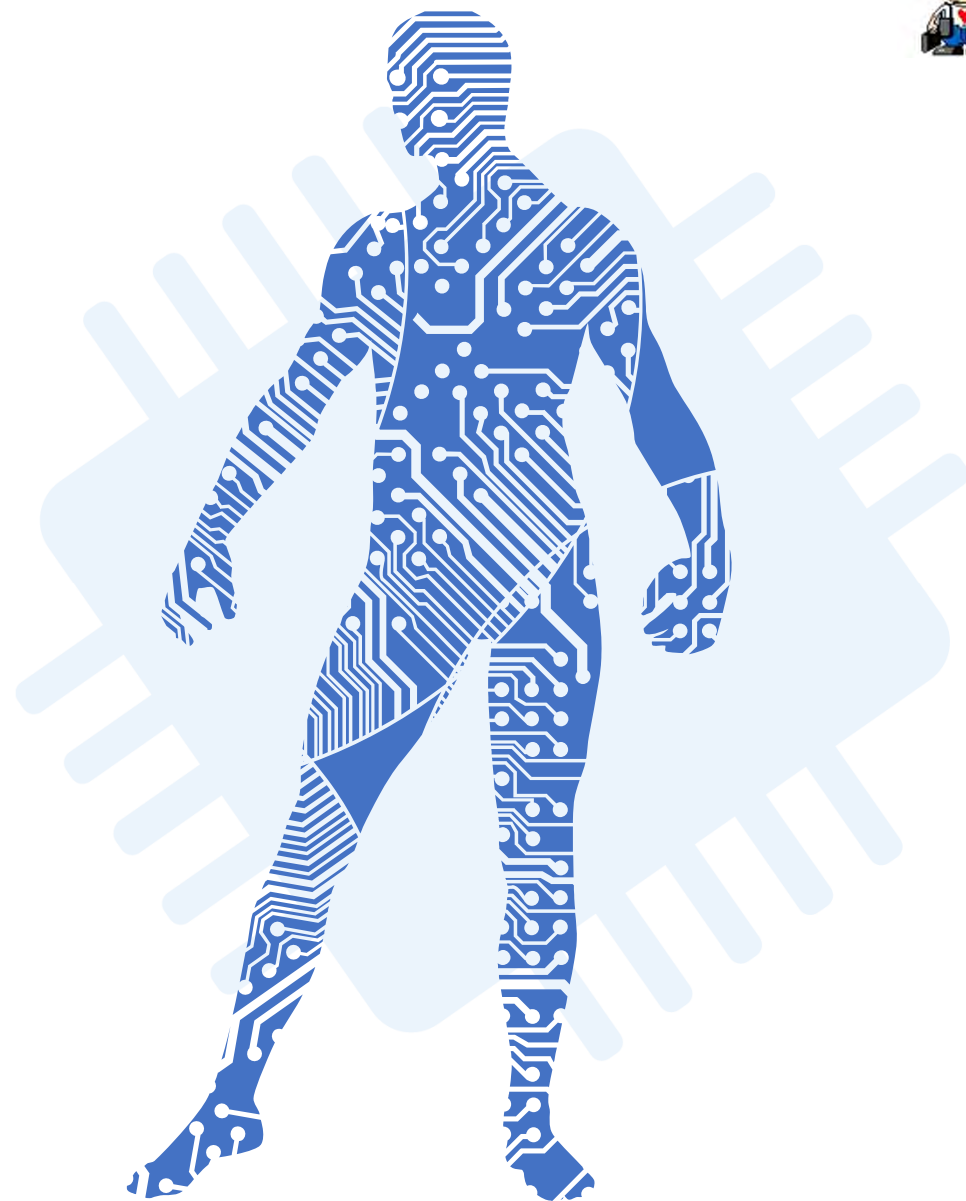
## 第 10 章 單純貝氏分類器 ( Naïve Bayes Classifier )

講師：紀俊男



# 本章大綱

- 理論說明
- 資料前處理
- 實作單純貝氏分類器
- 交叉驗證
- 將模型視覺化
- 檢查「單純」前提是否成立
- 重點整理





# AI

## 理論説明



# 為何叫做「單純 ( Naïve ) 」



- 單純貝氏分類器 **假設前提**

各個「自變數  $x_i$ 」獨立！

**獨立 = 單純** vs. **不獨立 = 不單純**



# 使用到的數學定理



- 貝氏定理 ( Bayes' Theorem )

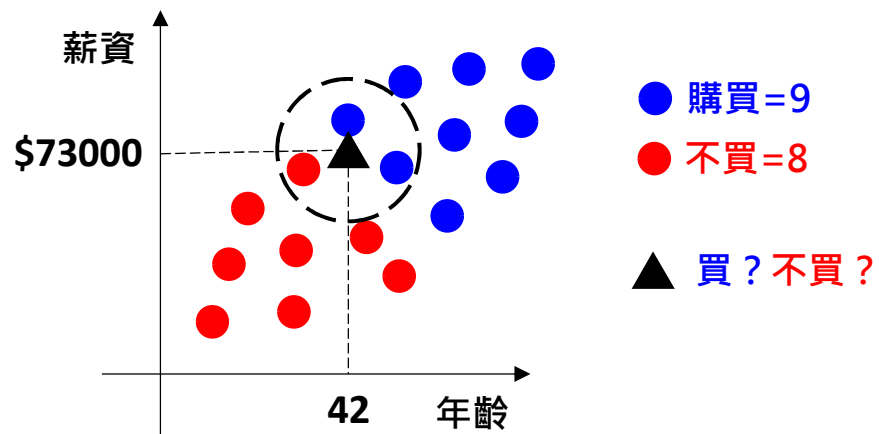
$$P(\textcolor{red}{Y}|\textcolor{blue}{X})P(\textcolor{blue}{X}) = P(\textcolor{blue}{X}|\textcolor{red}{Y})P(\textcolor{red}{Y})$$

$$P(\textcolor{red}{Y}|\textcolor{blue}{X}) = \frac{P(\textcolor{blue}{X}|\textcolor{red}{Y})P(\textcolor{red}{Y})}{P(\textcolor{blue}{X})}$$

# 「貝氏定理」解說



$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$



$$P(X = (42, \$73000)) = \text{○}$$

$$= \text{與 } \blacktriangle \text{ 相似的機率} = \frac{3}{8+9} = \frac{3}{17}$$

$$P(Y = ? | X = (42, \$73000))$$

$$= \begin{cases} P(Y = \text{會買} | X = (42, \$73000)) \\ P(Y = \text{不買} | X = (42, \$73000)) \end{cases}$$

$$= \begin{cases} \frac{P(X = (42, \$73000) | Y = \text{會買})P(Y = \text{會買})}{P(X = (42, \$73000))} \\ \frac{P(X = (42, \$73000) | Y = \text{不買})P(Y = \text{不買})}{P(X = (42, \$73000))} \end{cases}$$

$$= \begin{cases} \frac{2/9 \times 9/17}{3/17} = \frac{2}{3} \text{ (會買)} \\ \frac{1/8 \times 8/17}{3/17} = \frac{1}{3} \text{ (不買)} \end{cases}$$

您覺得是  
買? 不買?



# AI

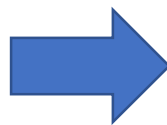
## 資料前處理



- 依照講師指示，**下載**並**瀏覽**資料集



Social\_Network\_Ads.csv



	A	B	C	D	E
1	User ID	Gender	Age	Salary	Purchased
2	15624510	Male	19	19000	0
3	15810944	Male	35	20000	0
4	15668575	Female	26	43000	0
5	15603246	Female	27	57000	0
6	15804002	Male	19	76000	0
7	15728773	Male	27	58000	0
8	15598044	Female	27	84000	0
9	15694829	Female	32	150000	1
10	15600575	Male	25	33000	0

目的：利用「性別、年齡、薪資」  
--> 推算「是否購買」



## 撰寫程式碼

```
1 import HappyML.preprocessor as pp
2
3 # Load Data
4 dataset = pp.dataset(file="Social_Network_Ads.csv")
5
6 # X, Y decomposition
7 X, Y = pp.decomposition(dataset, x_columns=[1, 2, 3], y_columns=[4])
8
9 # One-Hot Encoder
10 X = pp.onehot_encoder(ary=X, columns=[0], remove_trap=True)
11
12 # Feature Selection
13 from HappyML.preprocessor import KBestSelector
14 selector = KBestSelector()
15 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
16
17 # Split Training / TEsting Set
18 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
19
20 # Feature Scaling
21 X_train, X_test = pp.feature_scaling(fit_ary=X_train, transform_arys=(X_train, X_test))
```

## 資料前處理流程：

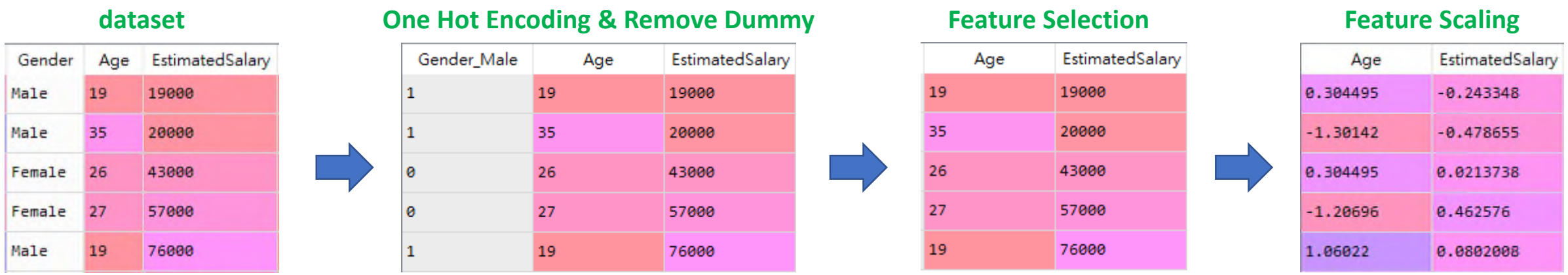
1. 載入資料
2. 切分自變數、應變數
3. 處理缺失資料  
( 無缺失資料 )
4. 類別資料數位化  
( + 移除虛擬變數陷阱 )
5. 特徵選擇
6. 切分訓練集、測試集
7. 特徵縮放  
( 僅縮放 X 即可 )

• 執行結果

```
The Significant Level: 0.05

--- The p-values of Feature Importance ---
TRUE <0.05 0.00000000e+00 (EstimatedSalary)
TRUE <0.05 4.04303193e-100 (Age)
FALSE >0.05 5.44126248e-01 (Gender_Male)

Number of Features Selected: 2
```





- 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
4 dataset = pp.dataset(file="Social_Network_Ads.csv")
5
6 # X, Y decomposition
7 X, Y = pp.decomposition(dataset, x_columns=[1, 2, 3], y_columns=[4])
8
9 # One-Hot Encoder
10 X = pp.onehot_encoder(ary=X, columns=[0], remove_trap=True)
11
12 # Feature Selection
13 from HappyML.preprocessor import KBestSelector
14 selector = KBestSelector()
15 X = selector.fit(x_ary=X, y_ary=Y, verbose=True, sort=True).transform(x_ary=X)
16
17 # Split Training / TEsting Set
18 X_train, X_test, Y_train, Y_test = pp.split_train_test(x_ary=X, y_ary=Y)
19
20 # Feature Scaling
21 X_train, X_test = pp.feature_scaling(fit_ary=X_train, transform_arys=(X_train, X_test))
```





# AI

## 實作單純貝氏分類器

• 程式碼

1

```
1 from sklearn.naive_bayes import GaussianNB
2
3 model = GaussianNB()
4 model.fit(X_train, Y_train.values.ravel())
5
6 Y_pred = model.predict(X_test)
```

1. 引入「高斯分佈」的簡單貝氏分類器
  - **GaussianNB**: 自變數  $X_i$  呈「高斯分佈（常態）」時使用
  - **MultinomialNB**: 自變數  $X_i$  呈「多項分佈（多選一）」時使用
  - **BernoulliNB**: 自變數  $X_i$  呈「白努利分佈（二選一）」時使用
2. 訓練模型
3. 預測結果

• 執行結果

Y_test			Y_pred	
	0		0	
0	0		0	
1	1		1	
2	1		1	
3	0		0	
4	0		0	
5	1		1	
6	0		0	
7	0		0	
8	1		1	
9	1		0	
10	0		0	



- 請撰寫下列程式碼，並執行之：

```
1 from sklearn.naive_bayes import GaussianNB
2
3 model = GaussianNB()
4 model.fit(X_train, Y_train.values.ravel())
5
6 Y_pred = model.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。







- 原始程式碼解說 ( 1 ) :

/HappyML/classification.py

引入必要的套件 →

類別成員變數 →

建構函數，根據參數，  
初始化不同分佈的單純貝氏核心

```
1  from sklearn.naive_bayes import BernoulliNB, MultinomialNB, GaussianNB
2  import pandas as pd
3
4  class NaiveBayesClassifier:
5      __classifier = None
6      __y_columns = None
7
8      def __init__(self, type="gaussian"):
9          algorithm_dict = {
10              "bernoulli" : BernoulliNB(),
11              "multinomial" : MultinomialNB(),
12              "gaussian" : GaussianNB()
13          }
14          self.__classifier = algorithm_dict[type]
```

- 原始程式碼解說（2）：

/HappyML/classification.py

16	@property		
17	def classifier(self):		
18	return self.__classifier		
19		classifier 的 getter/setter	
20	@classifier.setter		
21	def classifier(self, classifier):		
22	self.__classifier = classifier		
23			
24	def fit(self, x_train, y_train):		
25	self.classifier.fit(x_train, y_train.values.ravel())		
26	self.__y_columns = y_train.columns	訓練模型	
27			
28	return self		
29			
30	def predict(self, x_test):		
31	return pd.DataFrame(self.classifier.predict(x_test), index=x_test.index, columns=self.__y_columns)	預測結果	



- 呼叫範例

```
1 from HappyML.classification import NaiveBayesClassifier
2
3 model = NaiveBayesClassifier()
4 Y_pred = model.fit(X_train, Y_train).predict(X_test)
```

- 執行結果

Y_test		Y_pred	
Index	Purchased	Index	Purchased
20	1	20	1
163	0	163	0
348	0	348	0
213	0	213	0
373	1	373	1



- 請先將前一小節「**實作單純貝氏分類器**」的程式碼**註解**掉。
- 請用「**變數觀察面板**」，**清除**掉所有變數。
- 輸入下列程式碼，並重跑整個程式：

```
1 from HappyML.classification import NaiveBayesClassifier
2
3 model = NaiveBayesClassifier()
4 Y_pred = model.fit(X_train, Y_train).predict(X_test)
```

- 比較 **Y\_test (真實值)** 與 **Y\_pred (預測值)** 的結果。





# AI

## 交叉驗證

- 何謂「K 次交叉驗證 ( K-Fold Cross Validation ) 」

5-fold CV

DATASET

Estimation 1	Test	Train	Train	Train	Train
Estimation 2	Train	Test	Train	Train	Train
Estimation 3	Train	Train	Test	Train	Train
Estimation 4	Train	Train	Train	Test	Train
Estimation 5	Train	Train	Train	Train	Test

accuracy 1

accuracy 2

accuracy 3

accuracy 4

accuracy 5

Total Accuracy

也可以求算  
Recall, Precision...等

➡

不受單次抽樣偏頗影響！

• 程式碼解說

引入必要套件

```
1 from sklearn.model_selection import cross_val_score
7 accuracies = cross_val_score(estimator=model.classifier, X=X, y=Y.values.ravel(), scoring="accuracy", cv=10, n_jobs=-1)
4 print("{} Folds Mean Accuracy: {}".format(k_fold, accuracies.mean()))
```

- 1. **esitmiator=**：輸入機器學習模型
- 2. **X=**：輸入自變數
- 3. **y=**：輸入應變數
- 4. **scoring=**：輸入此次求算的效能指標。  
( 可以輸入 "accuracy" , "recall" , "precision" , "f1" 等數值 )
- 5. **cv=**：Cross Validation 的次數。如：cv=10 → 10-Fold。
- 6. **n\_jobs=**：使用多少 CPU 核心 ( 如：1, 2, ... )。-1 代表全部使用。
- 7. **accuracies**：保存 10 份 Accuracy 的 NDAarray 陣列。
- 8. **.mean()**：取平均 ( Average )。

# 隨堂練習：以「標準函式庫」實作交叉驗證



- 請依照前一頁投影片的說明，以「交叉驗證」評估模型效能：

```
1 from sklearn.model_selection import cross_val_score
2
3 accuracies = cross_val_score(estimator=model.classifier, X=X, y=Y.values.ravel(), scoring="accuracy", cv=10, n_jobs=-1)
4 print("{} Folds Mean Accuracy: {}".format(k_fold, accuracies.mean()))
```

- 如果要計算「廣度」、「精度」、或「F1-score」，可以參考下列程式碼：

```
1 from sklearn.model_selection import cross_val_score
2
3 k_fold = 10
4 accuracies = cross_val_score(estimator=model.classifier, X=X, y=Y.values.ravel(), scoring="accuracy", cv=k_fold, n_jobs=-1)
5 print("{} Folds Mean Accuracy: {}".format(k_fold, accuracies.mean()))
6
7 recalls = cross_val_score(estimator=model.classifier, X=X, y=Y.values.ravel(), scoring="recall", cv=k_fold, n_jobs=-1)
8 print("{} Folds Mean Recall: {}".format(k_fold, recalls.mean()))
9
10 precisions = cross_val_score(estimator=model.classifier, X=X, y=Y.values.ravel(), scoring="precision", cv=k_fold, n_jobs=-1)
11 print("{} Folds Mean Precision: {}".format(k_fold, precisions.mean()))
12
13 f_scores = cross_val_score(estimator=model.classifier, X=X, y=Y.values.ravel(), scoring="f1", cv=k_fold, n_jobs=-1)
14 print("{} Folds Mean F1-Score: {}".format(k_fold, f_scores.mean()))
```



• 程式碼解說 ( 1 )

類別的成員變數

建構函數

k\_fold 成員變數  
的 getter

/HappyML/performance.py

引入所需的套件

1 from sklearn.model\_selection import cross\_val\_score

2

3 class KFoldClassificationPerformance:

4 {

5 {

6 {

7 {

8 {

9 {

10 def \_\_init\_\_(self, x\_ary, y\_ary, classifier, k\_fold=10, verbose=False):

11 {

12 {

13 {

14 {

15 {

16 {

17 {

18 {

19 {

• 程式碼解說 ( 2 )

k\_fold 成員變數  
的 setter

verbose 成員變數  
的 getter & setter

classifier 成員變數  
的 getter

/HappyML/performance.py

```
21 @k_fold.setter
22 def k_fold(self, k_fold):
23     if k_fold >=2:
24         self.__k_fold = k_fold
25     else:
26         self.__k_fold = 2
27
28 @property
29 def verbose(self):
30     return self.__verbose
31
32 @verbose.setter
33 def verbose(self, verbose):
34     if verbose:
35         self.__verbose = 10
36     else:
37         self.__verbose = 0
38
39 @property
40 def classifier(self):
41     return self.__classifier
```



# 以「快樂版函式庫」實作交叉驗證



## • 程式碼解說 ( 3 )

/HappyML/performance.py

確度

廣度

精度

F 值

```
43 def accuracy(self):
44     results = cross_val_score(estimator=self.classifier, X=self.__x_ary, y=self.__y_ary.values.ravel(), scoring="accuracy", cv=self.k_fold, verbose=self.verbose)
45     return results.mean()
46
47 def recall(self):
48     def recall_scorer(estimator, X, y):
49         return recall_score(y, estimator.predict(X), average="macro")
50
51     results = cross_val_score(estimator=self.classifier, X=self.__x_ary, y=self.__y_ary.values.ravel(), scoring=recall_scorer, cv=self.k_fold, verbose=self.verbose)
52     return results.mean()
53
54 def precision(self):
55     def precision_scorer(estimator, X, y):
56         return precision_score(y, estimator.predict(X), average="macro")
57
58     results = cross_val_score(estimator=self.classifier, X=self.__x_ary, y=self.__y_ary.values.ravel(), scoring=precision_scorer, cv=self.k_fold, verbose=self.verbose)
59     return results.mean()
60
61 def f_score(self):
62     def f1_scorer(estimator, X, y):
63         return fbeta_score(y, estimator.predict(X), beta=1, average="macro")
64
65     results = cross_val_score(estimator=self.classifier, X=self.__x_ary, y=self.__y_ary.values.ravel(), scoring=f1_scorer, cv=self.k_fold, verbose=self.verbose)
66     return results.mean()
```

• 呼叫範例

引入所需的套件

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

自變數、應變數

機器學習模型

K 值

是否逐步  
顯示過程

印出各種  
指標

• 執行結果

```
10 Folds Mean Accuracy: 0.8775719199499689
10 Folds Mean Recall: 0.8578717948717948
10 Folds Mean Precision: 0.8874708964782494
10 Folds Mean F1-Score: 0.8619361410762924
```



- 請輸入下列程式碼，並觀看您模型「K 次交叉驗證」效能：

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





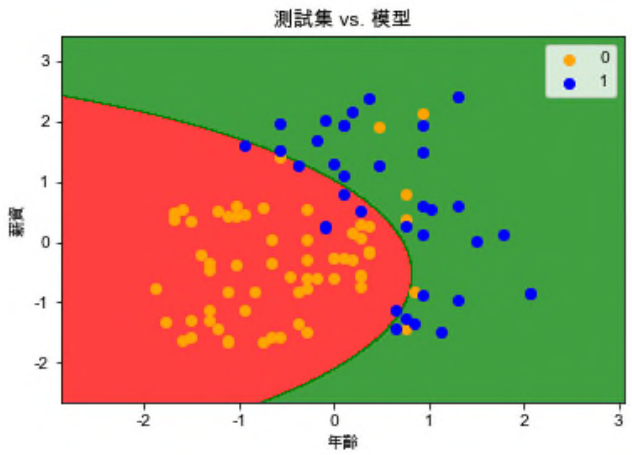
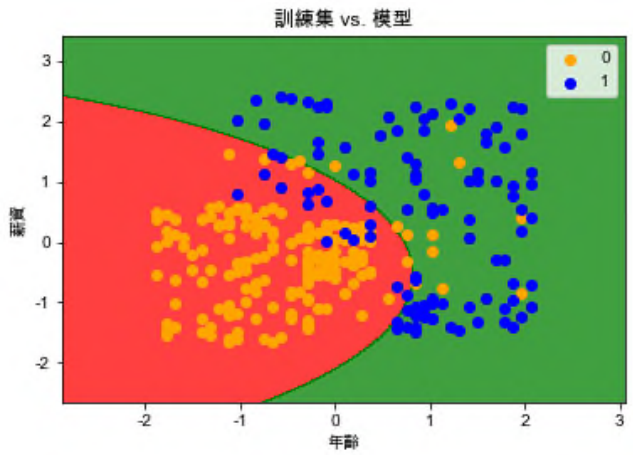
# AI

## 將模型視覺化

- 程式碼

```
1 import HappyML.model_drawer as md
2
3 md.classify_result(x=X_train, y=Y_train, classifier=model.classifier, title="訓練集 vs. 模型", font='DFKai-sb')
4 md.classify_result(x=X_test, y=Y_test, classifier=model.classifier, title="測試集 vs. 模型", font='DFKai-sb')
```

- 執行結果



注意！自變數 X 必須是兩個特徵！且經過縮放！



- 請輸入下列程式碼，並執行觀看您**模型**的好壞：

```
1 import HappyML.model_drawer as md
2
3 md.classify_result(x=X_train, y=Y_train, classifier=model.classifier, title="訓練集 vs. 模型", font='DFKai-sb')
4 md.classify_result(x=X_test, y=Y_test, classifier=model.classifier, title="測試集 vs. 模型", font='DFKai-sb')
```





檢查「單純」前提  
是否成立



# 檢查前提是否成立

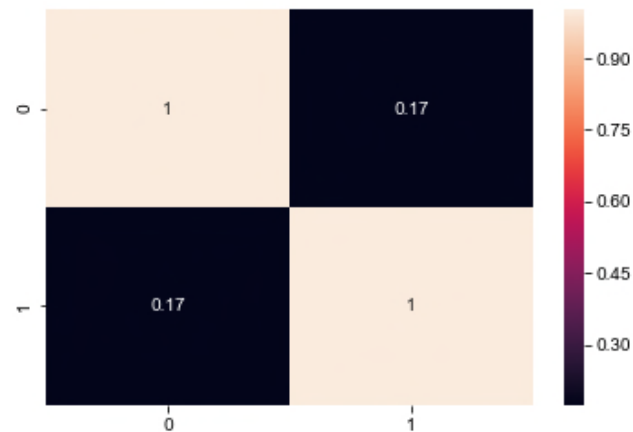


- 簡單貝氏分類器前提：自變數各自獨立

```
1 from HappyML.criteria import AssumptionChecker
2
3 checker = AssumptionChecker(x_train=X_train, x_test=X_test, y_train=Y_train, y_test=Y_test, y_pred=Y_pred)
4 checker.features_correlation(heatmap=True)
```

- 執行結果

```
*** Check for Correlation of Features ***
--- Features Correlation Matrix ---
      0      1
0  1.0000  0.1735
1  0.1735  1.0000
No Correlation (>=0.8) Found!
```







# 隨堂練習：檢查前提是否成立



- 請輸入、並執行下列程式碼，檢查**自變數**是否各自**獨立**？

```
1 from HappyML.criteria import AssumptionChecker
2
3 checker = AssumptionChecker(x_train=X_train, x_test=X_test, y_train=Y_train, y_test=Y_test, y_pred=Y_pred)
4 checker.features_correlation(heatmap=True)
```



- 說明：
  - 下載「美國國家糖尿病、消化、與腎臟病研究中心」的資料庫
    - <https://www.kaggle.com/kandij/diabetes-dataset>
  - 各欄位說明如下：
    - **Pregnancies**：懷孕次數（次）
    - **Glucose**：口服葡萄糖耐量試驗中血漿葡萄糖濃度
    - **BloodPressure**：舒張壓（mmHg）
    - **SkinThickness**：三頭肌組織皮褶厚度（mm）
    - **Insulin**：兩小時血清胰島素（ $\mu\text{U}/\text{ml}$ ）
    - **BMI**：身體質量指數
    - **DiabetesPedigreeFunction**：糖尿病系統功能
    - **Age**：年齡（歲）
    - **Outcome**：是否有糖尿病



# 課後作業：糖尿病資料庫



## 要求

- 使用 **KBestSelector** 挑選「**夠顯著**」的特徵，並**印出**你選中了哪些特徵？  
(提示：使用預設參數 best\_k= "auto" 即可)
- 印出「**K 次交叉驗證**」的「**確度、廣度、精度、F-Score**」，說明你的模型有多好？
- 檢查您的模型是否符合「**單純貝氏**」分類器的**前提**(自變數獨立)？
- 用 **KBestSelector** 將特徵壓至 **2 個**，另外訓練出一個單純貝氏分類器。
- 用這個「僅有兩個特徵」的**資料集與模型**，繪製出「**訓練集**」與「**測試集**」的分類結果。

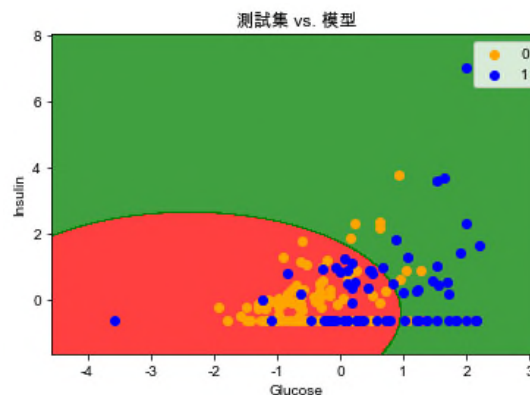
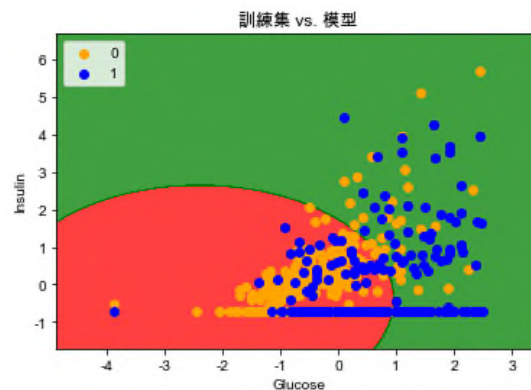
## 輸出

```
The Significant Level: 0.05

--- The p-values of Feature Importance ---
TRUE <0.05 0.00000000e+00 (Insulin)
TRUE <0.05 5.48728628e-309 (Glucose)
TRUE <0.05 2.51638830e-41 (Age)
TRUE <0.05 1.32590849e-29 (BMI)
TRUE <0.05 4.55261043e-26 (Pregnancies)
TRUE <0.05 3.15697650e-13 (SkinThickness)
TRUE <0.05 2.71819252e-05 (BloodPressure)
TRUE <0.05 2.02213728e-02 (DiabetesPedigreeFunction)

Number of Features Selected: 8
10 Folds Mean Accuracy: 0.7564935064935066
10 Folds Mean Recall: 0.7185099715099715
10 Folds Mean Precision: 0.7348495448677373
10 Folds Mean F1-Score: 0.7235726493134441
*** Check for Correlation of Features ***
--- Features Correlation Matrix ---
              Pregnancies ...      Age
Pregnancies      1.0000 ...  0.5530
Glucose          0.1505 ...  0.2718
BloodPressure    0.1524 ...  0.2719
SkinThickness   -0.0670 ... -0.0680
Insulin         -0.0567 ... -0.0379
BMI              0.0688 ...  0.0719
DiabetesPedigreeFunction 0.0031 ...  0.0688
Age              0.5530 ...  1.0000

[8 rows x 8 columns]
No Correlation (>=0.8) Found!
```



## 加分題：

- 這個題目用「**單純貝氏**」好？還是「**邏輯迴歸**」模型好？





- 單純貝氏分類器成立前提

- 自變數各自獨立

- 應該了解的數學定理

- 貝氏定理： $P(Y|X)P(X) = P(X|Y)P(Y)$

- 三種單純貝氏分類器的分佈模型

- 高斯分佈：GaussianNB - 自變數呈常態分佈時使用
- 多項分佈：MultinomialNB - 自變數呈多項分佈（多擇一）時使用
- 白努利分佈：BernoulliNB - 自變數呈白努利分佈（二擇一）時使用

- K 次交叉驗證

- 可以不受抽樣影響、公正客觀評判模型的確度、廣度、精度、F 值。
- 使用 `sklearn.model_selection.cross_val_score()` 實作

- 何時使用「單純貝氏分類器」

- 小樣本（千筆以下）
  - 靠樣本推算母體機率分佈，再拿最可能的母體機率分佈預測未來。
  - 屬「模型生成式（Generative）」推估法，受樣本影響較小。
- 觀察樣本點，有一個「決定邊界」，但該邊界彎彎曲曲，並非任何數學方程式可以擬合

