

dataset 搜集与整理

Balance
InBalance

排演算法建模 (model)

KNN
貝氏演算法
SVM
deep learning

產生混淆矩陣

用測試資料建的

如果是用訓練資料，資料本身就是“已知”的

算 accuracy 沒意義！

Training (75%) (70%)

隨機取

Testing (25%) (30%)

隨機取，做交叉驗證

為避免依賴
特定資料產生
的偏差

什麼是混淆矩陣？

用來判定一個模型表現的好、壞

		Predict		
		Positive	Negative	
Real	Positive	A TP	B FN	recall: $\frac{A}{A+B}$ 召回率
	Negative	C FP	D TN	specificity: $\frac{D}{C+D}$
Precision 精確率 $\frac{A}{A+C}$		negative predictive value $\frac{D}{B+D}$		Accuracy 正確率 $\frac{A+D}{A+B+C+D}$

多少正樣本被正確識別

1. recall: 真陽性比率

2. specificity: 真陰性比率

3. 真陽性 Precision: $\frac{A}{A+C}$

4. 真陰性: $\frac{D}{B+D}$

5. error-rate (錯誤率)

TP (True Positive): 正確預測正樣本

TN (True Negative): 正確預測負樣本 $\Rightarrow \frac{B+D}{A+B+C+D}$

FP (False Positive): 錯誤預測正樣本

FN (False Negative): 錯誤預測負樣本

6. F1-measure (準確率 & 召回率的
調和平均)

Ex: 判斷是否有 COVID-19

TP: 真的陽性且判斷陽性

TN: 真的陰性且判斷陰性

FP: 真的陰性但判斷陽性

FN: 真的陽性但判斷陰性 * 嚴重

* Overfitting: 雖然能完美的符合測試“訓練資料”，

但測試得太過精確，在新的測試資料中會有更高的 error rate

$$= \frac{2a}{2a+b+c}$$