

§ 2. 模型评估与选择

1. 经验误差与过拟合

$$\text{error rate } E = \frac{a}{m}$$

m 个样本中有 a 个样本分类错误

$$\text{accuracy} = 1 - \frac{a}{m}$$

overfitting 过拟合: 学习器在训练样本学得“太好”, 泛化能力下降
underfitting 欠拟合: 训练样本的一般性质尚未学好

2. 评估方法

1) hold-out 留出法

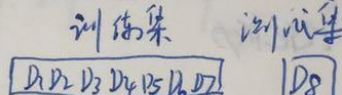
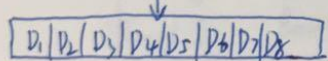
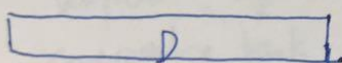
1000 $\begin{cases} S & 700 \text{ 个} \\ T & 300 \text{ 个} \end{cases}$

注意训练/测试集的划分要多次
保持数据分布的一致性

为保证 fidelity, 将 2/3 - 4/5 样本用于训练, 其余测试

2) 交叉验证法 cross validation

Leave-one-out 留一法 (特例)



训练集

测试集

D_8

D_8

平均 \rightarrow 结果

3) bootstrapping 自助法: 每次取一个样本, 将其他只放入 D' . 重复 m 次后, 得到包含 m 个样本的数据集 D'

3. 性能度量 performance measure

1) mean squared error: $E(f; D) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2$

accuracy: $acc(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) = y_i) = 1 - E(f; D)$

2) precision/recall rate

Confusion Matrix

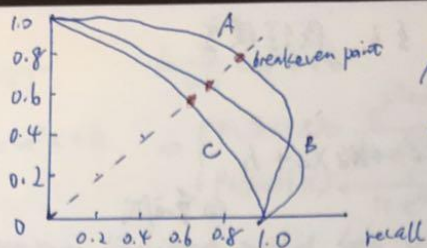
Actual	Prediction	
	+	-
+	TP	FN
-	FP	TN

Precision = $\frac{TP}{TP + FP}$

Recall = $\frac{TP}{TP + FN}$

P-R plot

precision



A is better than C

✓ break-even point

$$F_1 = \frac{2PR}{P+R} = \frac{2 \cdot TP}{\text{sample size} + TP + TN}$$

$$F_2 = \frac{(1+R^2) \cdot P \cdot R}{(R^2 \cdot P) + R}$$

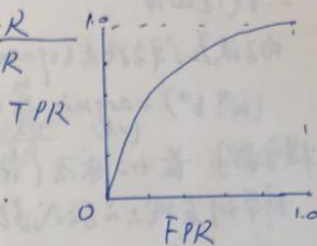
3) ROC, AUC

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{TN + FP}$$

ROC: (0,1)

AUC: area under curve



4) 估计模型错误率与代价函数

$$\text{cost-sensitive error rate} \quad E(f; D; \text{cost}) = \frac{1}{n} \left(\sum_{x_i \in D^+} \mathbb{I}(f(x_i) \neq y_i) \times \text{cost } 0/1 + \sum_{x_i \in D^-} \mathbb{I}(f(x_i) \neq y_i) \times \text{cost } 1/0 \right)$$

4. 比较检验

1) 假设检验

2) 交叉验证与检验

3) McNemar

Contingency table

B	+	-
+	e_{00}	e_{01}
-	e_{10}	e_{11}

$$\chi^2 = \frac{(|e_{01} - e_{10}| - 1)^2}{e_{01} + e_{10}}$$

4) Friedman, Nemenyi: 后验检验

Friedman - 多组比较, 并法排序

Nemenyi: 后验 - 进一步区分各并法

5) 偏差与方差

$$\text{估计误差} = \text{偏差} + \text{方差} + \text{噪声} = \text{bias}^2(x) + \text{var}(x) + \epsilon$$

$$\text{var}(x) = E_D [f(x; D) - \bar{f}(x)]^2$$

$$\epsilon^2 = E_D [(y_D - y)]^2$$

$$\text{bias}^2(x) = (\bar{f}(x) - y)^2$$