

## Evaluating Model (Validation Techniques)

### 1. Resubstitution

所有的数据都被用于训练。

无法评估模型对于没有出现过的数据的效果。

### 2. Hold out validation

将数据集划分为两个互斥的子集，即训练集与测试集。如按照 6:4, 7:3, 8:2 的比例划分：

- 1) 划分可能导致两个子集中类别不平衡或数据分布不一致  $\Rightarrow$  分层采样 (stratified sampling)
- 2) 一次划分可能带来误差  $\Rightarrow$  重复多次划分，取多次结果的平均。
- 3) 每次都会有部分数据没有被用于训练，导致信息丢失。

### 3. k-Fold cross-validation

数据集被划分为  $k$  个相同大小的子集，每次使用  $k-1$  个集作为训练集，剩下的一个作为测试集。这样重复  $k$  次训练，最终的误差可以取  $k$  次训练的平均。

- 1) 为了保证数据一致性，所以对  $k$  个集采用分层采样方法。
- 2)  $k$  的取值通常为 10。
- 3)  $k$  个子集的划分同样有随机性，因此重复划分  $\Rightarrow$   $p$  次  $k$  折交叉验证。

### 4. Leave-One-Out Cross-Validation (LOOCV)

$k$ -Fold cross validation 的特殊情况，此时  $k = N$ ，每次只有一个样本作为测试集。

### 5. Random subsampling

每次从整个数据集中采样  $m$  个数据作为测试集，剩下的数据作为训练集。

## 6. Bootstrapping

前面提及的算法都存在训练集比  $D$  小的情况，会引入由训练规模不同导致的误差。

给定包含  $n$  个样本的数据集  $D$ ，每次从  $D$  中有放回地采样  $n$  个样本组成训练集  $D_1$ 。

重复这个步骤  $m$  次，即可得到  $m$  个包含  $n$  个样本的训练集。每次采样剩下的  $n$  个样本

组成对应的测试集：

Completed dataset:  $\{x_1, x_2, x_3, x_4, x_5\}$

1:  $\{x_1, x_2, x_4, x_4, x_3\}$   $\{x_5\}$

2:  $\{x_4, x_5, x_5, x_1, x_1\}$   $\{x_2, x_3\}$

一个样本在  $m$  次采样中始终不被采样到的概率是： $(1 - \frac{1}{n})^m$ ，有：

$$\lim_{m \rightarrow \infty} (1 - \frac{1}{n})^m \approx \frac{1}{e} \approx 0.368.$$

适用于数据集较小，难以有效划分训练/测试集时很有用。

# Evaluation Metrics

## classification

$$prevalence = \frac{n_{pos}}{n}$$

数据集中正样本的比例

### 1. Accuracy

$$accuracy = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(f(x_i) = y_i)$$

不适用于类别不平衡情况，无法知道每个类别的情况。

### 2. Confusion Matrix, Precision, Recall, F-measure.

(binary)

|                   | Predict positive    | Predict negative    |
|-------------------|---------------------|---------------------|
| Actually positive | True positive (TP)  | False negative (FN) |
| Actually negative | False positive (FP) | True negative (TN)  |

正类的 precision

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

⇒ 预测为正类的样本中有多少确实是正类。

正类的 recall

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

⇒ 所有正类样本有多少被预测正确。

$$sensitivity = recall = \frac{TP}{N_p} \quad \rightarrow \text{正类的准确率, true positive rate}$$

$$specificity = \frac{TN}{N_n} \quad \rightarrow \text{负类的准确率, true negative rate}$$

$$\text{type-I error} = \text{false positive rate} = \frac{FP}{TN + FP} = 1 - specificity \quad \rightarrow \text{把假的说成真的}$$

$$\text{type-II error} = \text{false negative rate} = \frac{FN}{TP + FN} = 1 - sensitivity$$



$$\text{balanced accuracy} = \frac{1}{2} (\text{sensitivity} + \text{specificity})$$

↓  
每个类别正确率的平均。

$$= \frac{1}{2} (\text{recall-pos} + \text{recall-neg})$$

$$= \frac{1}{2} (\text{au-pos} + \text{au-neg})$$

$$F\text{-measure} = \frac{(1 + \beta^2) P \times R}{\beta^2 P + R}$$

F-measure

$$F_\beta = \frac{(1 + \beta^2) P \times R}{\beta^2 P + R}$$

⇒ P和R的调和平均。β > 1时, recall更重要。

α < 1时, precision更重要

$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

⇒ precision和recall同样重要

balanced accuracy 和 F-measure 适用于类别不平衡问题。

1)  $N \gg P$ 时, recall-neg - 定是很低, 因此  $F_1$  更合适。

2)  $p \gg N$ 时, recall-pos 和 pre-pos 都很低, 因此 balanced accuracy 更合适。

### Average Precision (AP)

precision 只考虑了返回结果的准确性却没有考虑结果的顺序。在排序相关应用中(搜索), 顺序是很重要的。

1) 按照样本预测为正类别的概率从大到小排序。

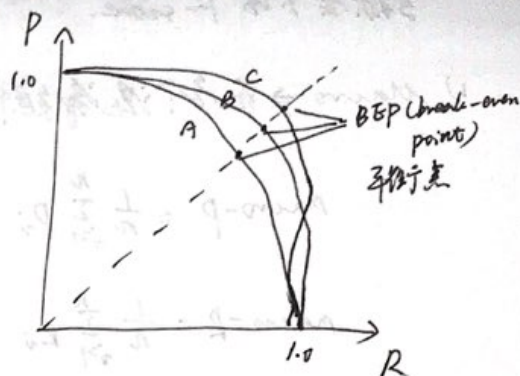
2) 以每个位置作为划分点, 计算 precision。

3) AP = 所有 precision 的平均。

4) mAP (mean average precision) = 多次查询的 AP 值的平均。

## P-R Curve

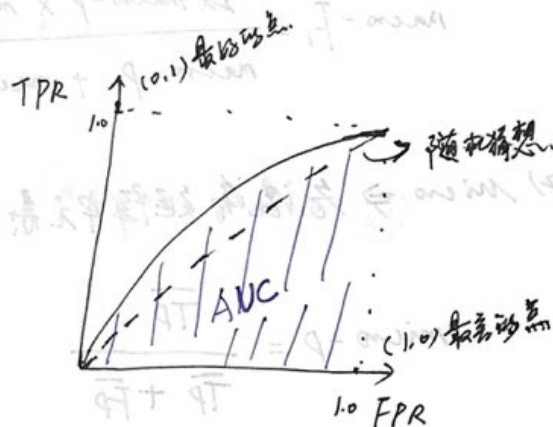
- 1) 将样本预测为正的的概率从大到小排序
- 2) 分别以每个位置作为阈值点, 得到对应的 P, R.
- 3) BEP 越大, 说明性能越好.



## ROC - Curve.

$$\begin{cases} \text{TPR} = \frac{TP}{TP + FN} & \Rightarrow \text{true positive rate} \\ \text{FPR} = \frac{FP}{FP + TN} & \Rightarrow \text{false positive rate} \end{cases}$$

AUC 越大, 性能越好.



多标签下的 F-score.

1) Macro  $\Rightarrow$  在各个混淆矩阵上分别计算 P 和 R, 有:

$$\text{macro-P} = \frac{1}{n} \sum_{i=1}^n P_i$$

$$\text{macro-R} = \frac{1}{n} \sum_{i=1}^n R_i$$

$$\text{macro-F}_1 = \frac{2 \times \text{macro-P} \times \text{macro-R}}{\text{macro-P} + \text{macro-R}}$$

2) Micro  $\Rightarrow$  各混淆矩阵中元素求平均, 只算一次 P, R.

$$\text{micro-P} = \frac{\overline{TP}}{\overline{TP} + \overline{FP}}$$

$$\text{micro-R} = \frac{\overline{TP}}{\overline{TP} + \overline{FN}}$$

$$\text{micro-F}_1 = \frac{2 \times \text{micro-P} \times \text{micro-R}}{\text{micro-P} + \text{micro-R}}$$

log loss: (logistic regression loss / cross-entropy loss)  
(熵?)

binary:  ~~$L = -\log P(y|p)$~~

1° binary:  $L = -[y \log p + (1-y) \log (1-p)]$    
 → 样本属于类别1的概率.

$$L = \frac{1}{N} \sum_{i=0}^N L_i = \frac{1}{N} \sum_{i=0}^N -[y_i \log p_i + (1-y_i) \log (1-p_i)]$$

likelihood =  $p^y \cdot (1-p)^{(1-y)}$    
  $\xrightarrow{-\log} -[y \log p + (1-y) \log (1-p)]$

2° multiclass

$$L = -\frac{1}{N} \sum_{i=0}^N \sum_{k=1}^K y_{i,k} \log p_{i,k}$$

likelihood =  $\prod_{i=0}^N \prod_{k=1}^K p_{i,k}^{y_{i,k}}$    
  $\xrightarrow{-\log} -\sum_{i=0}^N \sum_{k=1}^K y_{i,k} \log p_{i,k}$

$y_i$  是一个 1-of-k binary indicator (one-hot encoding)

Motivation 两个模型如果有相同的预测结果, 那么 accuracy, AU-ROC, AU-PRC 都是一致的.

$$\text{Gain} = y p(x) + (1-y) (1-p(x))$$

~~熵:  $H = -\sum p_i \log p_i$~~

$$\text{log loss} = -[y \log p + (1-y) \log (1-p)]$$

↓  
只有当结果预测正确才会有  $y \log p$ , 因此 log loss 对预测错误的结果惩罚非常严重.