



# 智能软件开发方向基础

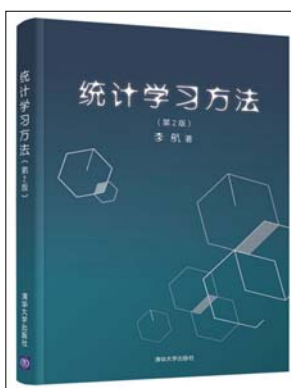
## 第三章 模型的选择与性能评价

张朝晖

2022~2023学年第二学期

序号	内容
1	概述
2	机器学习的基本概念
3	模型的选择与性能评价
4	数据的获取、探索与准备
5	近邻模型-----分类、回归
6	决策树模型-----分类、回归
7	集成学习-----分类、回归
8	(朴素)贝叶斯模型-----分类
9	聚类
10	特征降维及低维可视化(PCA, t-SNE)
11	总复习

### 课程主要参考书



### 主要内容

#### 1. 模型的学习能力与泛化能力

学习能力、泛化能力

学习误差(或训练误差)、测试误差

欠学习(欠拟合, underfitting)

过学习(过拟合, overfitting)

#### 2. 交叉验证

#### 3. (基于测试集的)模型评价

二分类/多分类模型、回归模型

### 1.1 学习能力

- 学习得到的模型关于**训练样本集**的预测能力
- 学习能力的评价:  
训练误差(或经验误差), 基于损失函数的经验风险

设学习得到的模型为:  $Y = \hat{f}(X)$   
训练样本集  $D_{\text{train}} = \{(x_i, y_i), i = 1, \dots, N_{\text{train}}\}$

训练误差:  $R_{\text{emp}}(\hat{f}) = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} L(y_i, \hat{f}(x_i))$

### 面向分类:

注意此处的指示函数

$$0-1 \text{ 损失 } L(Y, f(X)) = I(Y \neq f(X)) = \begin{cases} 1 & Y \neq f(X) \\ 0 & Y = f(X) \end{cases}$$

$$R_{\text{emp}}(\hat{f}) = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} L(y_i, \hat{f}(x_i)) = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} I(y_i, \hat{f}(x_i))$$

训练样本的  
预测错误率

### 面向回归:

平方损失函数  $L(Y, f(X)) = [Y - f(X)]^2$

$$R_{\text{emp}}(\hat{f}) = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} L(y_i, \hat{f}(x_i)) = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} [y_i - \hat{f}(x_i)]^2$$

均方误差MSE



## 1.2 泛化能力 (generalization ability)

➤ 学得模型关于 **未知样本** 的预测能力

➤ **泛化误差**:

设学习得到的模型为:  $Y = \hat{f}(X)$

该模型关于未知样本的预测误差, 即为 **泛化误差**, 它是学习得到的模型的 **期望风险**。

$$R_{\text{exp}}(\hat{f}) = \int_{\mathcal{X} \times \mathcal{Y}} L(y, \hat{f}(x)) P(x, y) dx dy$$

随机变量  $X, Y$  的函数的数学期望  
--统计意义上的损失值的平均

## 1.2 泛化能力 (generalization ability)-续

➤ **测试误差**

**泛化误差** 难以估计, 实际以 **测试误差** 评价 **模型泛化能力**

给定测试集  $D_{\text{test}} = \{(x_j, y_j), j = 1, \dots, N_{\text{test}}\}$

**测试误差**:

$$R_{\text{test}}(\hat{f}) = \frac{1}{N_{\text{test}}} \sum_{j=1}^{N_{\text{test}}} L(y_j, \hat{f}(x_j))$$

例: 对于“分类问题”的0-1损失, 测试误差就是模型关于测试集的 **预测错误率**:

$$\text{Err}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} I(y_i \neq \hat{f}(x_i))$$

与 **预测错误率** 对应的是 **预测正确率**:

$$\text{Acc}_{\text{test}} = 1 - \text{Err}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} I(y_i = \hat{f}(x_i))$$

注意此处的  
指示函数

指示函数举例:

$$I(a=b) = \begin{cases} 1 & \text{若 } a=b \\ 0 & \text{若 } a \neq b \end{cases} \quad I(a \neq b) = \begin{cases} 1 & \text{若 } a \neq b \\ 0 & \text{若 } a=b \end{cases} \quad I(a \geq b) = \begin{cases} 1 & \text{若 } a \geq b \\ 0 & \text{若 } a < b \end{cases}$$

## 1.3 过拟合与模型选择

当采用规模有限的训练集学习模型时, 因模型过于复杂, 使得学得模型对 **训练集预测性能很好**, 而关于 **测试集的预测性能很差**, 称该现象为 **过拟合 (overfitting)**。

例: 基于给定的训练集, 进行 **M次多项式函数拟合**

**M次多项式函数**:

$$f_M(x; \omega) = \omega_0 + \omega_1 x + \omega_2 x^2 + \dots + \omega_M x^M = \sum_{j=0}^M \omega_j x^j$$

当 **M** 给定时, 需要估计的参数向量为:

$$\omega = [\omega_0 \quad \dots \quad \omega_M]^T$$

训练集  $D_{\text{train}} = \{(x_i, y_i), i = 1, \dots, N_{\text{train}}\}$

测试集  $D_{\text{test}} = \{(x_k, y_k), k = 1, \dots, N_{\text{test}}\}$

采用 **经验风险最小化策略** 学习模型:

➤ **损失函数形式**----平均平方误差MSE

➤ **最小二乘准则**

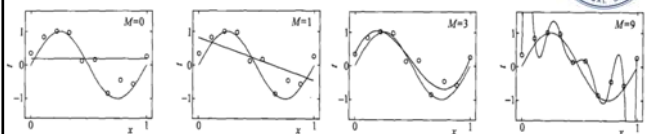
(即: 预测误差的平方和最小, 或MSE最小)

$$\begin{aligned} E(\omega) &= \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} [y_i - \hat{y}_i(\omega)]^2 \\ &= \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} [y_i - f_M(x; \omega)]^2 \end{aligned}$$

$$\omega^* = \underset{\omega}{\operatorname{argmin}} E(\omega)$$

如何求解?

讨论：不同模型复杂度 $M$ 的多项式函数拟合结果



随着模型复杂度(多项式阶数 $M$ )的增加：  
模型的训练误差逐渐减小，甚至趋0；  
但模型的测试误差先减小，后增加。

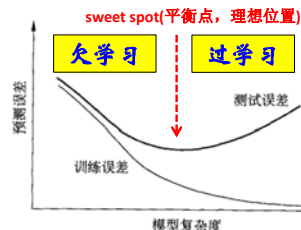
$M$ 是多项式函数拟合的超参数！

→ 多项式函数拟合，应选择合适的 $M$ 值。

对于多项式函数的拟合来说， $M$ 值的确定过程就是模型的选择过程。

右图：

训练误差、测试误差与模型复杂度之间的关系



随着模型复杂度的增加：  
训练误差逐渐减小，甚至为0；  
测试误差先减小，达到最小值后，又增大。

当模型过于复杂时，将产生过拟合。

为避免过拟合，应选择适当复杂度的模型。

常采用两种方式

➤ 方式1. 采用结构风险最小化策略，  
构建含有正则项的目标函数  
本讲暂略。

➤ 方式2. 基于交叉验证方式的模型选择

交叉验证不仅用于模型选择，还用于最终模型的评价。

## 主要内容

### 1. 模型的学习能力与泛化能力

### 2. 交叉验证

#### 2.1 数据集的划分方式

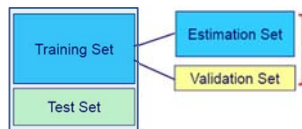
#### 2.2 基于交叉验证的模型选择与评价

### 3. (基于测试集的)模型评价

二分类/多分类模型、回归模型

## 数据集 $D$ 的两种层次划分及适用场合

以监督式模型学习为例



选择模型(超参数)后再对所有的训练数据进行训练

完整的模型学习过程

首先，模型的选择(待学习的模型的超参数调节)

训练集 { 估计集(学习)  
验证集(评价) }

然后，基于选择的模型结构，学习得到最终模型

{ 训练集(超参数固定之后，用于学习最终模型)  
测试集(评价最终模型) }

## 主要内容

### 1. 模型的学习能力与泛化能力

### 2. 交叉验证

#### 2.1 数据集的划分方式

#### 2.2 基于交叉验证的模型选择与评价

### 3. (基于测试集的)模型评价

二分类/多分类模型、回归模型、聚类模型的评价

## sklearn.model\_selection: Model Selection

**User guide:** See the Cross-validation: evaluating estimator performance, Tuning the hyper-parameters of an estimator and Learning curve sections for further details.

### Splitter Classes

<code>model_selection.GroupKFold(n_splits)</code>	K-fold iterator variant with non-overlapping groups.
<code>model_selection.GroupShuffleSplit(...)</code>	Shuffle-Group(s)-Out cross-validation iterator
<code>model_selection.KFold(n_splits, shuffle, ...)</code>	K-Folds cross-validator
<code>model_selection.LeaveOneGroupOut()</code>	Leave One Group Out cross-validator
<code>model_selection.LeavePGroupsOut(n_groups)</code>	Leave P Group(s) Out cross-validator
<code>model_selection.LeaveOneOut()</code>	Leave-One-Out cross-validator
<code>model_selection.LeavePOut(p)</code>	Leave-P-Out cross-validator
<code>model_selection.PredefinedSplit(test_fold)</code>	Predefined split cross-validator
<code>model_selection.RepeatedKFold(k, n_splits, ...)</code>	Repeated K-Fold cross validator.
<code>model_selection.RepeatedStratifiedKFold(k, ...)</code>	Repeated Stratified K-Fold cross validator.
<code>model_selection.ShuffleSplit(n_splits, ...)</code>	Random permutation cross-validator
<code>model_selection.StratifiedKFold(n_splits, ...)</code>	Stratified K-Folds cross-validator.
<code>model_selection.StratifiedShuffleSplit(...)</code>	Stratified ShuffleSplit cross-validator
<code>model_selection.StratifiedGroupKFold(...)</code>	Stratified K-Folds iterator variant with non-overlapping groups.
<code>model_selection.TimeSeriesSplit(n_splits, ...)</code>	Time Series cross-validator

### 思考:

- 如何基于交叉验证, 进行模型的选择?
- K-fold CV用于模型选择/评价的具体实现过程?
- K-fold CV 中的K值一般为多少?
- 何时使用LOOCV? 何时使用K-fold CV?
- 如何基于交叉验证, 进行模型的评价?
- 最终要使用的模型应怎么得到?



以监督式学习系统为例, 首先**考察最终模型泛化能力的评估方式**  
然后**再理解如何在生成最终模型之前, 用交叉验证法进行模型选择**。

给定已知答案的**数据集** $D = \{(x_i, y_i), i = 1, \dots, N\}$

{ 训练集  $D_{train}$  --模型的学习  
测试集  $D_{test}$  --模型泛化能力的评价

$$D = D_{train} \cup D_{test}$$

**数据集D**划分的几种典型实现方式

- [1]留法 (**hold-out**)
- [2]交叉验证 (**cross validation**)
- [3]自助法 (**bootstrapping**)

#### 如何划分:

- “**随机打乱**”——回归
- **分属随机打乱**——分类

[1]留法 (**hold-out**)、留法交叉验证 (**hold-out cv**)

$$D = D_{train} \cup D_{test}$$

$$\Phi = D_{train} \cap D_{test}$$

数据集**随机划分**尽量**保持数据分布的一致性**

- 单独一次随机划分, 估计结果不够稳定可靠
  - 应多次随机划分, 重复评估, **取结果的均值及标准差!!!**
- (**hold-out cross-validation**)

$D_{train}$	$D_{test}$
$\frac{2}{3} \sim \frac{3}{4}$	$\frac{1}{4} \sim \frac{1}{3}$

$D_{train}$	$D_{test}$
$\frac{2}{3} \sim \frac{3}{4}$	$\frac{1}{4} \sim \frac{1}{3}$

矛盾:

- **训练集规模应足够大**, 以便使模型的学习, 尽可能在模型学习过程中, 让其见识更为丰富的样本多样性; 但会导致测试集内样本多样性降低, 基于测试集的评价结果不够稳定、准确, 无法较好地近似模型的泛化能力。
- **测试集规模应足够大**, 以便使测试误差更接近泛化误差; 但会导致训练集规模降低, 使得训练得到的模型与基于整个数据集得到的模型差别较大, 降低了评估结果的保真性。

[2]**k-倍交叉验证**

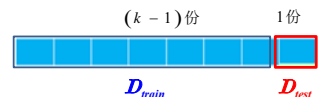
(**k-fold cross validation, k-fold CV**)—rotation estimation

$$D = D_{train} \cup D_{test} \quad \Phi = D_{train} \cap D_{test}$$

数据集的随机划分尽量**保持数据分布的一致性**

随机打乱  $D$ , 均分成  $k$  等份

$k = ?$  3, 5, 10



- { **单轮**  $k$ -倍交叉验证
- { **多轮**  $k$ -倍交叉验证

**悲观、有偏估计**

**留一法交叉验证** (**leave-one-out cross-validation, LOO-CV**)

**无偏估计**

例：基于单轮K-折交叉验证错误率，评价某分类模型的预测性能

STEP1. 将样本集D分层随机打乱，均分成互斥的K等份： $D = D_1 \cup \dots \cup D_K$

STEP2. for  $j=1,2,\dots,K$ , do:

以样本集 $D_j$ 为测试集 $D_{test}$ ，其余部分构成训练集 $D_{train}=D \setminus D_j$

以 $D_{train}$ 学习得到模型 $model_j$ 评价 $D_{test}$ ，得到测试错误率 $Error_j$

STEP3. 估计上述K个错误率的均值及标准差:

$$\mu(Error) = \frac{1}{K} \sum_{j=1}^K Error_j$$

$$\sigma(Error) = \left[ \frac{1}{K} \sum_{j=1}^K (Error_j - \mu(Error))^2 \right]^{\frac{1}{2}}$$

STEP4. 输出K-折交叉验证的错误率估计结果： $\mu(Error) \pm \sigma(Error)$

注意：对于分类问题，面向交叉验证的数据集划分应采用分层方式(即：Stratified k-fold cross validation)，确保训练集、测试集的类别分布一致。

例：Comparison between Support Vector Machines, the Kernel Fisher Discriminant (KFD), a single radial basis function classifier (RBF), AdaBoost (AB), and regularized AdaBoost (ABR) on 13 different benchmark datasets (see text). Best result in bold face, second best in *italics*.

平均错误率  $\pm$  错误率标准差

	SVM	KFD	RBF	AB	ABR
Banana	11.5 $\pm$ 0.07	<b>10.8<math>\pm</math>0.05</b>	<b>10.8<math>\pm</math>0.06</b>	12.3 $\pm$ 0.07	<i>10.9<math>\pm</math>0.04</i>
B. Cancer	<i>26.0<math>\pm</math>0.47</i>	<b>25.8<math>\pm</math>0.46</b>	27.6 $\pm$ 0.47	30.4 $\pm$ 0.47	26.5 $\pm$ 0.45
Diabetes	<i>23.5<math>\pm</math>0.17</i>	<b>23.2<math>\pm</math>0.16</b>	24.3 $\pm$ 0.19	26.5 $\pm$ 0.23	23.8 $\pm$ 0.18
German	<b>23.6<math>\pm</math>0.21</b>	<i>23.7<math>\pm</math>0.22</i>	24.7 $\pm$ 0.24	27.5 $\pm$ 0.25	24.3 $\pm$ 0.21
Heart	<b>16.0<math>\pm</math>0.33</b>	<i>16.1<math>\pm</math>0.34</i>	17.6 $\pm$ 0.33	20.3 $\pm$ 0.34	16.5 $\pm$ 0.35
Image	<i>3.0<math>\pm</math>0.06</i>	3.3 $\pm$ 0.06	3.3 $\pm$ 0.06	<b>2.7<math>\pm</math>0.07</b>	<b>2.7<math>\pm</math>0.06</b>
Ringnorm	<b>1.7<math>\pm</math>0.01</b>	<b>1.5<math>\pm</math>0.01</b>	1.7 $\pm$ 0.02	1.9 $\pm$ 0.03	<i>1.6<math>\pm</math>0.01</i>
F. Sonar	<b>32.4<math>\pm</math>0.18</b>	<i>33.2<math>\pm</math>0.17</i>	34.4 $\pm$ 0.20	35.7 $\pm$ 0.18	34.2 $\pm$ 0.22
Splice	10.9 $\pm$ 0.07	10.5 $\pm$ 0.06	<i>10.0<math>\pm</math>0.10</i>	10.1 $\pm$ 0.05	<b>9.5<math>\pm</math>0.07</b>
Thyroid	4.8 $\pm$ 0.22	<b>4.2<math>\pm</math>0.21</b>	4.5 $\pm$ 0.21	<i>4.4<math>\pm</math>0.22</i>	4.6 $\pm$ 0.22
Titanic	<b>22.4<math>\pm</math>0.10</b>	23.2 $\pm$ 0.20	23.3 $\pm$ 0.13	<i>22.6<math>\pm</math>0.12</i>	<i>22.6<math>\pm</math>0.12</i>
Twonorm	3.0 $\pm$ 0.02	<b>2.6<math>\pm</math>0.02</b>	2.9 $\pm$ 0.03	3.0 $\pm$ 0.03	<i>2.7<math>\pm</math>0.02</i>
Waveform	<i>9.9<math>\pm</math>0.04</i>	<i>9.9<math>\pm</math>0.04</i>	10.7 $\pm$ 0.11	10.8 $\pm$ 0.06	<b>9.8<math>\pm</math>0.08</b>

河北师范大学软件学院  
Software College of Hebei Normal University

例：基于留一法折交叉验证正确率，评价某分类模型的预测性能

给定样本集 $D = \{(x_i, y_i), i = 1, \dots, N\}$

STEP1. 初始化计数器Counter=0

STEP2. for  $j=1,2,\dots,N$ , do:

以第j个样本 $(x_j, y_j)$ 为测试样本，其余样本构成训练集

$D_{train} = D \setminus \{(x_j, y_j)\}$

以 $D_{train}$ 学习得到模型 $model_j$ 预测 $x_j$ 的类别输出为 $\hat{y}_j$

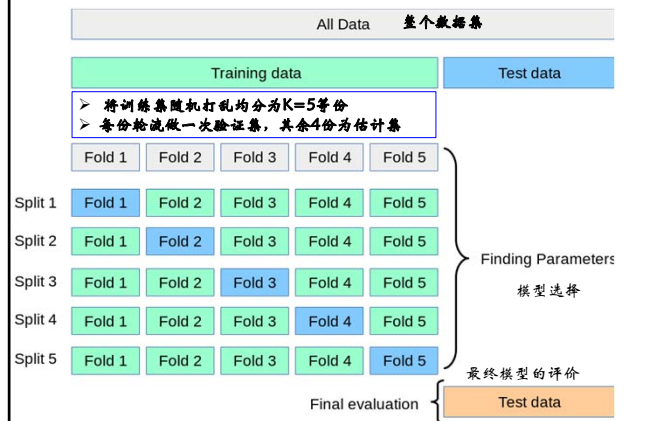
若 $\hat{y}_j$ 与 $y_j$ 一致，则Counter=Counter+1

STEP3. 输出留一法交叉验证的正确率:

$$Acc = \frac{Counter}{N} \times 100\%$$

河北师范大学软件学院  
Software College of Hebei Normal University

基于单轮K-折交叉验证的模型选择



[3]自助法 (bootstrapping)

bootstrap sampling

数据集 $D = \{(x_i, y_i), i = 1, \dots, N\}$  -- 初始数据集

训练集 $D_{train}$  -- 模型的学习

方式1. 对初始数据集D有放回的随机抽取N次，每次抽取1个样本  
得自助数据集 $D'$

方式2. 从原始数据集D内随机抽取样本“若干(?)”

测试集 $D_{test}$  -- 模型泛化能力的评价

对数据集D中没有被抽取到的样本集 $D_{test} = D \setminus D'$

初始数据集内，样本未被抽取的概率 $\lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^N = \frac{1}{e} \approx 0.368$

自助法估计，也称“包外估计”(out-of-bag estimation)，例：包外错误率

河北师范大学软件学院  
Software College of Hebei Normal University

自助法常用于集成学习(如：bagging、随机森林)时，集成模型的性能评价。基本过程：

➤ 每个版本的训练集 $D_{train}$ 对应一个个体模型的学习

➤ 针对样本集D内每个样本，当其作为某个版本训练集的包外样本时，基于相应的个体模型，对该包外测试样本的输出预测，记录其预测结果及其作为包外测试样本的次数。

➤ 综合样本集D内每个样本的预测结果

例：对于分类问题，以投票方式决定每个样本的最终预测类

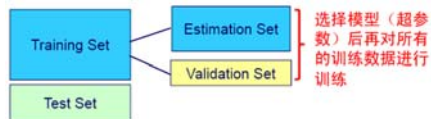
➤ 最终，生成集成模型的最终评价结果

河北师范大学软件学院  
Software College of Hebei Normal University



思考：如何基于交叉验证方式，进行模型的选择？

例：基于给定训练集  $D_{train} = \{(x_i, y_i), i = 1, \dots, N\}$ ，针对  $C=3$  的分类问题，采用交叉验证方式进行KNN分类模型的选择

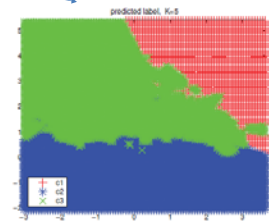
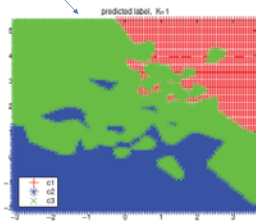


训练集

近邻数K值不同时的分类结果

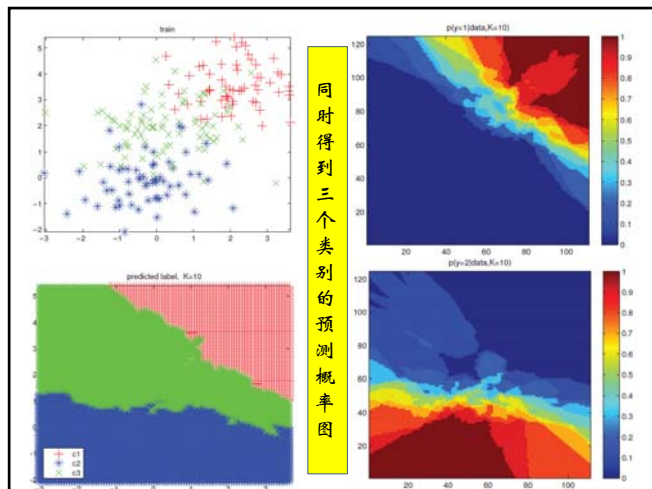
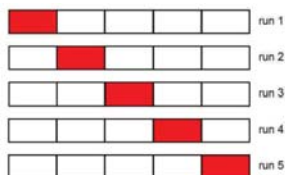
K=1

K=5



基于单轮5-fold cross validation 进行近邻数K值的优选

训练集=估计集∪验证集



同时得到三个类别的预测概率图

## 主要内容

1. 模型的学习能力与泛化能力
2. 交叉验证
3. (基于测试集的)模型评价

### 3.1 基于测试集的二分类模型评价

两种情况：(1)非对称类别：是vs.非

(2)对称类别：第1类vs.第2类

### 3.2 基于测试集的多分类模型的评价

### 3.3 基于测试集的实值函数预测模型评价

[https://scikit-learn.org/stable/modules/model\\_evaluation.html](https://scikit-learn.org/stable/modules/model_evaluation.html)

Scoring Classification	Function	Comment
'accuracy'	metrics.accuracy_score	
'balanced_accuracy'	metrics.balanced_accuracy_score	
'top_k_accuracy'	metrics.top_k_accuracy_score	
'average_precision'	metrics.average_precision_score	
'neg_brier_score'	metrics.brier_score_loss	
'f1'	metrics.f1_score	for binary targets
'f1_micro'	metrics.f1_score	micro-averaged
'f1_macro'	metrics.f1_score	macro-averaged
'f1_weighted'	metrics.f1_score	weighted average
'f1_samples'	metrics.f1_score	by multilabel sample
'neg_log_loss'	metrics.log_loss	requires predict_proba support
'precision' etc.	metrics.precision_score	suffixes apply as with 'f1'
'recall' etc.	metrics.recall_score	suffixes apply as with 'f1'
'jaccard' etc.	metrics.jaccard_score	suffixes apply as with 'f1'
'roc_auc'	metrics.roc_auc_score	
'roc_auc_ovr'	metrics.roc_auc_score	
'roc_auc_ovo'	metrics.roc_auc_score	
'roc_auc_ovr_weighted'	metrics.roc_auc_score	
'roc_auc_ovo_weighted'	metrics.roc_auc_score	

## Classification metrics

See the Classification metrics section of the user guide for further details.

<code>metrics.accuracy_score(y_true, y_pred, *, labels=None)</code>	Accuracy classification score.
<code>metrics.auc(x, y)</code>	Compute Area Under the Curve (AUC) using the trapezoidal rule.
<code>metrics.average_precision_score(y_true, y_score)</code>	Compute average precision (AP) from prediction scores.
<code>metrics.balanced_accuracy_score(y_true, y_score)</code>	Compute the balanced accuracy.
<code>metrics.brier_score_loss(y_true, y_prob, *, labels=None)</code>	Compute the Brier score loss.
<code>metrics.classification_report(y_true, y_pred, *, labels=None, target_names=None, digits=2)</code>	Build a text report showing the main classification metrics.
<code>metrics.cohen_kappa_score(y1, y2, *, labels=None)</code>	Cohen's kappa: a statistic that measures inter-annotator agreement.
<code>metrics.confusion_matrix(y_true, y_pred, *, labels=None)</code>	Compute confusion matrix to evaluate the accuracy of a classification.
<code>metrics.dcg_score(y_true, y_score, *, k=10)</code>	Compute Discounted Cumulative Gain.
<code>metrics.det_curve(y_true, y_score, *, labels=None)</code>	Compute error rates for different probability thresholds.
<code>metrics.f1_score(y_true, y_pred, *, labels=None)</code>	Compute the F1 score, also known as balanced F-score or F-measure.
<code>metrics.fbeta_score(y_true, y_pred, *, beta)</code>	Compute the F-beta score.
<code>metrics.hamming_loss(y_true, y_pred, *, labels=None)</code>	Compute the average Hamming loss.
<code>metrics.hinge_loss(y_true, y_pred, *, labels=None)</code>	Average hinge loss (non-regularized).
<code>metrics.jaccard_score(y_true, y_pred, *, labels=None)</code>	Jaccard similarity coefficient score.
<code>metrics.log_loss(y_true, y_prob, *, labels=None)</code>	Log loss, aka logistic loss or cross-entropy loss.
<code>metrics.matthews_corrcoef(y_true, y_pred, *, labels=None)</code>	Compute the Matthews correlation coefficient (MCC).

<code>metrics.matthews_corrcoef(y_true, y_pred, *, labels=None)</code>	Compute the Matthews correlation coefficient (MCC).
<code>metrics.multilabel_confusion_matrix(y_true, y_score, *, labels=None)</code>	Compute a confusion matrix for each class or sample.
<code>metrics.ndcg_score(y_true, y_score, *, k=10)</code>	Compute Normalized Discounted Cumulative Gain.
<code>metrics.precision_recall_curve(y_true, y_score)</code>	Compute precision-recall pairs for different probability thresholds.
<code>metrics.precision_recall_fscore_support(y_true, y_score, *, labels=None)</code>	Compute precision, recall, F-measure and support for each class.
<code>metrics.precision_score(y_true, y_pred, *, labels=None)</code>	Compute the precision.
<code>metrics.recall_score(y_true, y_pred, *, labels=None)</code>	Compute the recall.
<code>metrics.roc_auc_score(y_true, y_score, *, labels=None)</code>	Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores.
<code>metrics.roc_curve(y_true, y_score, *, labels=None)</code>	Compute Receiver operating characteristic (ROC).
<code>metrics.top_k_accuracy_score(y_true, y_score, *, k=1)</code>	Top-k Accuracy classification score.
<code>metrics.zero_one_loss(y_true, y_pred, *, labels=None)</code>	Zero-one classification loss.

## 二分类的情况之一： 非对称类别 (是 vs. 非)

(人脸 vs. 非人脸；垃圾邮件 vs. 非垃圾邮件)

- 什么是混淆矩阵？如何计算混淆矩阵？
- 如何利用混淆矩阵生成若干评价指标？
- 什么是ROC曲线？如何利用有限规模测试样本绘制ROC曲线？
- 如何基于ROC曲线进行模型的评价与比较？
- 如何计算ROC曲线下面积(即：AUC值，也称AUC<sub>ROC</sub>)？
- 什么是P-R曲线？如何利用有限规模测试样本绘制P-R曲线？
- 如何基于P-R曲线实现模型的定性或定量评价？
- 如何计算P-R曲线下面积(即：AUC<sub>PR</sub>)？
- 哪种曲线更适合非对称类别分类模型的性能评价？为什么？

## 方式1. 基于混淆矩阵的两类别分类模型的性能评价

## 方式2. 基于ROC曲线/P-R曲线的性能评价

## 方式1. 基于混淆矩阵的两类别分类模型的性能评价

混淆矩阵：confusion matrix

### (1) 样本的两种自然类别状态

通常设定感兴趣的一类为正类

类别标号	类别各种名称			
1	Positive (P)	正	阳性	Case Samples (病理样本)
0或-1或2	Negative (N)	负	阴性	Control Samples (对照样本)

## (2) 两类决策的混淆矩阵(confusion matrix)

正确分类；错误分类；决策阈值

自然状态 True Value Actual Value	预测输出(Predicted Outcome)	
	Positive (Predicated 1)	Negative (Predicated 0 or -1)
Positive (True 1)	<b>a</b> True Positive (TP) 真阳性 真正类 Hits	<b>b</b> False Negative (FN) 假阴性 Misses 假负类
Negative (True 0 或 -1)	<b>c</b> False Positive (FP) 假阳性 False Alarms	<b>d</b> True Negative (TN) 真阴性 True Rejections 真负类

## 概率形式(正确分类; 错误分类)---条件概率

自然状态 True Value Actual Value	预测输出(Predicated Outcome)	
	Positive (Predicated 1)	Negative (Predicated 0 or -1)
Positive (True 1)	$S_n$ -灵敏度 $P(\text{Predicated } P   \text{True } P)$	$\beta$ -假阴性率 $P(\text{Predicated } N   \text{True } P)$
Negative (True 0 或 -1)	$\alpha$ -假阳性率 $P(\text{Predicated } P   \text{True } N)$	$Sp$ -特异度 $P(\text{Predicated } N   \text{True } N)$

$$P(\text{Predicated } P | \text{True } P) + P(\text{Predicated } N | \text{True } P) = 1$$

$$P(\text{Predicated } P | \text{True } N) + P(\text{Predicated } N | \text{True } N) = 1$$

## (3)基于两类别混淆矩阵的评价指标

- [1](阳性类) 查准率,精度(Precision,P)  $Precision = \frac{TP}{TP + FP}$
- [2](阳性类) 查全率,召回率(Recall,R),灵敏度(Sensitivity)  
命中率, 真阳性率(TruePositiveRate)  $S_n = Recall = \frac{TP}{TP + FN}$
- [3] 特异度(Specificity), 真阴性率(TrueNegativeRate)  $S_p = \frac{TN}{TN + FP}$
- [4] 假阴性率(FalseNegativeRate),漏报率(MissedRate)  $\beta = \frac{FN}{FN + TP}$   
第2类错误率
- [5] 假阳性率(FalsePositiveRate),虚警率(FalseAlarmRate)  $\alpha = \frac{FP}{TN + FP}$   
第1类错误率

## [6](阳性类的) $F_\beta$ -Score以及 $F_1$ -Score

----Precision与Recall的调和平均

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

$0 < \beta < 1$ , 更看重P

$\beta > 1$ , 更看重R

$$\beta = 1 \text{ 时, } F_1 = \frac{2}{\frac{1}{R} + \frac{1}{P}} = \frac{2P \cdot R}{P + R} = \frac{2TP}{2TP + FP + FN} = \frac{2TP}{\text{测试样本总数} + TP + TN}$$

## [7]马修相关系数(Matthews Correlation Coefficient,MCC)

类别不均衡的两类别分类, 公平度量模型预测能力

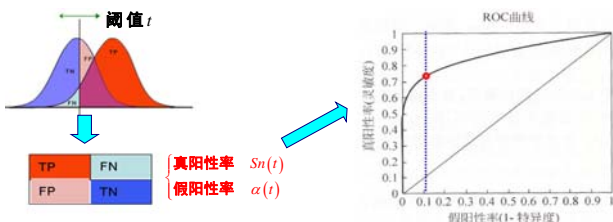
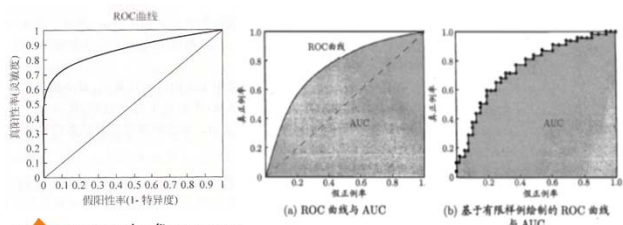
$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FN)(TN + FP)(TP + FP)(TN + FN)}}$$

方式1.基于混淆矩阵的两类别分类模型的性能评价

方式2.基于ROC曲线/P-R曲线的性能评价

## (1)ROC曲线(Receiver Operating Characteristic Curves)及绘制

The ROC curve is a graphical plot of the sensitivity, or true positive rate (y-axis), vs. false positive rate (1 - specificity or 1 - true negative rate) (x-axis), for a binary classifier system as its discrimination threshold is varied.



决策阈值  $\Rightarrow$

特异度(Specificity, 真阴性率)  $Sp(t) = \frac{TN(t)}{FP(t) + TN(t)}$

假阴性率(第二类错误率)  $\beta(t) = \frac{FN(t)}{FN(t) + TP(t)}$

假阳性率(第一类错误率, 误报率)  $\alpha(t) = 1 - Sp(t) = \frac{FP(t)}{FP(t) + TN(t)}$

真阳性率(Sensitivity, 敏感度, 召回率)  $S_n(t) = 1 - \beta(t) = \frac{TP(t)}{FN(t) + TP(t)}$



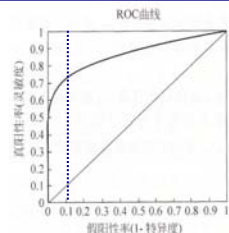
## (2)基于ROC曲线的分类器性能评价

第一, 根据分类器的设计指标要求, 由ROC曲线确定真实分类器的工作点

真条件密度已知, 改变决策阈值, 得ROC曲线:

ROC曲线上每个点, 均对应一对(灵敏度, 特异度);

根据分类器对两种指标要求, 确定曲线工作点



ROC曲线越往左上靠, ROC曲线越优,

存在一个平衡点, 使这个平衡点就是max(TPR-FPR)所对应的分类器阈值

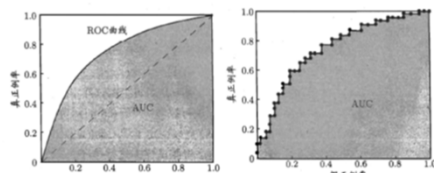
## (2)基于ROC曲线的分类器性能评价

第二, 类别数=2的分类器性能评价-AUC

AUC-ROC曲线下的面积AUC (Area Under ROC Curves)

ROC曲线越靠近坐标系的左上角, AUC值越大, 分类器性能越好

$$AUC = \frac{1}{2} \sum_{i=1}^{N-1} ((\alpha_{i+1} - \alpha_i)(R_{i+1} + R_i))$$



河北师范大学软件学院  
Software College of Hebei Normal University

例: 疾病诊断

“positive(阳性类)”-疾病

“negative(阴性类)”-正常

(1)不希望漏诊疾病

希望结果: 真阳性率=1

不希望结果: 真阴性率降低(极端值=0) 假阳性率升高(极端值=1)

(2)不希望误判为疾病

希望结果: 真阴性率=1, 假阳性率=0

不希望结果: 真阳性率降低(极端值=0) 假阴性率升高(极端值=1)

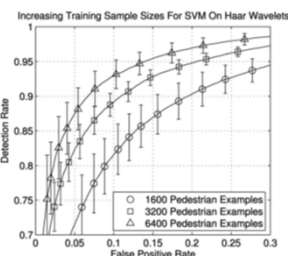
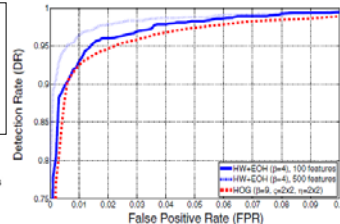
ROC曲线靠近左上角

河北师范大学软件学院  
Software College of Hebei Normal University

## (2)基于ROC曲线的分类器性能评价

可用于不同分类器性能比较

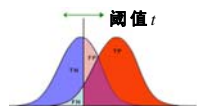
可针对同种分类器, 进行特征选择、样本规模等优劣比较。



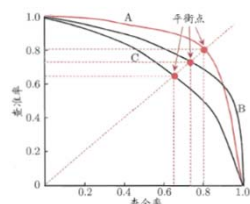
可以使用AUC值, 单独评价每个特征的真鉴别能力

可以结合多个特征的多条ROC曲线横向比较各特征的真鉴别性

## (3)P-R曲线(Precision-Recall Curve)及其AUC-PR



$$\begin{cases} \text{precision}(t) = \frac{TP(t)}{TP(t) + FP(t)} \\ \text{recall}(t) = \frac{TP(t)}{TP(t) + FN(t)} \end{cases}$$



定性比较: 如果一个学习器的P-R曲线被另一学习器的P-R曲线完全包住, 可断言后者性能优于前者。例如: A、B学习器优于学习器C。

定量评价方式1: 借助不同P-R曲线下方的面积值, 进行横向比较

定量评价方式2: 结合平衡点或者F1值。

平衡点(BEP, break even point)是P-R曲线与直线P=R的交点。

P=R值越大, 学习器性能越好。

F1 = 2 \* P \* R / (P + R)。F1值越大, 该学习器性能越好。

二分类的情况之二:  
对称类别(第1类vs.第2类)

(例: 猫vs.狗)

什么是混淆矩阵? 该混淆矩阵与“非对称的二分类”混淆矩阵的区别?

如何产生混淆矩阵? 如何基于混淆矩阵产生评价指标?

什么是ROC曲线? 如何利用有限规模测试样本绘制ROC曲线?

对于对称类别的二分类问题, 一个测试集可以同时得到几条ROC曲线?

如何结合这些ROC曲线进行分类模型的性能评价?

河北师范大学软件学院  
Software College of Hebei Normal University

## (1) 混淆矩阵(confusion matrix)

第1类；第2类；决策阈值

$$A = [a_{ij}]_{2 \times 2}$$

真实类别	预测类别	
	第1类	第2类
第1类	$a_{11}$	$a_{12}$
第2类	$a_{21}$	$a_{22}$

$a_{ij}$ ----参与测试的样本中，真实类别为第i类，但预测为第j类的样本数

$$D_{\text{test}} = \{(x_i, y_i), i = 1, \dots, N_{\text{test}}\}$$

总体预测错误率：

$$Err_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} I(y_i \neq \hat{y}_i) \times 100\%$$

总体预测正确率：

$$Acc_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{i=1}^{N_{\text{test}}} I(y_i = \hat{y}_i) \times 100\%$$

## (2) 基于两类别混淆矩阵的评价指标

[1] 总体正确率, 准确率, 准确度 (Overall Accuracy)  $Accuracy = \frac{\sum_{i=1}^2 \sum_{j=1}^2 a_{ij}}{\sum_{i=1}^2 \sum_{j=1}^2 a_{ij}} \times 100\%$

[2] 总体错误率 (ErrorRate)  $Error = \frac{a_{12} + a_{21}}{\sum_{i=1}^2 \sum_{j=1}^2 a_{ij}} \times 100\% = 1 - Accuracy$

[3] 第i类查准率, 精度 (Precision, P)  $P_i = \frac{a_{ii}}{\sum_{j=1}^2 a_{ji}} \quad i = 1, 2$

[4] 第i类查全率, 召回率 (Recall, R), 正确率  $R_i = \frac{a_{ii}}{\sum_{j=1}^2 a_{ij}} = Acc_i \quad i = 1, 2$

[5] 第i类  $F_\beta$ -Score 以及  $F_1$ -Score --  $P_i$  与  $R_i$  的调和平均

$$F_\beta^{(i)} = \frac{(\beta^2 + 1) P_i \cdot R_i}{\beta^2 P_i + R_i} \quad i = 1, 2$$

$0 < \beta < 1$ , 更看重  $P_i$

$\beta > 1$ , 更看重  $R_i$

$$\beta = 1 \text{ 时, } F_1^{(i)} = \frac{2 P_i \cdot R_i}{P_i + R_i}$$

[6] 宏查准率、宏查全率  $Macro\_P = \frac{1}{2} \sum_{i=1}^2 P_i \quad Macro\_R = \frac{1}{2} \sum_{i=1}^2 R_i$

[7]  $C=2$  个类别的宏平均  $F_1$ -Score (Macro-averaging  $F_1$ )

$$Macro\_F_1 = \frac{2 \times Macro\_P \times Macro\_R}{Macro\_P + Macro\_R}$$

[8] 各类预测正确率的几何平均--类别不平衡的评价

$$Acc_{G-mean} = \left( \prod_{i=1}^2 Acc_i \right)^{\frac{1}{2}}$$

[9] 马修相关系数

(Matthews Correlation Coefficient, MCC)

类别不平衡的两类别分类，公平度量模型预测能力

$$MCC = \frac{a_{11} \cdot a_{22} - a_{12} \cdot a_{21}}{\sqrt{(a_{11} + a_{12})(a_{21} + a_{22})(a_{11} + a_{21})(a_{12} + a_{22})}}$$

## 主要内容

1. 模型的学习能力与泛化能力

2. 交叉验证

3. (基于测试集的)模型评价

3.1 基于测试集的二分类模型评价

3.2 基于测试集的多分类模型的评价

3.3 基于测试集的实值函数预测模型评价



## 多类别分类(C>2)

- 如何结合测试集的类别标签信息以及预测标签，产生多类别分类的混淆矩阵？
- 如何结合C阶混淆矩阵，产生若干评价指标？
- 一个C阶混淆矩阵可生成C个“第i类vs.非第i类”二阶混淆矩阵
- 一个测试集可以同时得到几条ROC曲线？
- 如何结合这些ROC曲线进行分类模型的性能评价？

## (1) 混淆矩阵(confusion matrix)

C阶矩阵： $A = [a_{ij}]_{C \times C}$

$a_{ij}$ ----测试样本集内，真实类别为第i类，但预测为第j类的样本数

真实类别为第i类的样本总数： $\sum_{j=1}^C a_{ij}$

预测类别为第j类的样本总数： $\sum_{i=1}^C a_{ij}$

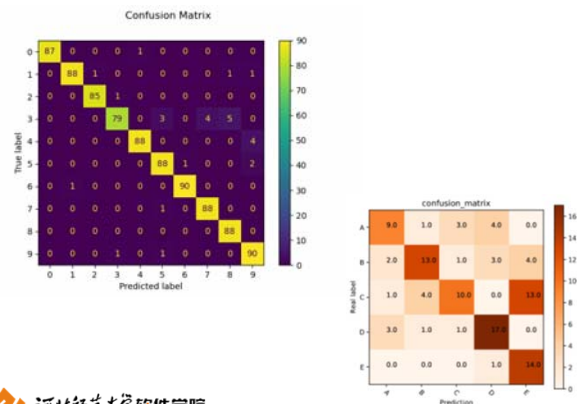
测试样本总数： $\sum_{i=1}^C \sum_{j=1}^C a_{ij}$

例：基于混淆矩阵的手写体数字识别

混淆矩阵为10行\*10列

true class i	class j predicted by a classifier									
	'0'	'1'	'2'	'3'	'4'	'5'	'6'	'7'	'8'	'9'
'0'	97	0	0	0	0	0	1	0	0	1
'1'	0	98	0	0	1	0	0	1	0	0
'2'	0	0	96	1	0	1	0	1	0	0
'3'	0	0	2	95	0	1	0	0	1	0
'4'	0	0	0	0	98	0	0	0	0	2
'5'	0	0	0	1	0	97	0	0	0	0
'6'	1	0	0	0	0	1	98	0	0	0
'7'	0	0	1	0	0	0	0	98	0	0
'8'	0	0	0	1	0	0	1	0	96	1
'9'	1	0	0	0	3	1	0	0	0	95

sklearn.metrics.ConfusionMatrixDisplay



测试集  $D_{test} = \{(x_i, y_i), i = 1, \dots, N_{test}\}$

总体预测错误率：

$$Err_{test} = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} I(y_i \neq \hat{y}_i) \times 100\%$$

总体预测正确率：

$$Acc_{test} = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} I(y_i = \hat{y}_i) \times 100\%$$

## (2) 基于多类别混淆矩阵的评价指标

[1] 总体正确率, 准确率, 准确度 (Overall Accuracy)  $Accuracy = \frac{\sum_{i=1}^C a_{ii}}{\sum_{i=1}^C \sum_{j=1}^C a_{ij}} \times 100\%$

[2] 总体错误率 (ErrorRate)  $Error = 1 - Accuracy$

[3] 第i类查准率, 精度 (Precision, P)  $P_i = \frac{a_{ii}}{\sum_{j=1}^C a_{ij}} \quad i = 1, \dots, C$

[4] 第i类查全率, 召回率 (Recall, R), 正确率  $R_i = \frac{a_{ii}}{\sum_{j=1}^C a_{ji}} = Acc_i \quad i = 1, \dots, C$

[5] 宏查准率、宏查全率  $Macro\_P = \frac{1}{C} \sum_{i=1}^C P_i \quad Macro\_R = \frac{1}{C} \sum_{i=1}^C R_i$

[6]第*i*类  $F_\beta$ -Score 以及  $F_1$ -Score --  $P_i$  与  $R_i$  的调和平均

$$F_\beta^{(i)} = \frac{(\beta^2 + 1)P_i \cdot R_i}{\beta^2 P_i + R_i} \quad i = 1, \dots, C$$

$0 < \beta < 1$ , 更看重  $P_i$

$\beta > 1$ , 更看重  $R_i$

$$\beta = 1 \text{ 时, } F_1^{(i)} = \frac{2P_i \cdot R_i}{P_i + R_i}$$

[7]  $C$  个类别的宏平均  $F_1$ -Score (Macro-averaging  $F_1$ )

$$Macro\_F_1 = \frac{2 \times Macro\_P \times Macro\_R}{Macro\_P + Macro\_R}$$

[8] 由  $C$  阶混淆矩阵, 得到  $C$  个“第  $i$  类 vs. 非第  $i$  类”混淆矩阵  
将这  $C$  个混淆矩阵取平均, 得到平均混淆矩阵

微查准率/微精度 (Precision)  $micro\_P = \frac{\overline{TP}}{\overline{TP} + \overline{FP}}$

微查全率/微召回率 (Recall)  $micro\_R = \frac{\overline{TP}}{\overline{TP} + \overline{FN}}$

微  $F_1$ -Score  $micro\_F_1 = \frac{2micro\_P \cdot micro\_R}{micro\_P + micro\_R}$

[9] Kappa 系数 (Cohen's kappa)

Cohen's kappa measures the agreement between two raters who each classify  $N$  items into  $C$  mutually exclusive categories.

$$K = \frac{p_0 - p_e}{1 - p_e} = \frac{\frac{\sum_{i=1}^C \sum_{j=1}^C a_{ij}}{C^2} - \left( \frac{\sum_{i=1}^C \sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}}{\left( \sum_{i=1}^C \sum_{j=1}^C a_{ij} \right)^2} \right)}{1 - \left( \frac{\sum_{i=1}^C \sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}}{\left( \sum_{i=1}^C \sum_{j=1}^C a_{ij} \right)^2} \right)} \quad -1 \leq K \leq 1$$

考虑两个极端:  
没有错分; 完全错分

理想情况, 完全一致,  $K = 1$ ;  
 $K$  值越小, 越不一致, 甚至为负。

基于 kappa 系数一致性的不同等级

- > [0.0-0.20] 轻微一致性 (slight)
- > [0.20-0.40] 一般一致性 (fair)
- > [0.40-0.60] 中等的一致性 (moderate)
- > [0.6-0.80] 高度的一致性 (substantial)
- > [0.8-1] 几乎完全一致 (almost perfect)

关于  $p_0, p_e$  的具体说明:

(1)  $p_0$  -- 实际一致率 (总体预测正确率)  $p_0 = \frac{\sum_{i=1}^C a_{ii}}{\sum_{i=1}^C \sum_{j=1}^C a_{ij}}$

(2)  $p_e$  -- 理论一致率 (机遇一致性)  
 $1 - p_e$  -- 非机遇一致性

$p_e$  表示一个样本被抽取的类别与决策类别一致的概率

$$p_e = \sum_{i=1}^C p_{ei} = \frac{\sum_{i=1}^C \left( \frac{\sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}}{\sum_{i=1}^C \sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}} \right)}{\sum_{i=1}^C \left( \frac{\sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}}{\sum_{i=1}^C \sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}} \right)} = \sum_{i=1}^C \left( \frac{\sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}}{\sum_{i=1}^C \sum_{j=1}^C a_{ij} \sum_{k=1}^C a_{ki}} \right)$$

第  $i$  类样本被观测到的概率

样本决策为第  $i$  类的概率

被抽取的类别与决策类别一致的概率

观测样本为第  $i$  类并同被决策为第  $i$  类的概率

## 主要内容

### 1. 模型的学习能力与泛化能力

### 2. 交叉验证

### 3. (基于测试集的) 模型评价

#### 3.1 基于测试集的二分类模型评价

#### 3.2 基于测试集的多分类模型的评价

#### 3.3 基于测试集的真实函数预测模型评价

说明: 聚类模型的评价在后续聚类模型学习中给出。

## Regression metrics

See the Regression metrics section of the user guide for further details.

<code>metrics.explained_variance_score(y_true, y_pred)</code>	Explained variance regression score function. The max_error metric calculates the maximum residual error.
<code>metrics.max_error(y_true, y_pred)</code>	The max_error metric calculates the maximum residual error.
<code>metrics.mean_absolute_error(y_true, y_pred)</code>	Mean absolute error regression loss.
<code>metrics.mean_squared_error(y_true, y_pred)</code>	Mean squared error regression loss.
<code>metrics.mean_squared_log_error(y_true, y_pred)</code>	Mean squared logarithmic error regression loss.
<code>metrics.median_absolute_error(y_true, y_pred)</code>	Median absolute error regression loss.
<code>metrics.mean_absolute_percentage_error(y_true, y_pred)</code>	Mean absolute percentage error (MAPE) regression loss.
<code>metrics.r2_score(y_true, y_pred)</code>	$R^2$ (coefficient of determination) regression score function.
<code>metrics.mean_poisson_deviance(y_true, y_pred)</code>	Mean Poisson deviance regression loss.
<code>metrics.mean_gamma_deviance(y_true, y_pred)</code>	Mean Gamma deviance regression loss.
<code>metrics.mean_tweedie_deviance(y_true, y_pred)</code>	Mean Tweedie deviance regression loss.
<code>metrics.d2_tweedie_score(y_true, y_pred)</code>	$D^2$ regression score function, percentage of Tweedie deviance explained.
<code>metrics.mean_pinball_loss(y_true, y_pred)</code>	Pinball loss for quantile regression.

测试样本集  $D_{Test} = \{(x_i, y_i), i = 1, \dots, N\}$

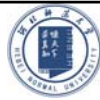
### 1. 均方误差 (Mean Squared Error, MSE)

`sklearn.metrics.mean_squared_error`

$$MSE(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N [y_i - \hat{y}_i]^2$$

### 2. 均方根误差 (Root Mean Squared Error, RMSE)

$$RMSE(y, \hat{y}) = \sqrt{\frac{1}{N} \sum_{i=1}^N [y_i - \hat{y}_i]^2} = \sqrt{MSE}$$



### 3. 平均绝对误差 (Mean Absolute Error, MAE)

`sklearn.metrics.mean_absolute_error`

$$MAE(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

### 4. 中值绝对误差 (Median Absolute Error)

最大绝对误差 (Max Absolute Error)

`sklearn.metrics.median_absolute_error`  
`sklearn.metrics.max_error`

$$\text{Median\_AE}(y, \hat{y}) = \text{median}_{i \in \{1, \dots, N\}} \{|y_i - \hat{y}_i|\}$$

$$\text{Max\_AE}(y, \hat{y}) = \max_{i \in \{1, \dots, N\}} \{|y_i - \hat{y}_i|\}$$



### 5. 均方根对数误差 (Root Mean Squared Logarithmic Error, RMSLE)

`sklearn.metrics.mean_squared_log_error`

$$\begin{aligned} RMSLE(y, \hat{y}) &= \sqrt{\frac{1}{N} \sum_{i=1}^N [\log(y_i + 1) - \log(\hat{y}_i + 1)]^2} \\ &= \sqrt{\frac{1}{N} \sum_{i=1}^N \left[ \log \frac{1 + y_i}{1 + \hat{y}_i} \right]^2} \end{aligned}$$

该评价指标的应用前提:

- 样本的“标签值”为重尾分布 (Heavy-tailed distribution)  
此时, 先进行对数运算再取RMSE, 可有效避免模型评价时, 少数极大值的影响。
- 此外, 应确保标签值与预测值为非负。



### 7. 对称平均绝对百分比误差

(Symmetric Mean Absolute Percentage Error, SMAPE)

$$SMAPE(y, \hat{y}) = 100\% \times \frac{1}{N} \sum_{i=1}^N 2 \left| \frac{y_i - \hat{y}_i}{|y_i| + |\hat{y}_i|} \right|$$

该指标的取值范围  $[0, +\infty)$

- 若 SMAPE = 0, 则预测模型为完美模型;
- 若 SMAPE > 100 %, 则预测模型为劣质模型。

注意: 当真值及预测值同为0时, 出现分母0除问题, 上述评价指标将不可用。

### 6. 平均绝对百分比误差 (Mean Absolute Percentage Error, MAPE)

`sklearn.metrics.mean_absolute_percentage_error`

$$MAPE(y, \hat{y}) = 100\% \times \frac{1}{N} \sum_{i=1}^N \frac{|y_i - \hat{y}_i|}{\max\{\epsilon, |y_i|\}}$$

为避免被0除, 引入小正数

该指标的取值范围  $[0, +\infty)$

- 若 MAPE = 0, 则预测模型为完美模型;
- 若 MAPE > 100 %, 则预测模型为劣质模型。

注意: 当真值存在取值0情况时, 出现分母0除问题, 上述评价指标将不可用。



### 8. 决定系数

(可决系数 coefficient of determination:  $R^2$ )

`sklearn.metrics.r2_score`

$$\begin{aligned} R^2(y, \hat{y}) &= 1 - \frac{MSE(y, \hat{y})}{Var(y)} = 1 - \frac{\frac{1}{N} \sum_{i=1}^N [y_i - \hat{y}_i]^2}{\frac{1}{N} \sum_{i=1}^N [y_i - \bar{y}]^2} \\ &= 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^N [y_i - \hat{y}_i]^2}{\sum_{i=1}^N [y_i - \bar{y}]^2} \end{aligned}$$

决定系数取值在  $(-\infty, 1]$  之间:

- 越接近于1, 说明模型的预测效果越好;
- 越接近于0, 说明模型的预测效果越差;
- 若取值为负值, 说明模型的效果非常差。



## 9. 解释方差 (Explained Variance)

[sklearn.metrics.explained\\_variance\\_score](#)

$$\text{explained\_variance\_score}(y, \hat{y}) = 1 - \frac{\text{Var}(y - \hat{y})}{\text{Var}(y)}$$

解释方差取值在 $(-\infty, 1]$ 之间:

- 越接近于1, 说明模型的预测效果越好;
- 越小于1, 说明模型的预测效果越差.



## 本章思考题(1)

1. 什么是模型的学习能力? 泛化能力?
2. 区分: 训练集、测试集、估计集、验证集。
3. 结合分类与回归问题, 给定已知答案的样本集, 如何基于K-fold CV或LOOCV进行模型的评价? **如何得到最终的模型? K值一般多大?**
4. 给定已知答案的训练样本集, 请结合分类与回归问题, 如何基于K-fold CV或LOOCV进行模型的选择? K值一般多大?
5. 区分hold out, bootstrapping, K-fold CV与LOOCV的使用场合。



## 本章思考题(2)

1. 对于**两类别(非对称类别/对称类别)**、**多类别**的分类问题, 如何分类模型对测试集的预测结果, 得到相应的混淆矩阵? **上述情况下的混淆矩阵有什么区别?** 如何计算混淆矩阵? 如何以图的方式可视化混淆矩阵?
2. 如何针对**上述三种情况**, 利用混淆矩阵生成若干评价指标?
3. 什么是ROC曲线? 如何利用有限规模测试样本绘制ROC曲线? 如何基于ROC曲线进行模型的评价与比较? 如何计算ROC曲线下面积(即: AUC值, 也称AUC\_ROC)?
4. 什么是P-R曲线? 如何利用有限规模测试样本绘制P-R曲线? 如何基于P-R曲线实现模型的定性或定量评价? 如何计算P-R曲线下面积(即: AUC\_PR)?
5. 哪种曲线更适合非对称类别分类模型的性能评价? 为什么?



## 本章思考题(3)

1. 理解并掌握回归模型的常用评价指标:
  - **决定系数、方差解释比**
  - **平均绝对误差、最大绝对误差、中值绝对误差**
  - **平均平方误差、均方根误差**
  - **均方根对数误差及应用场合**
  - **平均绝对百分比误差、对称平均绝对百分比误差**
  - ...

