

# 模式识别与机器学习

## 01 绪论 (第2部分)



大连理工大学 人工智能学院

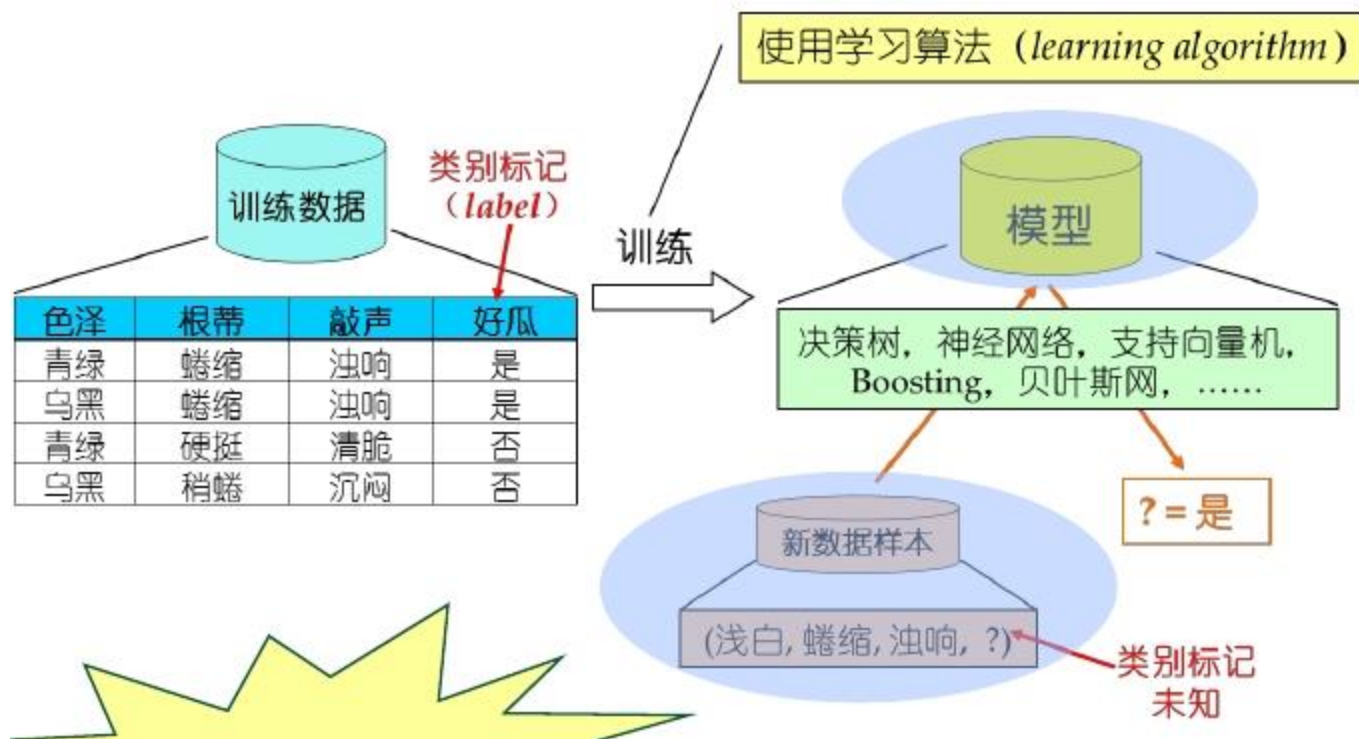
School of Artificial Intelligence, Dalian University of Technology



# 模型评估与选择

---

## Model Evaluation and Selection



泛化能力强!

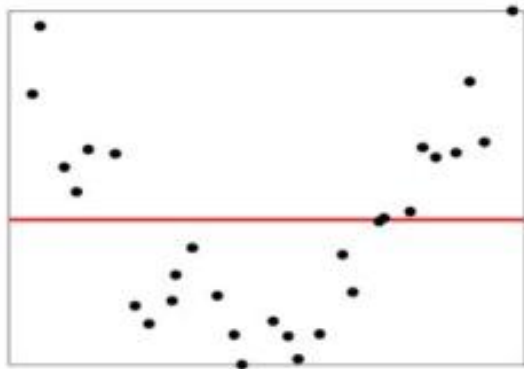
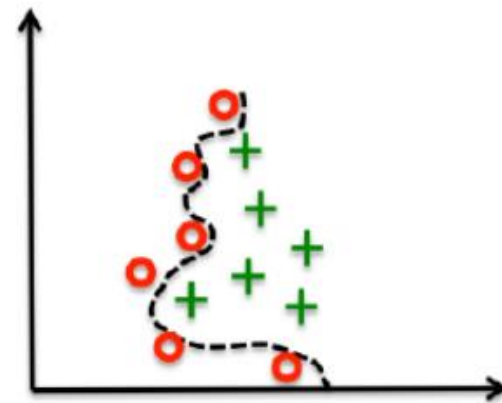
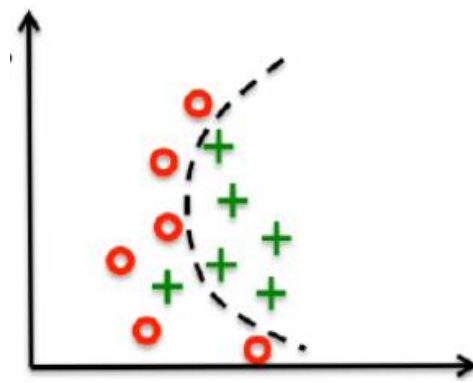
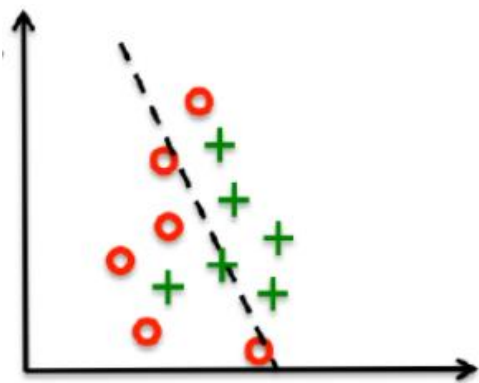
能很好地适用于未知样本  
例如, 错误率低、精度高

- 误差：样本真实输出与预测输出之间的差异
  - 训练(经验)误差：训练集上
  - 测试误差：测试集
  - 泛化误差：除训练集外所有样本
    - 泛化误差越小越好
    - 经验误差是否越小越好？

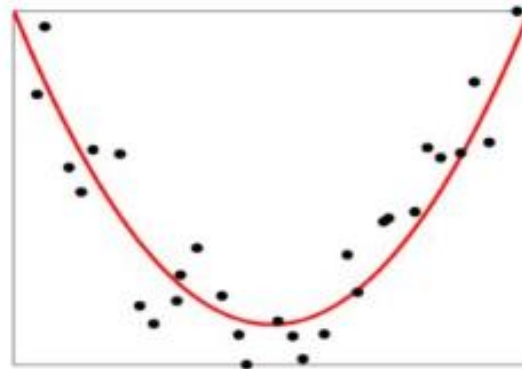
NO! 因为会出现“过拟合”(overfitting)

- 欠拟合：对训练样本的一般性质尚未学好
- 过拟合：学习器把训练样本学习的“太好”，将训练样本本身的特点当做所有样本的一般性质，导致泛化性能下降

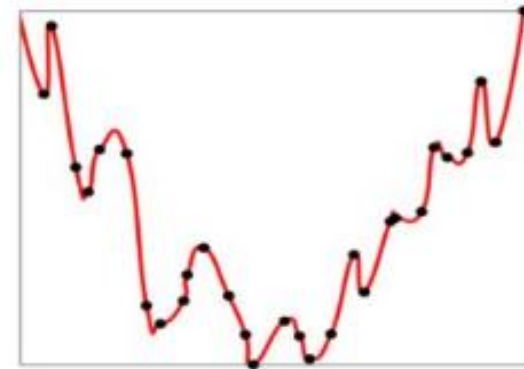




欠拟合



拟合



过拟合

三个关键问题:

□ 如何获得测试结果?

评估方法

□ 如何评估性能优劣?

性能度量

□ 如何判断实质差别?

比较检验

关键：怎么获得“测试集” (test set) ？

我们假设测试集是从样本真实分布中独立采样获得，  
将测试集上的“测试误差”作为泛化误差的近似，**测试集应该与训练集“互斥”**

常见方法：

- ❑ 留出法 (hold-out)
- ❑ **交叉验证法 (cross validation)**
- ❑ 自助法 (bootstrap)



通常将包含个 $m$ 样本的数据集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$  拆分成训练集  $S$  和测试集  $T$ :

- 留出法:
  - 直接将数据集划分为两个互斥集合
  - 训练/测试集划分要尽可能保持数据分布的一致性
  - 一般若干次随机划分、重复实验取平均值
  - 训练/测试样本比例通常为2:1~4:1

- $k$ -fold cross validation
- 将数据集分层采样划分为 $k$ 个大小相似的互斥子集，每次用 $k-1$ 个子集的并集作为训练集，余下的子集作为测试集，最终返回 $k$ 个测试结果的均值， $k$ 最常用的取值是10.

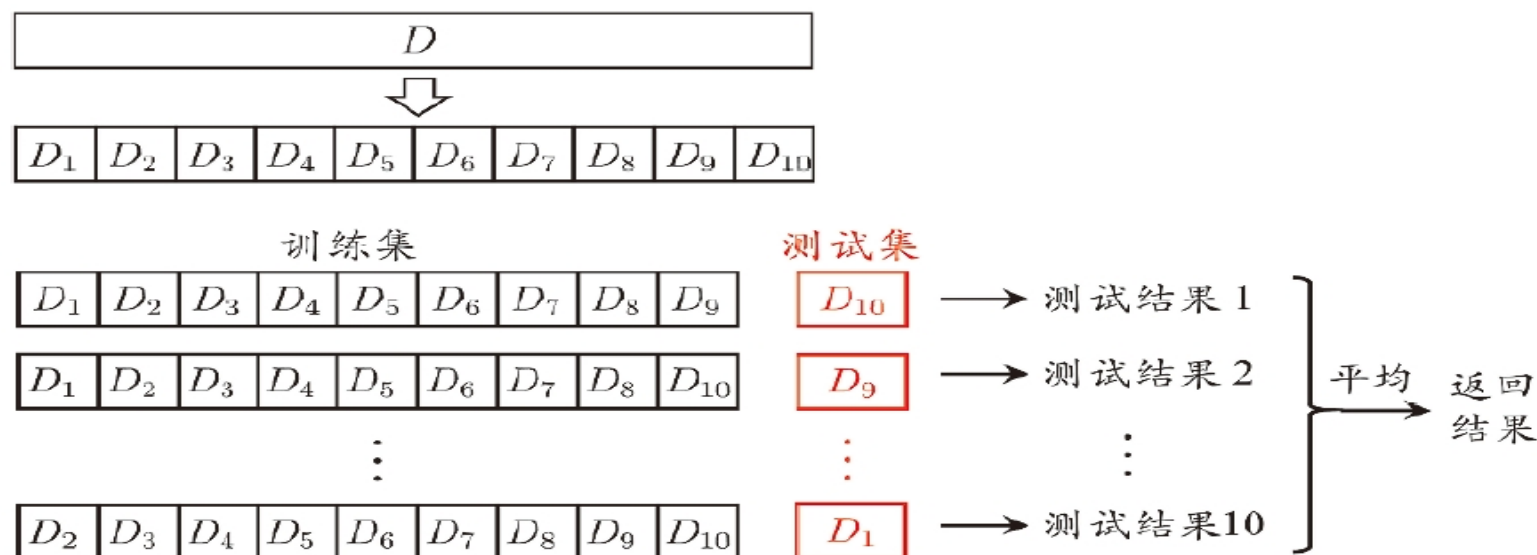


图 2.2 10 折交叉验证示意图

与留出法类似，将数据集 $D$ 划分为 $k$ 个子集同样存在多种划分方式，为了减小因样本划分不同而引入的差别， $k$ 折交叉验证通常随机使用不同的划分重复 $p$ 次，最终的评估结果是这 $p$ 次 $k$ 折交叉验证结果的均值，例如常见的“10次10折交叉验证”

假设数据集 $D$ 包含 $m$ 个样本，若令  $k = m$ ，则得到留一法：

- 不受随机样本划分方式的影响 (leave-one-out, LOO)
- 结果往往比较准确
- 当数据集比较大时，计算开销难以忍受

## □ 自助法 (Bootstrapping) :

以自助采样法（**即有放回的采样或重复采样**）为基础，对数据集  $D$  有放回采样  $m$  次得到训练集  $D'$ ， $D \setminus D'$  用做测试集。

- 实际模型与预期模型都使用  $m$  个训练样本
- 约有1/3的样本没在训练集中出现 (如果 $m$ 趋于无穷，应为 $e^{-1}$ .**why?**)
- 从初始数据集中产生多个不同的训练集，对集成学习有很大的好处
- 自助法在数据集较小、难以有效划分训练/测试集时很有用；
- 由于改变了数据集分布可能引入估计偏差，所以在数据量足够时，留出法和交叉验证法更常用。

算法的参数：一般由人工设定，亦称“超参数”

模型的参数：一般由学习确定

调参过程相似：先产生若干模型，然后基于某种评估方法进行选择

参数调得好不好对性能往往对最终性能有关键影响

区别：训练集 vs. 测试集 vs. 验证集 (validation set)

算法参数选定后，要用“训练集+验证集”重新训练最终模型

三个关键问题:

□ 如何获得测试结果?

评估方法

□ 如何评估性能优劣?

性能度量

□ 如何判断实质差别?

比较检验

- ◆ 在预测任务中，给定样例集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$
  - ◆ 评估学习器的性能  $f$  也即把预测结果  $f(\mathbf{x})$  和真实标记比较.
- 
- 性能度量(performance measure)是衡量模型泛化能力的评价标准，反映了任务需求
  - 使用不同的性能度量往往会导致不同的评判结果

对于回归(regression) 任务常用均方误差：

$$E(f; D) = \frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2$$

什么样的模型是“好”的，不仅取决于算法和数据，还取决于任务需求

对于分类任务, 错误率和精度是最常用的两种性能度量:

- 错误率: 分错样本占样本总数的比例
- 精度: 分对样本占样本总数的比率

□ 错误率:

$$E(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) \neq y_i)$$

□ 精度:

$$\begin{aligned} \text{acc}(f; D) &= \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(\mathbf{x}_i) = y_i) \\ &= 1 - E(f; D) . \end{aligned}$$

什么样的模型是“好”的, 不仅取决于算法和数据, 还取决于任务需求



信息检索、Web搜索等场景中经常需要衡量正例被预测出来的比率或者预测出来的正例中正确的比率，此时查准率和查全率比错误率和精度更适合。

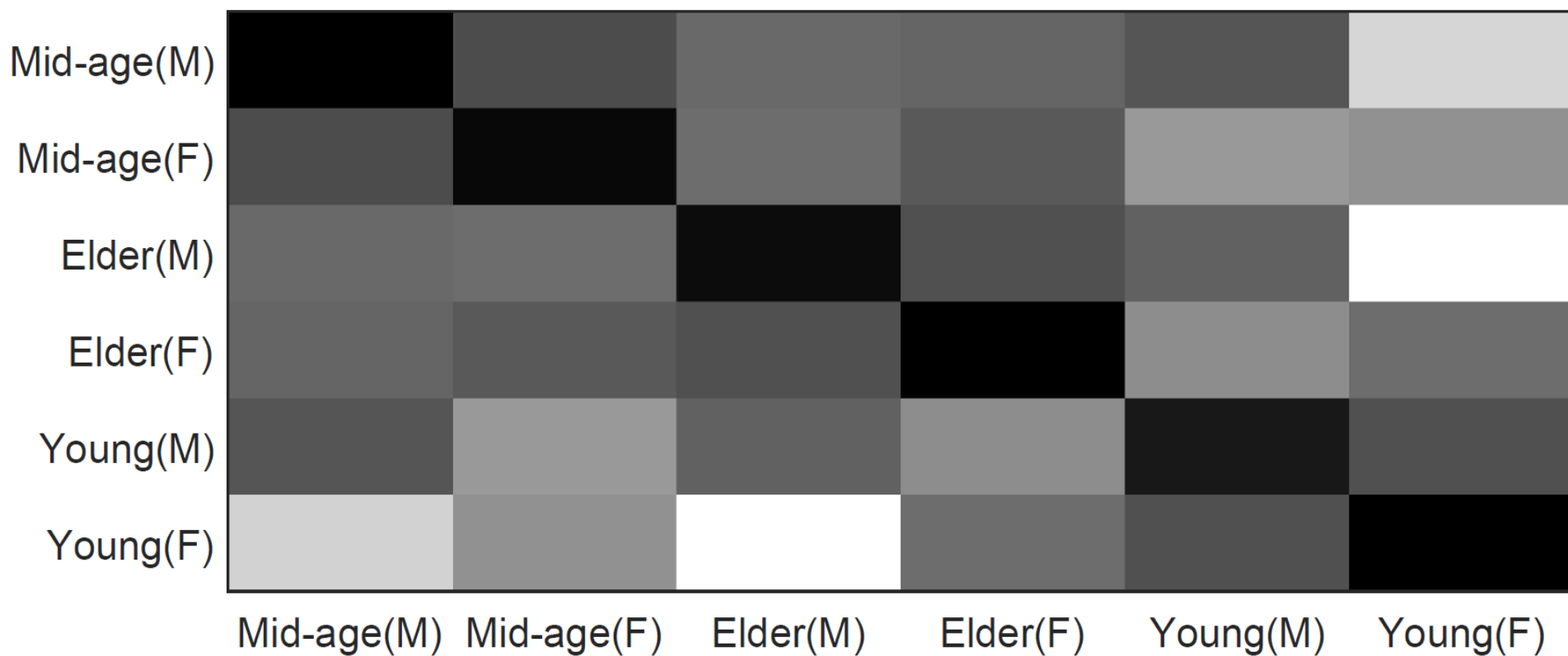
统计真实标记和预测结果的组合可以得到“混淆矩阵”

分类结果混淆矩阵

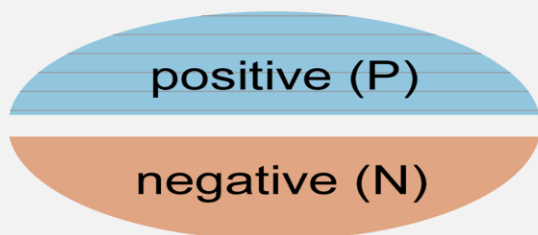
| 真实情况 | 预测结果       |            |
|------|------------|------------|
|      | 正例         | 反例         |
| 正例   | $TP$ (真正例) | $FN$ (假反例) |
| 反例   | $FP$ (假正例) | $TN$ (真反例) |

$$\text{查准率 } P = \frac{TP}{TP + FP}$$

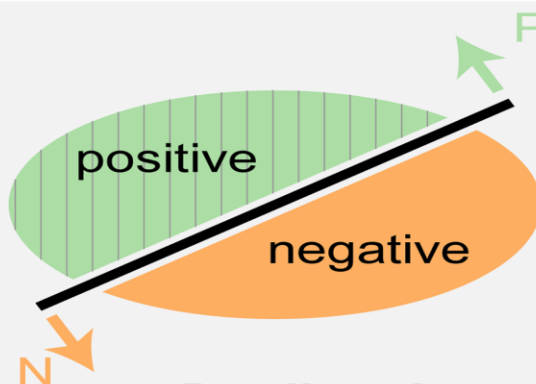
$$\text{查全率 } R = \frac{TP}{TP + FN}$$



**A**



**Actual**

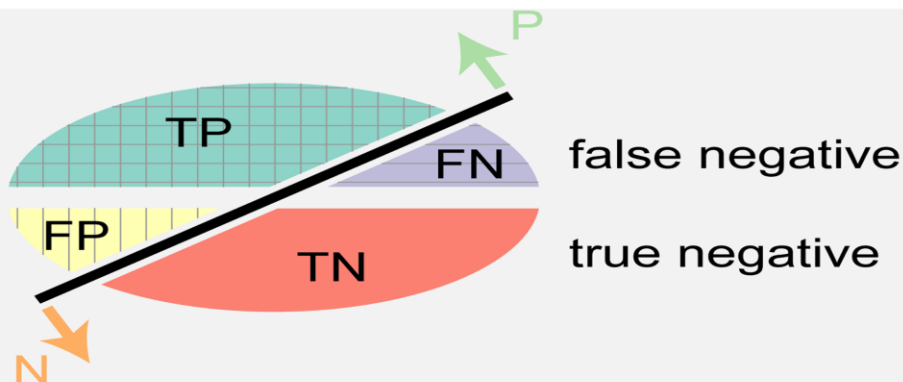


**Predicted**

**B**

true positive

false positive



**Four outcomes**

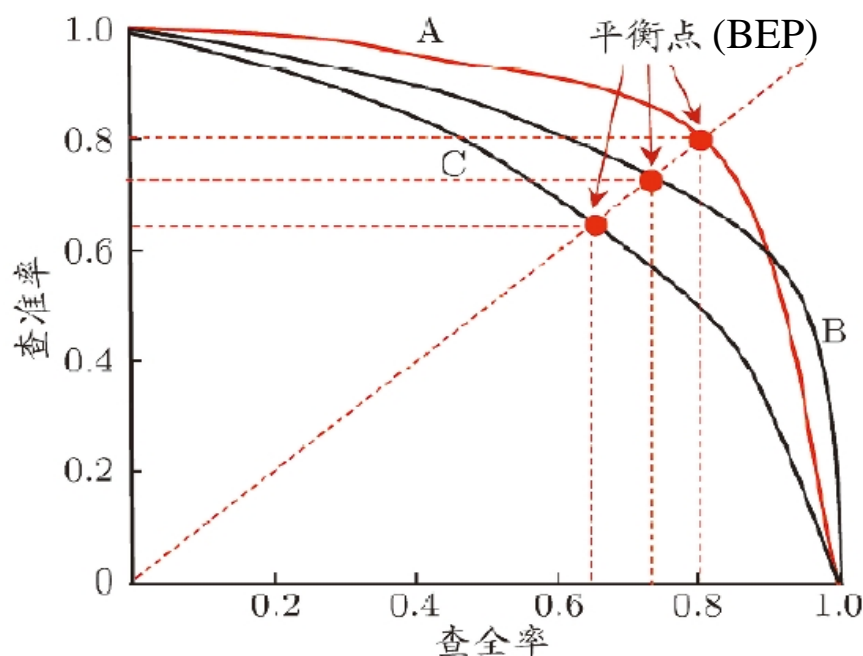
□ 查准率  
(precision) :

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

□ 查全率  
(Recall) :

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

- 根据学习器的预测结果按正例可能性大小对样例进行排序，并逐个把样本作为正例进行预测，则可以得到查准率-查全率曲线，简称“**P-R曲线**”。
- 平衡点**是曲线上“查准率=查全率”时的取值，可用来用于度量P-R曲线有交叉的分类器性能高低。



## P-R曲线:

- 学习器 A 优于 学习器 C
- 学习器 B 优于 学习器 C
- 学习器 A ?? 学习器 B

## BEP:

- 学习器 A 优于 学习器 B
- 学习器 A 优于 学习器 C
- 学习器 B 优于 学习器 C

比 BEP 更常用的 F1 度量：**F1度量的由来是加权调和平均**

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

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

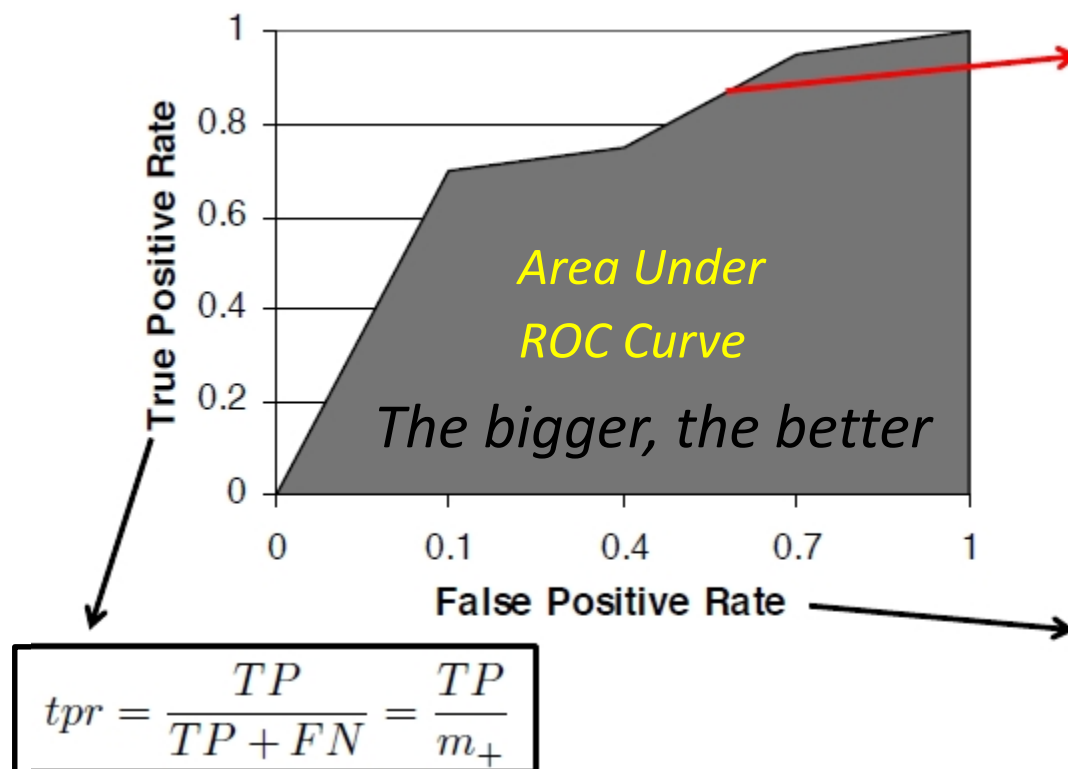
$\beta = 1$ : 标准F1

$\beta > 1$ : 偏重查全率(逃犯信息检索)

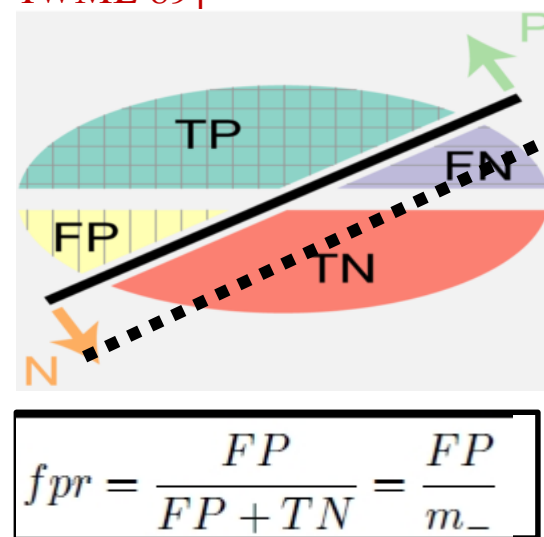
$\beta < 1$ : 偏重查准率(商品推荐系统)

调和平均数主要特点：易受极端值的影响，且受极小值的影响比受极大值的影响更大。

## AUC: Area Under the ROC Curve



ROC (Receiver Operating Characteristic) Curve [Green & Swets, Book 66; Spackman, IWML'89]



➤ 更多性能指标讨论: <http://charleshm.github.io/2016/03/Model-Performance/>

- 关于性能比较：
  - 测试性能并不等于泛化性能
  - 测试性能随着测试集的变化而变化
  - 很多机器学习算法本身有一定的随机性

**直接选取相应评估方法在相应度量下比大小的方法不可取！**

假设检验为学习器性能比较提供了重要依据，基于其结果我们可以推断出若在测试集上观察到学习器A比B好，则A的泛化性能是否在统计意义上优于B，以及这个结论的把握有多大。

参考《机器学习》周志华 第2.4节

“误差”包含了哪些因素？

换言之，从机器学习的角度看，

“误差”从何而来？



对回归任务，泛化误差可通过“偏差-方差分解”拆解为：

$$E(f; D) = \underbrace{bias^2(x)}_{\text{red}} + \underbrace{var(x)}_{\text{blue}} + \underbrace{\varepsilon^2}_{\text{green}}$$

期望输出与真实  
输出的差别

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

同样大小的训练集  
的变动，所导致的  
性能变化

$$var(x) = \mathbb{E}_D \left[ (f(x; D) - \bar{f}(x))^2 \right]$$

训练样本的标记与  
真实标记有区别

表达了当前任务上任何学习算法  
所能达到的期望泛化误差下界

$$\varepsilon^2 = \mathbb{E}_D \left[ (y_D - y)^2 \right]$$

泛化性能是由学习算法的能力、算法在不同数据集的稳定性以及学习任务本身的难度共同决定

一般而言，偏差与方差存在冲突：

- ❑ 训练不足时，学习器拟合能力不强，偏差主导
- ❑ 随着训练程度加深，学习器拟合能力逐渐增强，方差逐渐主导
- ❑ 训练充足后，学习器的拟合能力很强，方差主导

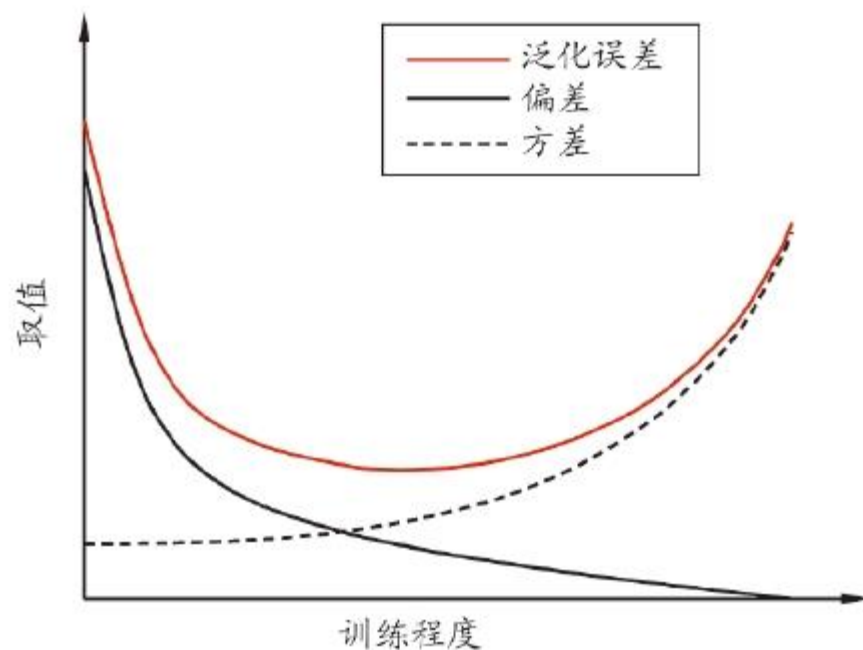


图 2.9 泛化误差与偏差、方差的关系示意图

➤ 更多偏差-方差的讨论：<https://zhuanlan.zhihu.com/p/38853908>



## 三个关键问题:

### □ 如何获得测试结果?

评估方法: 留出法、交叉验证法、自助法

### □ 如何评估性能优劣?

性能度量: 均方误差、错误率、精度、查准率、查全率、ROC曲线等

### □ 如何判断实质差别?

比较检验



# PART FOUR

参考资源

Resource

## ➤ 在线课程:

吴恩达 (Andrew Ng): 机器学习

<https://www.bilibili.com/video/av50747658?from=search&seid=9857444623866219885>

李宏毅 (Hung-Yi Lee): 机器学习

<https://www.bilibili.com/video/av10590361?from=search&seid=10930215693466580647>

## ➤ 公众号:

- VALSE
- 机器之心
- 极市平台
- 人工智能前沿讲习班
- 深度学习大讲堂
- SIGAI
- ...