

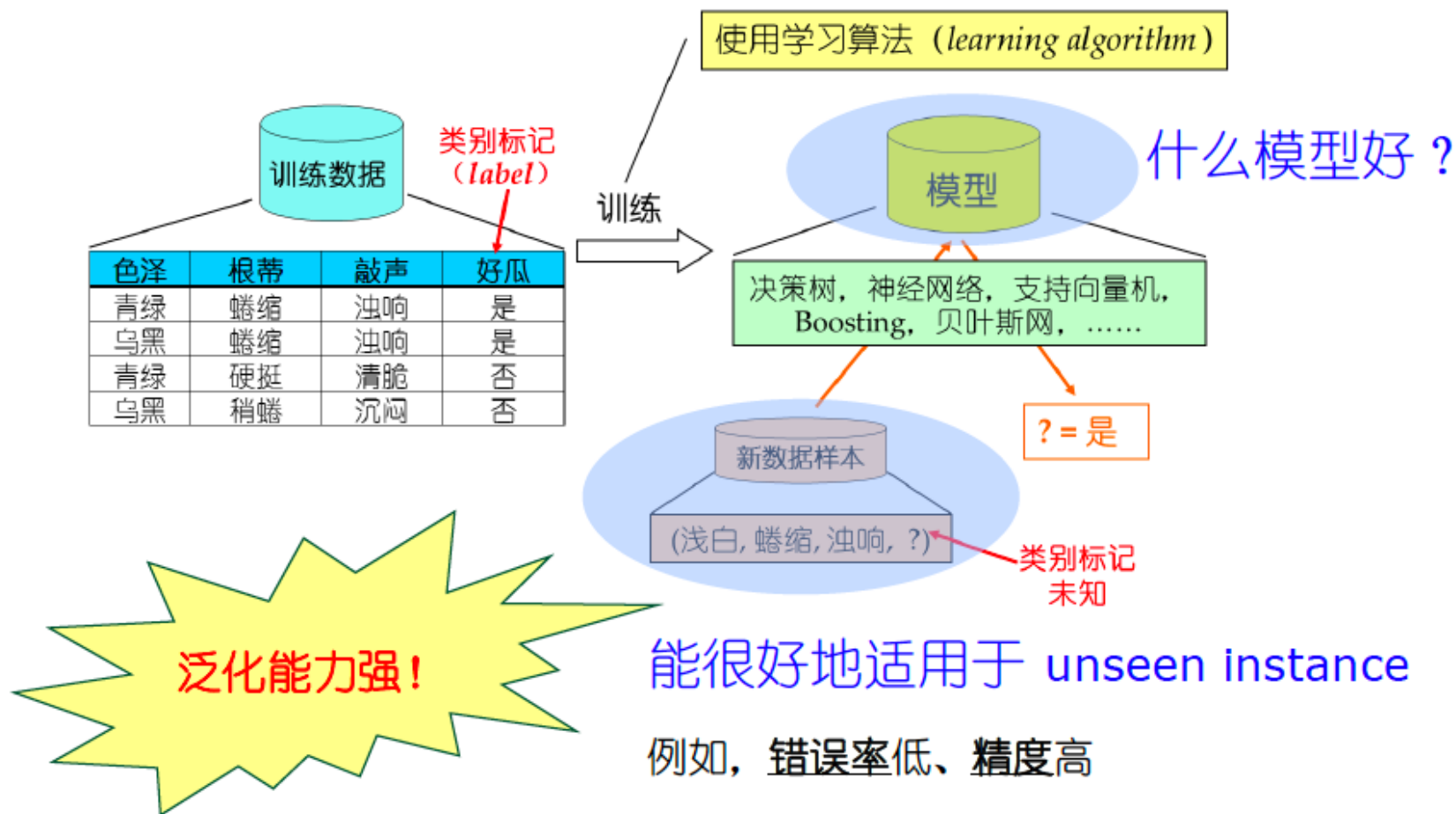


# 第二讲 模型评估和选择

## 高级机器学习



# 典型的机器学习过程



然而, 我们手上没有 unseen instance, .....

# 泛化误差 vs 经验误差

---

泛化误差：在“未来”样本上的误差

经验误差：在训练集上的误差，亦称“训练误差”

- 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)

# 过拟合 ( overfitting ) vs 欠拟合 ( underfitting )

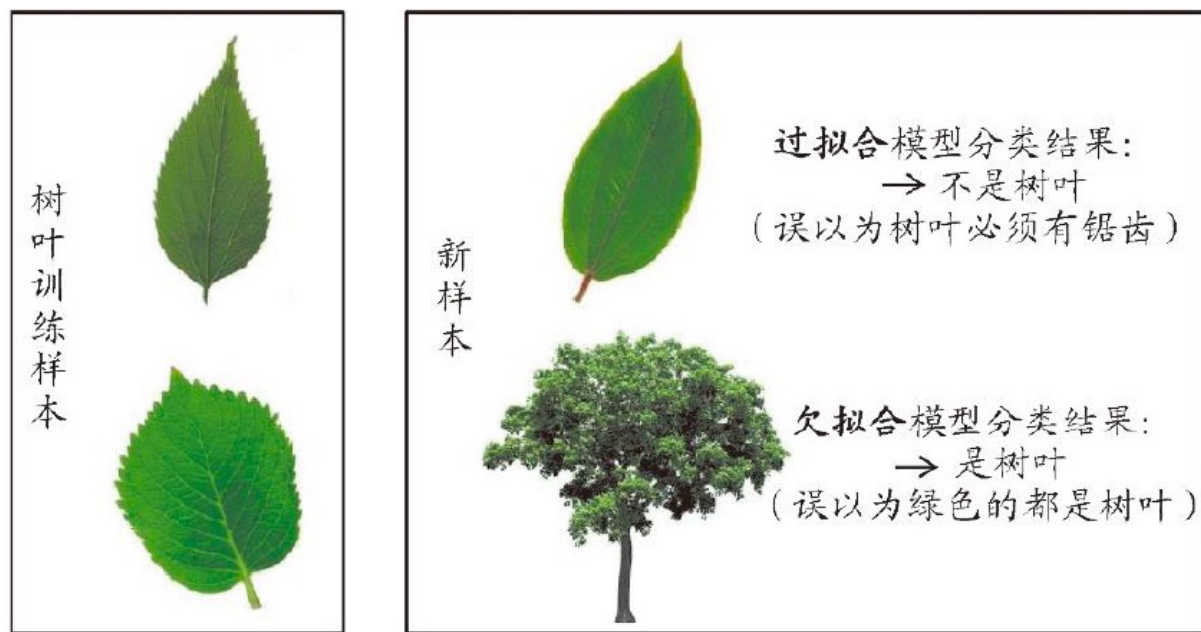


图 2.1 过拟合、欠拟合的直观类比

1. 在不存在“严重”过拟合或欠拟合现象时，经验误差和泛化误差可以进行相互界定；
2. 绝对地消除过拟合或欠拟合不可能；
3. 余下讨论中，假定算法没有严重过拟合现象。

# 模型选择 - 经验误差

---

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

# 评估方法

---

可是, 我们只有一个包含  $m$  个样例的数据集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ , 既要训练, 又要测试, 怎样才能做到呢? 答案是: 通过对  $D$  进行适当的处理, 从中产生出训练集  $S$  和测试集  $T$ . 下面介绍几种常见的做法.

- 留出法 (hold-out)
- 交叉验证法
- 留一法

越往下或更准确, 但  
更复杂, 开销越高

# 评估方法

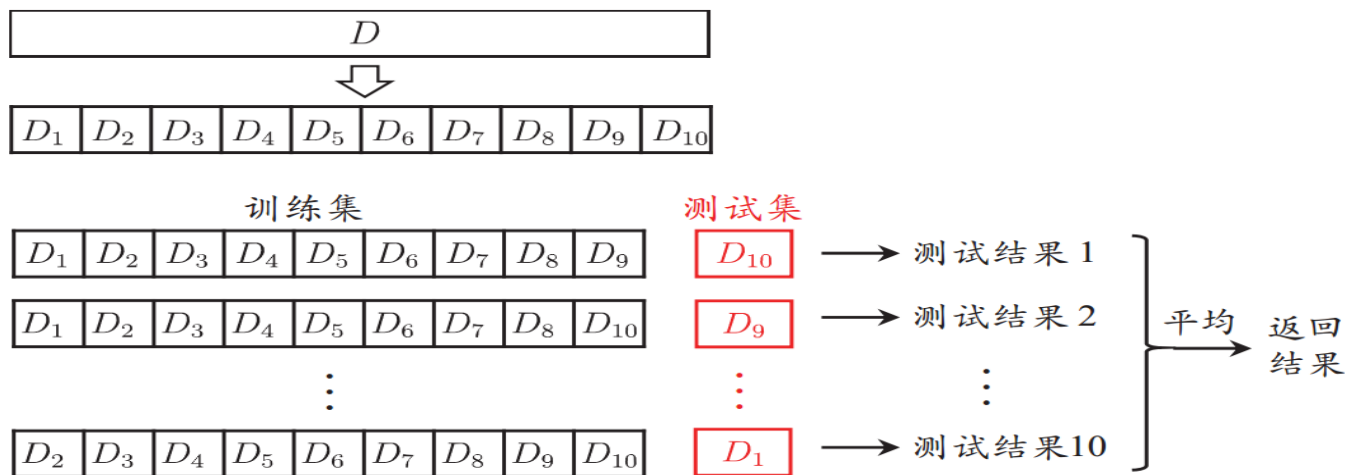
---

- 留出法 (hold-out):
  - 直接将数据集划分为两个互斥集合——训练和测试集
  - 训练/测试集划分要尽可能保持数据分布的一致性
  - 分层采样 ( stratified sampling ) : 保持类别比例一致
  - 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 , 效果还不错

# 评估方法

- 交叉验证法（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。



10折交叉验证示意图

# 评估方法

---

**P次k折交叉验证**：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(LOO)**：

- 不受随机样本划分方式的影响
- 结果往往比较准确
- 当数据集比较大时，计算开销难以忍受，实际中不采用

# 模型选择

---

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

# 性能度量

---

性能度量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

回归任务最常用的性能度量是“均方误差”：

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

对于分类任务，错误率和精度是最常用的两种性能度量：

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

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

# 性能度量

信息检索、Web搜索等场景中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)

查准率  $P = \frac{TP}{TP + FP}$

查全率  $R = \frac{TP}{TP + FN}$

# 性能度量

---

比P-R曲线平衡点更常用的是**F1**度量:

$$F1 = \frac{2 \times P \times R}{P + R} = \frac{2 \times TP}{\text{样例总数} + TP - TN}$$

比F1更一般的形式  $F_\beta$

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

$\beta = 1$ : 标准F1

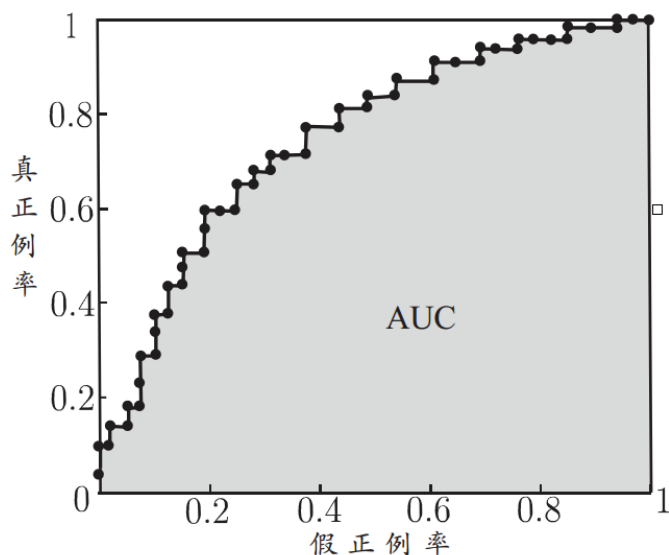
$\beta > 1$ : 偏重查全率(逃犯信息检索)

$\beta < 1$ : 偏重查准率(商品推荐系统)

# 性能度量

**ROC曲线**: 遍历真比例率和假正例率, 形成ROC曲线

**AUC值**: ROC曲线下面积大小



基于有限样例绘制的 ROC 曲线  
与 AUC

假设ROC曲线由 $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ 的点按序连接而形成, 则: AUC可估算为:

$$\text{AUC} = \frac{1}{2} \sum_{i=1}^{m-1} (x_{i+1} - x_i) \cdot (y_i + y_{i+1})$$

**AUC衡量样本预测的排序质量**

$$\text{AUC} = 1 - \ell_{\text{rank}} .$$

# 模型选择

---

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
-

# 性能评估

---

- 性能比较的若干事实：
  - 测试性能不等于泛化性能
  - 测试性能随着测试集的变化而变化
  - 机器学习算法本身有一定的随机性

模型A 5次评估性能：91%，92%，93%，92%，92%

模型B 5次评估性能：86%，100%，99%，89%，91%

模型C 5次评估性能：85%，98%，98%，88%，90%

模型A、B、C平均性能分别是92%，93%，92.2%

你会选用哪个模型？为什么？

# 性能评估

---

- 性能比较的若干事实：
  - 测试性能不等于泛化性能
  - 测试性能随着测试集的变化而变化
  - 机器学习算法本身有一定的随机性

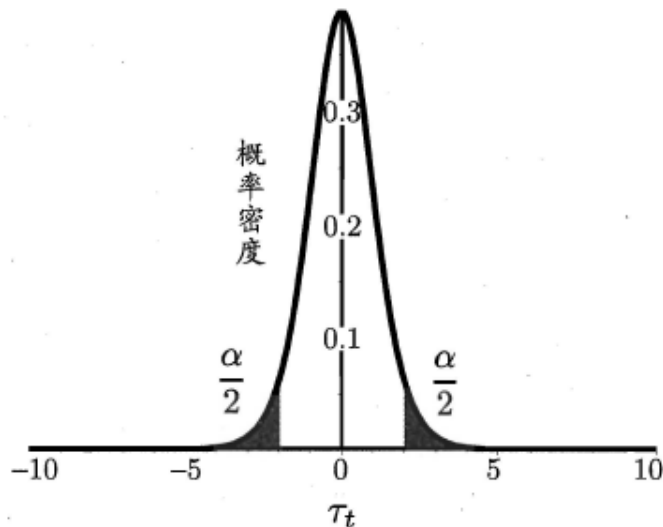
直接根据性能大小判断模型优劣可能会有错误，需要缓解上述不确定性。

**假设检验：**为学习器的性能比较提供了数学理论依据。利用假设检验可推断学习器A的泛化性能是否在统计意义上优于B，以及该结论的把握有多大

# T-检验

实施多次留出法或者交叉验证法进行多次训练/测试，得到学习器A或者B多次结果。性能比较的常用假设检验可采用“T-检验”

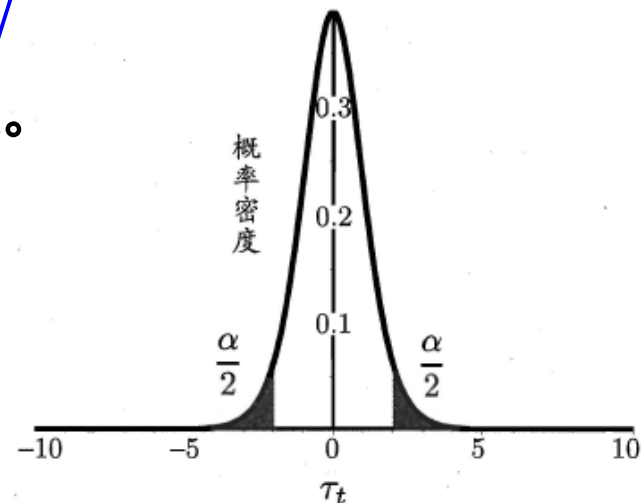
假定得到了k个测试错误率， $\hat{\epsilon}_1, \hat{\epsilon}_2, \dots, \hat{\epsilon}_k$ ，假设 $\epsilon = \epsilon_0$ 对于显著度 $\alpha$ ，若 $[t_{-\alpha/2}, t_{\alpha/2}]$ 位于临界范围 $|\mu - \epsilon_0|$ 内，则假设不能被拒绝，即可认为泛化错误率 $\epsilon = \epsilon_0$ ，其置信度为 $1 - \alpha$ 。



# T-检验

对两个学习器A和B,若k折交叉验证得到的测试错误率分别为 $\epsilon_1^A, \dots, \epsilon_k^A$ 和 $\epsilon_1^B, \dots, \epsilon_k^B$ ,可用k折交叉验证“成对t检验”进行比较检验。若两个学习器的性能相同,则他们使用相同的训练/测试集得到的测试错误率应相同,即 $\epsilon_i^A = \epsilon_i^B$ 。

要点：**成对T-检验必须要在相同的训练/测试集进行**，否则不具备数学统计意义。



# 理论问题：泛化风险和经验风险的差距

---

“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？

# 偏差-方差分解

对回归任务，泛化误差可通过“偏差-方差分解”拆解为：

$$E(f; D) = \underbrace{bias^2(\mathbf{x})}_{\text{期望输出与真实输出的差别}} + \underbrace{var(\mathbf{x})}_{\text{同样大小的训练集的变动, 所导致的性能变化}} + \underbrace{\varepsilon^2}_{\text{训练样本的标记与真实标记有区别}}$$

期望输出与真实输出的差别

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

同样大小的训练集的变动, 所导致的性能变化

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

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、数据的充分性以及学习任务本身的难度共同决定

# 小结

---

- 经验误差与过拟合
  - 过拟合, 欠拟合
- 评估方法
  - 留出法, 交叉验证
- 性能度量
  - 均方误差, 精度, 查准率/查全率/F1, AUC
- 比较检验
  - T-检验
- 偏差与方差
  - 了解基本原理