

解决的是近似正确的问题，研究的是学习算法：从数据中获得模型的算法

优化+学习+搜索找出最优假设

用 $1-\epsilon$ 表示趋近于1的量

有时候不需要很精准，比如预测考试分数

一般原则：奥卡姆剃刀：不要画蛇添足：若无必要，勿增实体

基本术语

- 训练：使用学习算法（决策树，神经网络.....），得到模型
- 预测：新的数据样本无标记，只包含描述特征，模型应用到新数据上得到标记
- 测试数据：用于评估模型的数据
- 独立同分布：测试和训练数据类似且独立

重要的假设

- 泛化：训练的模型能推广到新的问题上

我们更关注在未知的地方做的更好，测试集准确率高才重要

- 特征信息不充分：收集不到某项数据（雾天）
- 样本信息不充分：训练数据不够多（新用户怎么推荐）

- 横向：样例=instance+label示例+标记

样本的真相就是标记

- 纵向：属性，特征

假设是在属性的基础上建立的特征到标记的关系

- 监督学习：有标记，有特征描述
- 无监督学习：只有特征输入，无标记

聚类，密度估计，降维

- 先验概率：各类别出现概率
- 样例x的后验概率：x属于第i类的可能性

- 假设空间：每个点都对应了一条规则

机器学习是在假设空间中对假设搜索的过程

- 版本空间：与训练集一致的假设
- 归纳偏好：三个模型对训练集都100%，怎么选呢？--

任何有效的ML模型一定有偏好

NFL定理

如果某些问题上A更好，那一定有一种数据使B更好

把公式按照和谁有关重新排序

结果：所有算法一样好，不存在全局最优。方法只有和应用场景绑定才有好坏

真实目标函数 \mathcal{L}_a 的“训练集外误差”，即 \mathcal{L}_a 在训练集之外的所有样本上的误差为

$$E_{ote}(\mathcal{L}_a|X, f) = \sum_h \sum_{x \in \mathcal{X}-X} P(x) \mathbb{I}(h(x) \neq f(x)) P(h|X, \mathcal{L}_a), \quad (1.1)$$

其中 $\mathbb{I}(\cdot)$ 是指示函数, 若 \cdot 为真则取值 1, 否则取值 0.

考虑二分类问题, 且真实目标函数可以是任何函数 $\mathcal{X} \mapsto \{0, 1\}$, 函数空间为 $\{0, 1\}^{|\mathcal{X}|}$. 对所有可能的 f 按均匀分布对误差求和, 有

$$\begin{aligned} \sum_f E_{ote}(\mathcal{L}_a|X, f) &= \sum_f \sum_h \sum_{x \in \mathcal{X}-X} P(x) \mathbb{I}(h(x) \neq f(x)) P(h|X, \mathcal{L}_a) \\ &= \sum_{x \in \mathcal{X}-X} P(x) \sum_h P(h|X, \mathcal{L}_a) \sum_f \mathbb{I}(h(x) \neq f(x)) \\ &= \sum_{x \in \mathcal{X}-X} P(x) \sum_h P(h|X, \mathcal{L}_a) \frac{1}{2} 2^{|\mathcal{X}|} \\ &= \frac{1}{2} 2^{|\mathcal{X}|} \sum_{x \in \mathcal{X}-X} P(x) \sum_h P(h|X, \mathcal{L}_a) \end{aligned}$$

$$= 2^{|\mathcal{X}|-1} \sum_{x \in \mathcal{X}-X} P(x) \cdot 1. \quad (1.2)$$

式(1.2)显示出, 总误差竟然与学习算法无关! 对于任意两个学习算法 \mathcal{L}_a 和 \mathcal{L}_b , 我们都有

$$\sum_f E_{ote}(\mathcal{L}_a|X, f) = \sum_f E_{ote}(\mathcal{L}_b|X, f), \quad (1.3)$$

也就是说, 无论学习算法 \mathcal{L}_a 多聪明, 学习算法 \mathcal{L}_b 多笨拙, 它们的期望性能

化简恒等式

$$\sum_f \mathbb{I}(h(x) \neq f(x)) = 2^{|\mathcal{X}|-1}$$

$$\sum_h P(h|X, \mathcal{L}_a) = 1$$

选什么模型?

怎么配参数?

泛化误差: 未来样本上

经验误差: 训练样本上

过拟合：经验误差小但泛化误差大

所以训练到一定程度就够了，再继续就会学到不必要特征
不同算法有不同手段处理过拟合

- 1 评估方法：获得测试结果
- 2 性能度量：评估性能优劣，衡量泛化能力
- 3 假设检验：判断实质差别

一、评估方法

(1) 留出法：一刀切，互斥

- 可分层采样：保持数据分布一致性

如保持训练集和测试集中的正反例比例相同

- 单次结果不可靠，需要多次随机划分

每次一个训练/测试集，得到一个评估，所有评估结果取平均

(2) p次k折交叉验证法

常用 $k=10$

一共 pk 次训练

- p 次随机分组，全部数据分成 k 个互斥部分，每次选一个作为测试

尽量保证 k 组的数据分布一致性，即分层采样

- 返回 k 次训练测试后评估的均值
- 一共仅 k 个样本时：留一法：结果准确

最大化 k ， $k=m$ 时是留一法LOO，开销大但效果好

评价

- k 不是越大越好，开销
- 但大 k 会有更小的偏差
- 选择 k 时要最小化数据集之间的方差

5. 根据上述计算结果, 我们可以发现,

$$\mathbb{E}[\bar{x}_m^*] = \mathbb{E}[\bar{x}_m] = \mu \quad (50)$$

$$\text{var}[\bar{x}_m^*] = \frac{\sigma^2}{m} \left[2 - \frac{1}{m} \right] = \left(2 - \frac{1}{m} \right) \text{var}[\bar{x}_m] \approx 2 \text{var}[\bar{x}_m] \quad (51)$$

另外我们如果采用 k 折交叉验证法的方式采样, 类似地我们有,

$$\mathbb{E}[\bar{x}_m'] = \mathbb{E}[\bar{x}_m] = \mu \quad (52)$$

$$\text{var}[\bar{x}_m'] \approx \frac{k}{k-1} \text{var}[\bar{x}_m] \quad (53)$$

综上所述, 虽然通过自助采样法得到的样本均值仍然是总体均值的无偏估计, 但是其方差变为原来的接近两倍, 而这相当于使用 2 折交叉验证采样的效果, 所以一般来说, 自助法采样对数据分布的改变大于交叉验证法.

(3) 自助法

用于数据集小, 难以划分时

- 缺点: 改变数据集分布, 引起了估计偏差 (方差变为原来的近两倍)
- m次独立重复抽取
- 始终不被采到概率:

$$\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m} \right)^m \mapsto \frac{1}{e} \approx 0.368,$$

(4) 调参

算法参数 (超参数): 人工设定

算法参数设定后, 要用训练集+验证集重新训练最终模型

模型参数: 学习确定

- 范围+步长

计算开销+性能的折中方案

二、性能度量

□ 回归(regression) 任务常用均方误差:

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

1、精度&错误率

- 错误率

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

- 精度

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

2、查准率&查全率P-R图像

- FN：假的反例，应该是真的，预测成了假的
- 查准率R：所有标为正中的实际正例有多少
- 查全率R：所有的实际正例中标为正有多少
- 真正例率TPR=R
- 假正例率FPR：所有的实际反例中标为正有多少

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 与查全率 R 分别定义为

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

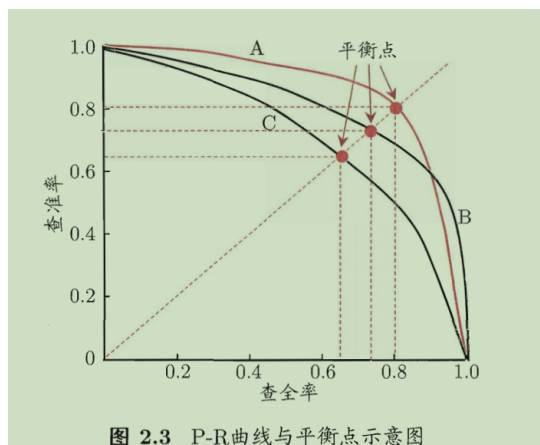
$$R = \frac{TP}{TP + FN} .$$

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

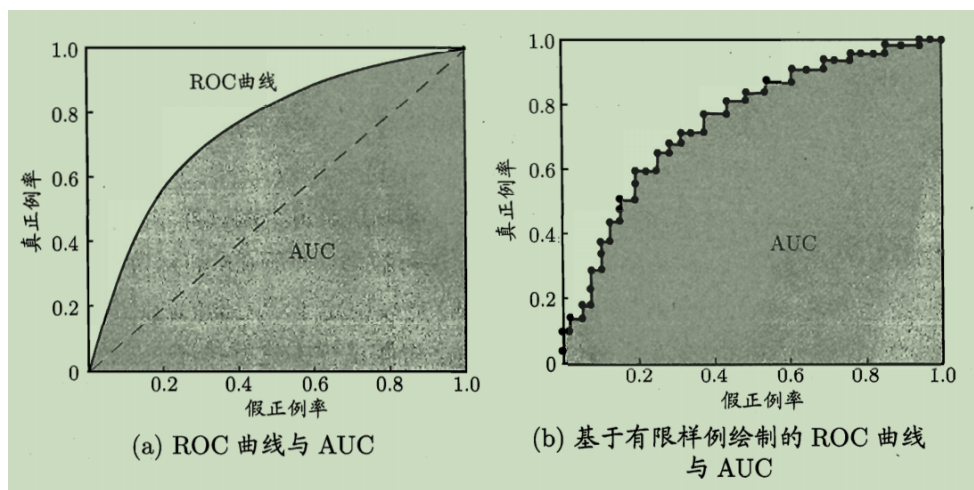
$$\text{FPR} = \frac{FP}{TN + FP} .$$

- 平衡点 (BEP) : $P=R$

解决P-R图交叉时不易选择的问题



3、ROC&AUC



- AUC: 曲线下面积

相当于一堆梯形的面积和

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

4、F1度量&Fβ度量

更常用，可调偏好

β越大，查全率R越重要

F1 度量：

$$F1 = \frac{2 \times P \times R}{P + R} \quad \frac{1}{F1} = \frac{1}{2} \cdot \left(\frac{1}{P} + \frac{1}{R} \right)$$
$$= \frac{2 \times TP}{\text{样例总数} + TP - TN}$$

若对查准率/查全率有不同偏好：

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

$\beta > 1$ 时查全率有更大影响; $\beta < 1$ 时查准率有更大影响

5、宏xx&微xx

多次测试

- 宏和微是不同的取平均方式

若能得到多个混淆矩阵：

(例如多次训练/测试的结果，多分类的两两混淆矩阵)

宏(macro-)查准率、查全率、F1

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

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

$$\text{macro-}F1 = \frac{2 \times \text{macro-}P \times \text{macro-}R}{\text{macro-}P + \text{macro-}R}.$$

微(micro-)查准率、查全率、F1

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

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

$$\text{micro-}F1 = \frac{2 \times \text{micro-}P \times \text{micro-}R}{\text{micro-}P + \text{micro-}R}.$$

6、错误非均等代价

将第*i*类样本预测为第*j*类样本的代价

- $m+$ 个正例， $m-$ 个反例
- $D+$ ， $D-$ 正反例集合

代价敏感(cost-sensitive)错误率：

$$E(f; D; cost) = \frac{1}{m} \left(\sum_{\mathbf{x}_i \in D^+} \mathbb{I}(f(\mathbf{x}_i) \neq y_i) \times cost_{01} \right. \\ \left. + \sum_{\mathbf{x}_i \in D^-} \mathbb{I}(f(\mathbf{x}_i) \neq y_i) \times cost_{10} \right)$$

三、比较检验

评估结果不能用于判断优劣：

- 1 测试性能 \neq 泛化性能
- 2 测试性能随测试集变化
- 3 很多ml算法本身有一定随机性

通过比较才能说某个方法好，直接说太抽象

(1) 交叉验证t检验

(2) McNemar检验

算法 B	算法 A	
	正确	错误
正确	e_{00}	e_{01}
错误	e_{10}	e_{11}

若我们做的假设是两学习器性能相同，则应有 $e_{01} = e_{10}$ ，那么变量 $|e_{01} - e_{10}|$ 应当服从正态分布，且均值为 1，方差为 $e_{01} + e_{10}$ 。因此变量

$$\tau_{\chi^2} = \frac{(|e_{01} - e_{10}| - 1)^2}{e_{01} + e_{10}} \quad (2.33)$$

服从自由度为 1 的 χ^2 分布，即标准正态分布变量的平方。给定显著度 α ，当以上变量值小于临界值 χ_{α}^2 时，不能拒绝假设，即认为两学习器的性能没有显著差别；否则拒绝假设，即认为两者性能有显著差别，且平均错误率较小的那个学习器性能较优。自由度为 1 的 χ^2 检验的临界值当 $\alpha = 0.05$ 时为 3.8415， $\alpha = 0.1$ 时为 2.7055。

(3) Friedman检验

表 2.5 算法比较序值表

数据集	算法 A	算法 B	算法 C
D_1	1	2	3
D_2	1	2.5	2.5
D_3	1	2	3
D_4	1	2	3
平均序值	1	2.125	2.875

序值是每个算法在某一数据集上的好坏排序

横轴为平均序值，每个算法圆点为其平均序值，线段为临界阈值的大小

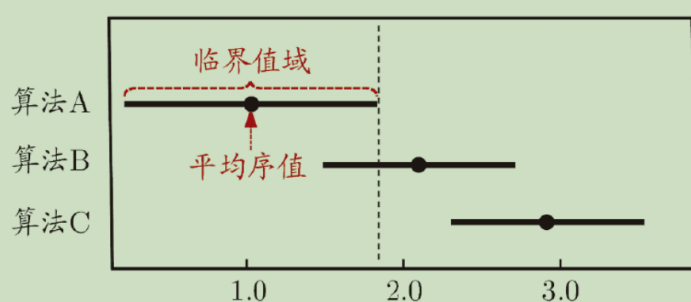


图 2.8 Friedman 检验图

若两个算法有交叠 (A 和 B)，则说明没有显著差别；
否则有显著差别 (A 和 C)，算法 A 显著优于算法 C

(4) Nemenyi检验

四、偏差&方差

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所共同决定的

- 偏差bias: 期望输出与真实标记的差别，刻画算法拟合能力

样本增多，偏差不会变

- 方差: 同样大小训练集的变动导致的学习性能变化，刻画数据扰动的影响
- 噪声: 当前任务任意算法所能达到的期望泛化误差下界，刻画学习问题难度

对回归任务，泛化误差可通过“偏差-方差分解”拆解为：

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

期望输出与真实输出的差别

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

同样大小的训练集的变动, 所导致的性能变化

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

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

训练样本的标记与真实标记有区别

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

偏差-方差窘境

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

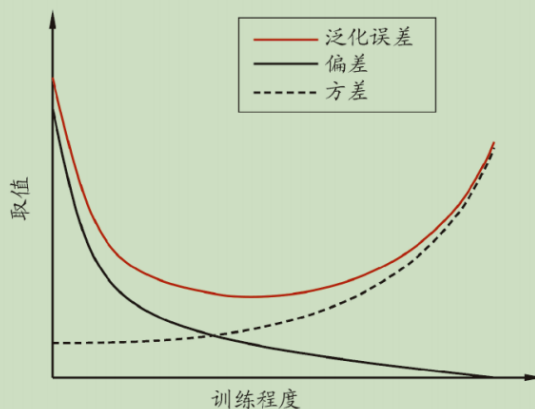


图 2.9 泛化误差与偏差、方差的关系示意图

偏差大是欠拟合（可以增加特征），方差大是过拟合