

第二章

1.2 泛化：在训练集上训练模型，模型应用
于新样本的能力

1.4

① 归纳偏好：在学习过程中对某种类型假设的偏好。（什么模型更好？假设）

可看作学习者在庞大假设空间中，对假设进行选择的启发式或“价值观”

“奥卡姆剃刀”

归纳偏好是否与问题匹配

② 没有免费的午餐

前提：所有“问题”出现机会相同，或所有问题同等重要，但很

有时只关注正在解决的问题

考虑是单问题

第二章 模型评估与选择

2.1 经验误差 过拟合

经验误差：在训练集上

泛化误差：在新样本上

2.2 评估方法

① 留出法：分层采样

多次使用留出法，每次随机划分得到不同的训练测试集，最后平均。

1) 测试集小：结果方差较大

2) 训练集小：结果偏差较大

② K折交叉验证

10次10折交叉验证=100次留出

③ 自助法

可重复采样，有放回采样

m个样本数据集D

每次随机选一个，再放回

重复m次得到包含m个样本的数据集

适用：数据集小，难以有效划分，利于集成学习
改变了分布，估计偏差

2. 性能度量

① 均方误差 $E = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$

② 准确率 召回率

真实 正 预测 反

TP 真正

FN 假正

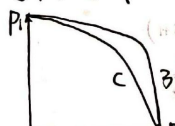
FP 假正

TN 真负

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

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

③ PR曲线



若完全与随机无异
则前者优于后者

2.4 比较检验

1) 假设检验

"假设": 对某些泛化错误率分布:

某种判断情景

测试错误率 $\xrightarrow{\text{假设}}$ 泛化错误率分布

泛化错误率 $\hat{\epsilon}$ 按测试错误率 $\hat{\epsilon}$

$$P(\hat{\epsilon}; \epsilon) = \binom{m}{\hat{\epsilon} \times m} \epsilon^{\hat{\epsilon} \times m} (1-\epsilon)^{m-\hat{\epsilon} \times m}$$

$$\frac{\partial P(\hat{\epsilon}; \epsilon)}{\partial \epsilon} = 0 \Rightarrow \epsilon = \hat{\epsilon} \text{ 时最大}$$

2.5 偏差与方差

$$\bar{f}(x) = E_D[f(x; D)] \quad \text{期望泛化}$$

$$\text{var}(x) = E_D[(f(x; D) - \bar{f}(x))^2] \quad \text{不同训练集方差}$$

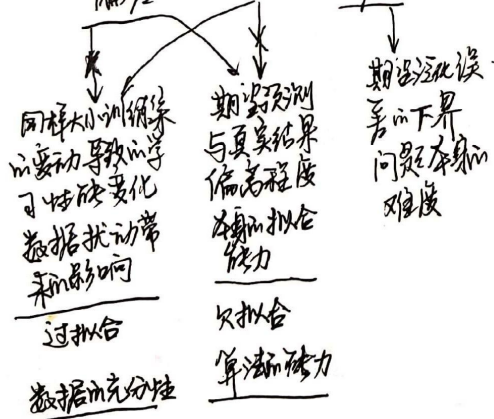
$$\epsilon^2 = E_D[(y_D - y)^2] \quad \text{噪声}$$

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

期望泛化误差:

$$E(f; D) = E_D[(f(x; D) - \bar{f}(x))^2] + (\bar{f}(x) - y)^2 + E_D[(y_D - y)^2]$$

$$= \underbrace{\text{bias}^2(x)}_{\text{偏差}} + \underbrace{\text{var}(x)}_{\text{方差}} + \underbrace{\epsilon}_{\text{噪声}}$$



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

调和平均 更重视较小值

对 P, R 重视程度不同.

相对于算术平均和几何平均.

1) 商品推荐: 尽可能少打扰用户 明确用户感兴趣: P 更重要

2) 逃犯信息检索: 不能漏掉: R 更重要

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

β 度量了 R 对 P 的相对重要性

$\beta > 1$: R 更重要

$\beta < 1$: P 更重要

③ ROC, AUC

将预测值和分类阈值比较

根据这个对测试样本排序.

以某一截断点 分为两部分, 前半为正面, 后半为负面.

根据阈值来设置截断点

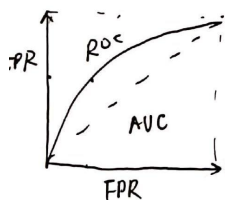
阈值 P: 靠前截断

阈值 R: 靠后截断

1) ROC 要查看工作特征

纵轴: 真正例率 $TPR = \frac{TP}{TP + FN}$

横轴: 假正例率 $FPR = \frac{FP}{FP + TP}$



根据预测结果排序 按顺序逐一把样本作为正例进行预测

(0,1) 对应于所有正例在负面

前面的理想模型.

真正例 $(x, y + \frac{1}{m})$

假正例 $(x + \frac{1}{m}, y)$

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

样本预测的排序质量, 排序误差

the steps you take don't need to be big
They just need to take you in the right direction

$$loss = \frac{1}{m \times n} \sum_{x' \in U^+} \sum_{x \in U^-} (I(f(x') < f(x)) + \frac{1}{2} I(f(x') = f(x)))$$

④ 代价敏感错误率, 代价曲线

患者 → 健康人

健康人 → 患者

权衡不同类型错误的代价. 非均匀代价

二分类任务为例:

设定代价矩阵: $cost_{ij}$ 表示将第 i 类样本

预测为第 j 类的代价 $cost_{ii} = 0$

	预测	
真实	0	1
0	0	$cost_{01}$
1	$cost_{10}$	0

非均匀代价下: "代价敏感"错误率

$$E(f; D; cost) = \frac{1}{n} \sum_{x_i \in U^+} (I(f(x_i) \neq y_i) \times cost_{01}) + \frac{1}{n} \sum_{x_i \in U^-} (I(f(x_i) \neq y_i) \times cost_{10})$$

非均匀代价下: ROC 曲线不行.

代价曲线:

横轴: 正例概率代价

$$P(+)cost = \frac{p \times cost_{01}}{p \times cost_{01} + (1-p) \times cost_{10}}$$

纵轴: [0,1] 为统一代价

$$cost_{norm} = \frac{FNR \times p \times cost_{01} + FPR \times (1-p) \times cost_{10}}{p \times cost_{01} + (1-p) \times cost_{10}}$$

p : 样例为正例的概率

