

第二章 模型评估与选择

王星 参考文献：周志华第二章

中国人民大学统计学院

March 29, 2021

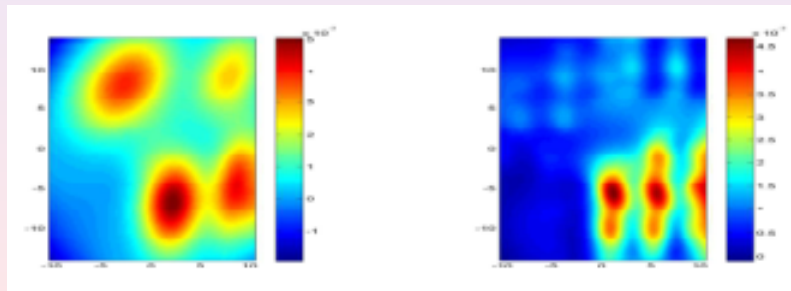
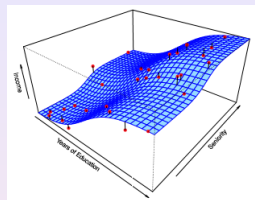
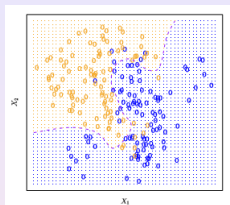
- 统计学习三要素
- 期望风险、经验风险和结构风险
- 泛化误差的几种近似计算（留出法，交叉验证和自助法）
- 性能指标（错误率、查准率，查全率，F1,ROC,AUC,代价敏感曲线CC）；
- 比较不同学习期的性能的假设检验
- 性能差异的理解：偏差与方差

什么是统计学习

- 统计学习是关于计算机基于数据构建概率统计模型并运用模型对数据进行预测与分析的一门学科—— 李航
- 统计学习(Statistical Learning)是一套以理解数据为目的的庞大的工具集，数据是信息的载体，但不全是信息。
- 统计学习中主要分为预测（prediction）和推断(inference)
 - 预测的例子：哪些人对这个产品会有兴趣？
 - **Example:** direct mailing prediction: Interested in predicting how much money an individual will donate based on observations from 90,000 people on which we have recorded over 400 different characteristics.
 - Don't care too much about each individual characteristic.
 - Just want to know: For a given individual should I send out a mailing?

- 推断的例子：有兴趣的人群有怎样的特征和分布, A type of relationship between y and the x 's.
 - **Example:** Which particular predictors actually affect the response?
 - 影响的方向 Is the relationship positive or negative?
 - 影响的模式 Is the relationship a simple linear one or is it more complicated etc.?
 - **Example:** housing inference: Wish to predict median house price based on 14 variables.
 - 影响的程度 Probably want to understand which factors have the biggest effect on the response and how big the effect is.
 - 影响的范围 For example how much impact does a river view have on the house value etc.
- 在这些例子中通常会有模式的存在性认识，随机特征的辨析，可能性的估计和参数的推断等一系列复杂的建模问题。

几个例子



统计学习三要素

- 方法=模型+策略+算法

- 模型: 假设空间中的函数, 如果假设空间用 \mathcal{F} 表示, 那么假设空间就是
 - 决策函数:

$$\mathcal{F} = \{f|Y = f(X)\},$$

如果模型是由参数控制的, 又称为参数空间:

$$\mathcal{F} = \{f|Y = f_{\theta}(X), \theta \in R^n\},$$

- 条件概率: $\mathcal{F} = \{P|P(Y|X)\}$;或 $\mathcal{F} = \{P|P_{\theta}(Y|X), \theta \in R^n\}$;
- 策略: 有了模型的假设空间, 需要考虑按照怎样的准则来学习到最优的模型
 - 损失函数(代价函数): 在假设空间 \mathcal{F} 上选取函数 f 作为决策函数, 对于给定的输入, 由 f 产生一个相应的输出 Y , 这个输出与真实的 Y 之间可能一致也可能不一致, 用损失函数或代价函数 $L(Y, f(X))$ 度量错误的程度。
- 算法: 算法是指学习模型的具体计算方法, 很多统计学习问题归为最优化问题, 如果最优化问题有显示解析解, 求解的方法比较直接, 如果没有显式解, 如何保证解的全局最优和高效性, 是算法中的一个重要问题。

损失函数

- 统计学习中的常用的损失函数如下：

- 0-1损失

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

- 平方损失函数(Quadratic Loss Function):

$$L(Y, f(X)) = (Y - f(X))^2;$$

- 绝对损失函数(Absolute Loss Function):

$$L(Y, f(X)) = |Y - f(X)|;$$

- 对数损失函数(Logarithmic Loss Function):

$$L(Y, P(Y|X)) = -\log P(Y|X).$$

- 损失函数性质：损失函数越小，模型越好，然而模型的输入和输出 (X, Y) 是随机变量，与联合分布 $P(X, Y)$ 有关系，所以通常用损失函数的期望来表达：
- 风险函数(期望损失):

$$R_{\text{exp}}(f) = E_P[L((Y, f(X)))] = \int_{x,y} L(y, f(x))P(x, y)dxdy;$$

损失函数 (2)

- 学习的目标是期望风险最小，也就是说：

$$f^* = \operatorname{argmin} R_{\text{exp}}(f)$$

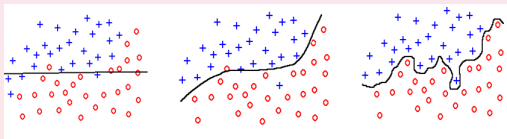
- 问题： $p(x, y)$ 未知。
- 给定一个训练数据 $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, 称模型 $f(X)$ 关于训练数据的平均损失为经验风险 (Empirical Risk), 记作：

$$R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)).$$

损失函数 (3)

[定理] x_i, y_i 独立同分布的情况下，经验风险是期望风险的无偏估计

$$\begin{aligned} \mathbb{E}[R_{\text{emp}}(f)] &= \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \right] \\ &= \frac{1}{N} \mathbb{E} \sum_{i=1}^N [L(y_i, f(x_i))] \quad x_i \text{ 独立} \\ &= \frac{1}{N} \mathbb{E} \sum_{i=1}^N L[(y_i, f(x_i))] \quad x_i \text{ 同分布} \\ &= R_{\text{exp}}(f) \end{aligned}$$



结构风险最小化

- Empirical Risk Minimization(ERM):

$$f^*(D_N) = \operatorname{argmin}_{f \in \mathcal{F}} L(f, D_N) = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

- 训练数据总是有限的，有噪声甚至数量不足，用训练数据上的经验风险估计期望风险并不理想，需要对经验风险作一定的矫正——结构风险最小化

- Structure Risk Minimization(SRM):

$$R_{\text{SRM}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f)$$

$J(f)$ 是模型的复杂度， f 越简单，复杂度 $J(f)$ 就越小， f 越复杂，复杂度 $J(f)$ 就越大。

$$f^*(D_N) = \operatorname{argmin}_{f \in \mathcal{F}} L(f, D_N) = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f)$$

- 算法是指学习模型的具体算法;
- 例如：最优化学习算法，梯度学习算法，深度学习算法等;

例1

- 对于二分类问题 $Y = \{\pm 0, 1\}$,
- 假设 $X \times Y$ 的联合分布由 (μ, η) 表示，其中 μ 是 X 的边缘分布， η 是 Y 的条件分布， $\eta(X) := P(Y = 1|X)$;
- 损失函数 $L[(\hat{y}, y) = 1|y \neq \hat{y}]$,
- 于是可以计算任意一个可测函数 f 的风险

$$R(f) = \text{EL}(f(X), Y) = P(f(X)) \neq Y).$$

- 定义 R^* 为最优风险， $R^* = \operatorname{argmin}_f R(f)$.
- 定义：0-1损失

$$f^*(x) = \begin{cases} 1, & \eta(X) \geq 1/2 \\ 0, & \text{otherwise} \end{cases}$$

定理 对任意的 $f : X \rightarrow Y$,

$$R(f) - R(f^*) = E(1[f(X) \neq f^*(X)]|2\eta(X) - 1|)$$

证明

$$\begin{aligned} & R(f) - R(f^*) \\ &= \mathbb{E}_{X,Y}\{(1[f(X) \neq Y]) - (1[f^*(X) \neq Y])\} \\ &= \mathbb{E}_{X,Y}\{1[f(X) \neq f^*(X)]([1[f(X) \neq Y] - 1[f^*(X) \neq Y]])\} \\ &= \mathbb{E}_{X,Y}\{1[f(X) \neq f^*(X)](2 \cdot 1[f^*(X) = Y] - 1)\} \\ &= \mathbb{E}_X\{1[f(X) \neq f^*(X)][2\mathbb{E}_{Y|X}1(f^*(X) = Y) - 1]\} \\ &= \mathbb{E}_X\{1[f(X) \neq f^*(X)][2(1/2 + |\eta(X) - 1/2|) - 1]\} \\ &= \mathbb{E}_X 1[f(X) \neq f^*(X)](|2\eta(X) - 1|) \end{aligned}$$

这个定理表明在二分类问题中，最优的决策函数 f^* 形式是由后验概率 $\eta(x)$.定义的

2.1 经验误差与过拟合

经验误差:假设学习到的模型是 $Y = \hat{f}(x)$

- **训练误差**(training error): 在训练集上预测的结果与训练集的真实结果的差异,也称为经验误差 (empirical error),代表学习器预测的结果与实际真实结果之间的差异:

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

N 是样本量

- **训练错误率**: 分类错误的样本数(a)占训练样本总数(N)的比例 $E = a/N$,称 $1 - a/N$ 为精度
- **泛化误差** (generalization error): 测试集或新样本上的预测误差

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

- 很明显,要使分类器尽可能的有用,应该要让泛化误差仅可能的小。然而现实环境中,很难知道新样本是什么样的,所以我们实际能做的只有努力使经验误差最小化,但经验误差又不合适。

在泛化误差未知的情况下，如何对模型进行评估呢？

- Larson于20世纪30年代就提出：在相同的数据上训练算法和评价算法的性能将会得到“过于乐观”的结果。
- 思路是：既然没法拿到新样本来进行泛化误差计算，那么可以从训练样本中取出一部分来，假装它是一组新样本，并用这些样本计算出来的误差作为泛化误差的近似。这组选出来的样本被称为“测试集”，测试集上的误差被称为测试误差。
- 留出法（hold-out）：将数据集 S 分为 $D = S \cup T, S \cap T = \emptyset, S$ 训练模型， T 拿来测试。 N_T 是 T 的样本量，泛化误差表示为：

$$\hat{R}_{\text{test}}^{HO}(\hat{f}) = \frac{1}{N_T} \sum_{x_i \in T} L(y_i, \hat{f}(x_i))$$

- 交叉验证法（cross validation）先将数据划分为 k 个大小相似的互斥子集， $D = D_1 \cup \dots \cup D_k, D_i \cap D_j = \emptyset, \forall i \neq j$ 分为多次Hold-out, 留 p 交叉, Repeated learning-testing (RLT) 交叉验证估计和 k -折交叉。
- 自助法（bootstrapping）：从 m 个样本的数据集 D ，随机采样生成训练集 D' 。用 D' 去估计泛化误差。

- 要有足够的样本量，以保证训练模型的效果；
- 在划分时注意保证数据分布的一致性（如：500个样本中正例和反例的比为2:3，则在训练集和测试集中正例和反例的比也要求为2:3），只需要采用随机分层抽样即可；
- 为降低随机划分的影响，重复划分训练集和测试集，对得到的多次结果取平均作为最后的结果
- 优点：Hold-out估计法的提出打破了传统的基于相同的数据进行训练和测试的分析，避免了训练和测试数据重叠引起的过拟合。
- 缺点: Hold-out估计过分依赖于某一次数据划分，数据的划分方式直接影响着估计的精度，经验上 $2/3 \sim 4/5$ 。

```
1 from sklearn.model_selection import train_test_split
2 #使用train_test_split划分训练集和测试集
3 train_X , test_X, train_Y ,test_Y = train_test_split(
4     X, Y, test_size=0.2,random_state=0)
5 ...
6 X为原始数据的自变量，Y为原始数据因变量；
7 train_X, test_X是将X按照8:2划分所得；
8 train_Y, test_Y是将Y按照8:2划分所得；
9 test_size是划分比例；
10 random_state设置是否使用随机数
```

交叉验证法1:多次Hold-out估计法

- 注意到hold-out估计依赖于数据的一次划分，容易受到数据划分随机性的影响，Geisser(1975)[1]提出了包含多次hold-out 估计平均的交叉验证方法的一个一般的表示，实现了从验证估计到交叉验证估计的过渡。因为多次数据划分导致数据之间有交叉，称为交叉验证

[1]GEISSER S. The predictive sample reuse method with applications[J]. Journal of the American Statistical Association,1975,70(35):320-328

- 如果记 S_1, \dots, S_K 为 D 的 K 个非空真子集, T_1, \dots, T_K 是分别对应的补集, 即 $S_i \cap T_i = \emptyset, S_i \cup T_i = D, \forall i = 1, \dots, K$ 此时泛化误差估计为:

$$\hat{R}_{\text{test}}^{MHO}(\hat{f}) = \frac{1}{K} \sum_{k=1}^K \hat{R}_{\text{test}(T_k)}^{HO}$$

- 优点:方法依赖于数据的多次重复划分, 可去掉划分随机性的影响(在统计分析中就是通过多次重复实验来减小方差.)
- 不足: 交叉验证数据划分方式有很多种, 随着样本量的增加, 划分方式的组合数也在急剧地增加, 现实的计算中很难穷尽, 计算复杂度非常高。

交叉验证法2: 留 p (leave- p -out:LPO)交叉验证估计

- 为减小多次Hold-out交叉验证中数据划分的组合数, Shao[2]提出了每次数据划分中测试样本个数都相同的留 p 交叉验证

[2]SHAO Jun. Linear model selection by cross-validation[J],Journal of the American Statistical Association, 1993, 88(422):486-494.

即训练样本容量 $N_S = N - p$, 测试集容量 $N_T = p$, 其中 $p \in 1, \dots, N - 1$ 。这样, 如果有 K 种划分 T_1, \dots, T_K , 测试集样本量均为 p , 留 p 交叉验证泛化误差估计为:

$$\hat{R}_{\text{test}}^{LPO}(\hat{f}) = \frac{1}{K} \sum_{k=1}^K \frac{1}{p} \sum_{x_j \in T_j} L(y_j, \hat{f}(x_j))$$

- 优点: 留 p 交叉验证相比于多次留出给出的交叉验证的划分组合数从 $\sum_{p=1}^{N-1} C_N^p$ 减少到了 C_N^p 。
- 不足: 对于较大的样本量 N , 组合数计算复杂度 C_N^p 仍然很大。
- $p = 1$ 时为留一(leave-one-out:LOO)交叉验证估计。

$$\hat{R}_{\text{test}}^{LOO}(\hat{f}) = \frac{1}{N} \sum_{x_j \in D_{-j}} L(y_j, \hat{f}(x_j))$$

交叉验证法3:RLT Repeated learning-testing

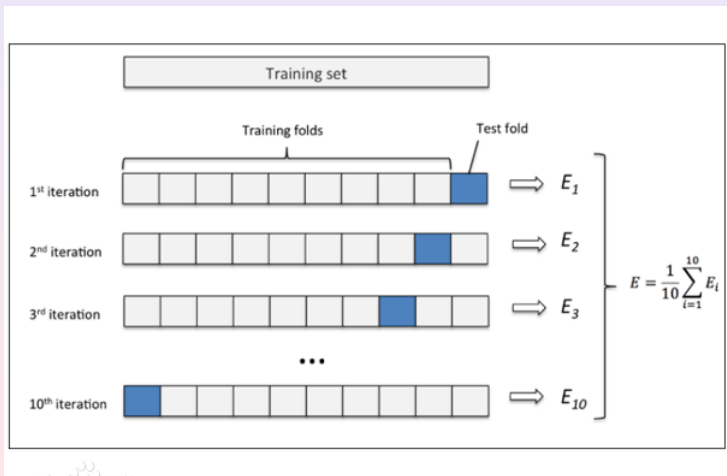
- 与留 p 交叉验证考虑所有数据划分相比,RLT交叉验证只基于部分数据划分进行,选择任意大于0 的 K 次数据划分进行泛化误差的估计.
- RLT交叉验证泛化误差估计为:

$$\hat{R}_{\text{test}}^{\text{RLT}}(\hat{f}) = \frac{1}{K} \sum_{k=1}^K \frac{1}{n_T} \sum_{x_j \in T_j} L(y_j, \hat{f}(x_j))$$

- 优点:RLT自从1981年被提出以后在实际应用中就被广泛使用,因为它有可接受的计算开销且操作简单。
- 缺点: K 的大小的选择一直是这个方法的最大问题,不同文献中有不同的结论,如文献[5] 中建议 $K = 15$,训练和测试集比例的选择也没有一个确定的结论,往往不同的研究者使用不同的训练和测试容量。
- [3]ARLOT S,CELISSE A. A survey of cross-validation procedures for model selection[J]. Statistics Surveys,2010,4:40-79.
- [4]ZHANG Ping. Model selection via multifold cross validation[J]. Annals of Statistics,1993,21(1):299-313.
- [5]NADEAU C, BENGIO Y. Inference for the generalization error[J]. Machine Learning,2003,52(3):239-281.

交叉验证法4: K -折交叉验证估计

将 D 划分为互斥的 k 个子集，每次训练用其中 $(k - 1)$ 个数据，用剩下的一个作为测试集。这样就能获得 k 组训练/测试集，对这 k 组分别进行训练和测试，最终返回 k 个测试结果的均值



从 m 个样本的数据集 D ，随机采样生成训练集 D' .用 D' 去顾及泛化误差。样本在 m 次采样中始终不被采集到的概率如下：

$$\lim_{m \rightarrow \infty} (1 - 1/m)^m = e = 0.368.$$

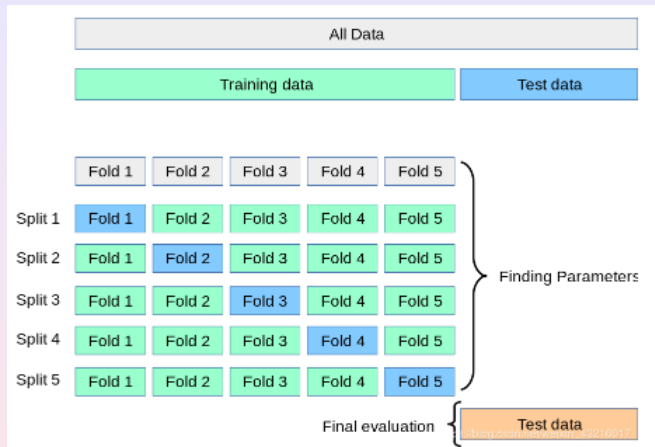
自助法在数据集较小、难以有效划分训练/测试集时很有用,自助法产生的数据集改变了初始数据集的分布，这会引入估计偏差。因此，在初始数据量足够时，留出法和交叉验证法更常用一些

3.6课后作业

对于数据Auto,

- (a)将数据集按照留出法分为训练集(比例0.8)和测试集(比例0.2), 对于函数 $mpg = 40 - 0.15 \times horsepower$, 编写程序估计平方损失下的泛化误差,自行设定试验次数, 进行泛化;
- (b)将数据集按照留 p 交叉验证的方式提取测试集, $K = 20$ 划分采取无放回抽样 $p = 10$,对于函数 $mpg = 40 - 0.15 \times horsepower$, 编写程序计算平方损失下的泛化误差;
- (c)建立一个二元变量, 每加仑里程量 ($mpg01$), 1表示每加仑汽油该型号车所跑里程数 (mpg) 在3/4 分位数以上, 0 表示每加仑汽油该型号车所跑里程数 (mpg) 在3/4分位数以下。将数据集按照留出法随机分为训练集(比例0.8)和测试集(比例0.2), 对于函数 $P(mpg01 = 1) = \frac{\exp\{3.85 - 0.01 \times weight\}}{1 + \exp\{3.85 - 0.01 \times weight\}}$ 编写程序估计平方损失下的泛化误差;
- (d)将数据集按照分层等比例分按照训练集(比例0.8)和测试集(比例0.2), 对于函数 $P(mpg01 = 1) = \frac{\exp\{3.85 - 0.01 \times weight\}}{1 + \exp\{3.85 - 0.01 \times weight\}}$ 编写程序估计平方损失下的泛化误差;
- (e)根据以上实验, 结合教材, 分析分层抽样和不同的抽样方式对泛化误差的影响;
- (f)书上的2.1.

- 不同的参数一般代表着不同的模型，特别是同一类的不同模型。一般参数分三种：
 - 一类是函数的参数，降低编程的难度，方便代码重用提高效率，比如，形参（左图中的参数表示颜色，线形）停止条件要给出估计的精度，参数传参时，可以通过位置参数传参或者关键字参数传参，位置可变参数可以接受多个可以接收多个实参，一个形参可以匹配任意个参数。
 - 一类是算法的参数，亦称“超参数”，如聚类所需要的簇数量 k ，人工选择一组候选。
 - 一类是模型的参数，如神经网络中每个节点的权重，回归中的每个变量的系数，让机器自己学习。
- 调参是机器学习的重点，也是决定模型性能的关键。一般调参过程中，会将训练数据再次划分为训练集和验证集（validation set）。具体包含关系如下：
（数据集（训练数据（训练集）（验证集））（测试集））



2.2 性能度量-回归问题和度量

- 回归问题的目标是在给定 d 维输入 (input) 变量 x 的情况下, 预测一个或者多个连续目标 (target) 变量 y 的值。
- 对于 $D = \{(x_1, y_1), (x_N, y_N), \dots\}$, 要评估学习器 f 的性能;
- 回归中使用平方损失The regression loss-function:
 $L(y, f(x)) = (f(x) - y)^2$;
- 期望损失:

$$\mathbb{E}(L) = \int \int \{f(x) - y\}^2 p(x, y) dx dy$$

- 最小化 $\mathbb{E}(L)$, 得到: $2 \int (f(x) - y) p(x, y) dy = 0$.
- 求解出 $f(x)$:

$$f(x) = \frac{\int y p(x, y) dy}{p(x)} = \int y p(y|x) dy = \mathbb{E}_y(y|x)$$

- 最优预测Solution: $f(x) = \int y p(y|x) dy = \mathbb{E}_y[y|x]$;

2.3 二分类问题的错误率与精度

- 等权错误率:

$$E(f, D) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(x_i) \neq y_i);$$

- 等权精度:

$$\text{ACC}(f, D) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(x_i) = y_i) = 1 - E(f, D);$$

- 一般错误率:

$$E(f, D) = \frac{1}{N} \int_{x \in D} \mathbb{I}(f(x) \neq y) p(x) dx;$$

- 一般精度:

$$\text{ACC}(f, D) = \frac{1}{N} \int_{x \in D} \mathbb{I}(f(x) = y) p(x) dx;$$

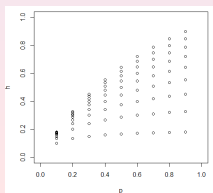
2.4 查准率和查全率

- **查准率P**表示在所有被认为是正确的样例中，有多少个被正确分类的样例；
- **查全率R**表示在所有本身即是正确的样例中，有多少个被正确分类的样例。
- 对于二分类问题，分类结果表达为混淆矩阵（confusion matrix）

Table 2. Confusion Matrix.

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

$$\text{precision} = \frac{TP}{TP+FP} \quad \text{recall} = \frac{TP}{TP+FN} \quad F1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (1)$$



P-R曲线

分类器依次逐个将测试样例作为正例进行预测，则可以计算出每例的查准率和查全率。以查准率为纵轴y轴、查全率为横轴x轴作图，就得到了查准率-查全率曲线，简称“P-R曲线”。

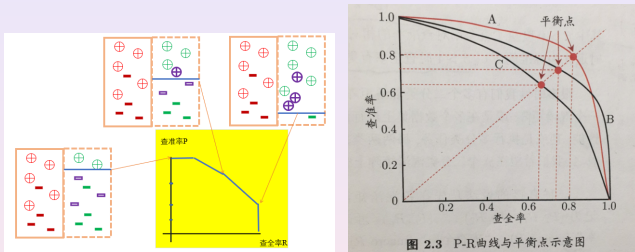


图 2.3 P-R曲线与平衡点示意图

- 若一个学习器的P-R曲线被另一个学习器的曲线完全“包住”，则可以断言后者的性能优于前者。例如上图中，学习器A的性能优于学习器C。
- 查准率和查全率是一对互为矛盾的度量，当查准率高时，查全率会偏低，当查全率高时，查准率会偏低，R一定会是零吗？
- **平衡点(Break-Even Point,简称BEP)**是一个用于比较学习器的性能度量，它是查准率=查全率时的取值，如图中的BEP是0.55。基于BEP值的高低可以比较出学习器的优劣。

P-R曲线2

$P(+)$ 大	\longrightarrow	\longrightarrow	$P(+)$ 小
+	+	-	-
+	+	+	-
+	-	+	-
+	+	-	-
$P=1, R=0$ $FP=0$	$FP \uparrow P \downarrow$ $FN \downarrow R \uparrow$	$FP \uparrow P \downarrow$ $FN=0, R \uparrow$	$R=1$

- 一些应用中，对查全率和查准率重视程度不一样，比如，推荐系统中，为更少打扰用户，被推荐的内容恰恰是用户感到有兴趣的，于是查准率很重要，而在治疗卵巢癌的手术时，为尽可能切除掉病灶，查全率更重要，根据查全率和查准率的不同偏好，使用 F_β 来表达对查全率和查准率的不同偏好：

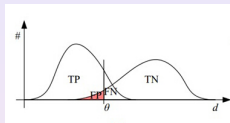
$$\frac{1}{F_\beta} = \frac{1}{1 + \beta^2} \times \left(\frac{1}{P} + \frac{\beta^2}{R} \right)$$

- 与算术平均 $\frac{P+R}{2}$ 和几何平均 $\sqrt{P \times R}$ 相比，调和平均更重视较小值。于是有：

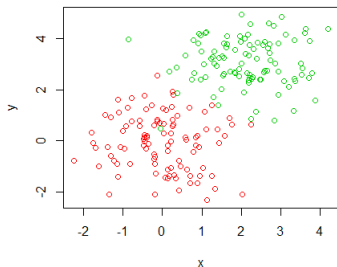
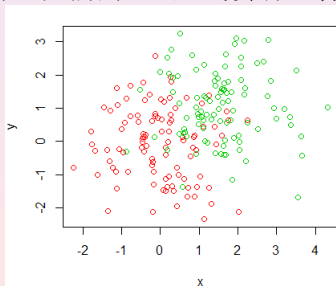
$$F_\beta = \frac{(1 + \beta^2) \times P \times R}{\beta^2 \times P + R}$$

思考题

$\beta > 1$, 查全率有更大影响, $\beta < 1$, 查准率有更大影响。 $\beta = 1$, 退化为 F_1



- 对于图中的两个图, 你认为如果用线性模型来划分两组数据, 他们的P-R曲线会有怎样的不同?



宏查准率、宏查准率、微查准率、微查全率

- 在前面的交叉验证法中对泛化误差的估计有很多，会有多次训练和测试，每次得到一个混淆矩阵，假设一共获得 $(P_1, Q_1), (P_2, Q_2), \dots, (P_K, Q_K)$ 个不同的查准率和查全率，
- 宏查全率，宏查准率：

$$\text{macro-P} = \frac{1}{K} \sum_{i=1}^K P_i; \quad \text{macro-R} = \frac{1}{K} \sum_{i=1}^K R_i;$$

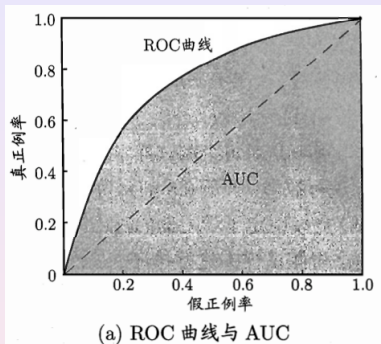
$$\text{macro-F1} = \frac{2 \times \text{macro-P} \times \text{macro-R}}{\text{macro-P} + \text{macro-R}}$$

- 微查全率，微查准率：

$$\text{micro-P} = \frac{\overline{TP}}{\overline{TP} + \overline{FP}}; \quad \text{micro-R} = \frac{\overline{TP}}{\overline{TP} + \overline{FN}};$$

$$\text{micro-R1} = \frac{2 \times \text{micro-R} \times \text{micro-R}}{\text{micro-P} + \text{micro-R}}$$

ROC(Receiver Operating Characteristic)与AUC



$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad \text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

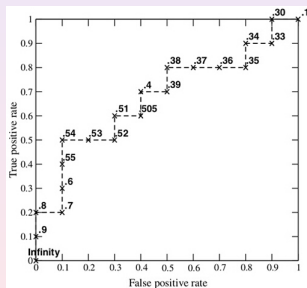
- Equal Error rate

如何绘制ROC与AUC

- 假设已经得出一系列样本被划分为正类的概率，然后按照大小排序，下表是一个示例，表中共有20个测试样本，“Class”一栏表示每个测试样本真正的标签（ p 表示正样本， n 表示负样本），“Score”表示每个测试样本属于正样本的概率。

Inst#	Class	Score	Inst#	Class	Score
1	p	.9	11	p	.4
2	p	.8	12	n	.39
3	n	.7	13	p	.38
4	p	.6	14	n	.37
5	p	.55	15	n	.36
6	p	.54	16	n	.35
7	n	.53	17	p	.34
8	n	.52	18	n	.33
9	p	.51	19	p	.30
10	n	.505	20	n	.1

- 接下来，从高到低，依次将“Score”值作为阈值threshold，当测试样本属于正样本的概率大于或等于这个threshold时，我们认为它是正样本，否则为负样本。举例来说，对于图中的第4个样本，其“Score”值为0.6，那么样本1，2，3，4都被认为是正样本，因为它们的“Score”值都大于等于0.6，而其他样本则都认为是负样本。每次选取一个不同的threshold，我们就可以得到一组FPR和TPR，即ROC曲线上的一点。这样一来，我们一共得到了20组FPR和TPR的值，将它们画在ROC曲线的结果如下图：



- 设当前点为 (x, y) ,当测试样本属于正样本时，下一个点是 $(x, y + \frac{1}{N^+})$, 否则就是 $(x + \frac{1}{N^-}, y)$

- AUC表示ROC曲线下的面积，用于比较学习器的性能优劣，一般面积更大的学习器的性能更为优良。AUC 可估算为

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

- 形式化地看，AUC 考虑的是样本预测的排序质量，因此它与排序误差有紧密联系。给定 N^+ 个正例和 N^- 个反例，令 D^+ 和 D^- 分别表示正、反例集合，则排序损失（loss）定义为

$$l_{\text{rank}} = \frac{1}{N^+N^-} \sum_{x^+ \in D^+} \sum_{x^- \in D^-} \left(\mathbb{I}(f(x^+) < f(x^-)) + \frac{1}{2} \mathbb{I}(f(x^+) = f(x^-)) \right).$$

它刻画了正例得分没有高出负例得分的那部分损失， l_{rank} 对应的是ROC曲线以上的面积

- $\text{AUC} = 1 - l_{\text{rank}}$;

ROC曲线的基本性质1/2

- 性质1: 任何ROC曲线必定经过原点和(1,1)

证明: 当阈值大于所有样本的预测值时, 所有的样本都会被归为负类, 这时正例数为0, 因此FPR和TPR都为0, 对应于ROC曲线的原点, 此时的分类方法是最“保守的”。当阈值小于所有样本的预测值时, 所有的样本都会被归为正类, 此时负例数为0, 因此 $TN = FN = 0$ 进而FPR和TPR都为1, 对应于ROC曲线的(1,1)点。

- 性质2: 位于上方的ROC曲线分类效果好于下方的ROC曲线

证明: 在给定样本分布的情况下, 假设有两个分类器A和B, A的ROC曲线高于另一个分类器B的ROC曲线, 说明在相同的假正例率的条件下, A 的真正例率大于B, 因此分类器A要优于分类器B。

ROC曲线的基本性质2/2

- **性质3:** 在样本有限的情况下, ROC曲线为由水平线和竖直线构成的折线

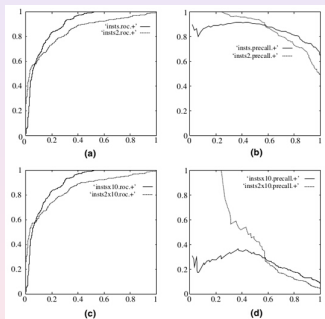
可以将所有样本的预测值由大到小排序。然后把分类器的阈值设为最大, 即所有样本都被归为反例, 这时对应于ROC曲线的原点。然后, 从高到低依次将阈值设为每个样本的预测值, 即依次把每个样本划分为正例。设前一个标记点的坐标为 (x, y) , 若当前为真正例, 则对应标记点的坐标为 $(x, y + 1/m^+)$, 若当前为假正例(负例), 则对应标记点的坐标为 $(x + 1/m^-, y)$, m^+, m^- 分别表示样本中正例和反例的总数, 这就证明了该性质。

- **性质4:** $y = x$ 曲线对应于随机猜测分类器在无穷多样本下的极限
证明: 对于随机猜测的情况, 被归为正例的样本中有一半是猜对的有一半是猜错的, 被归为负例的样本也一样。假设总共有 M 个样本归为正类, N 个样本归为负类,
则 $TN = FN = N/2, TP = FP = M/2$, 因此 $TPR = FPR = \frac{M}{M+N}$,
这对应于 $y = x$ 曲线。

- **性质5:** 任何一个有效的分类器都位于 $y = x$ 上方
证明: 由于任何有效的分类器都优于随机分类器, 由性质2和性质3可直接得到性质4。

ROC曲线的好特性

- **性质6: ROC曲线对样本的分布不敏感** 当测试集中的正负样本的分布不对等的时候, ROC曲线能够保持不变。在实际的数据集中经常会出现样本类不平衡, 即正负样本比例差距较大, 而且测试数据中的正负样本也可能随着时间变化。下图是ROC 曲线和Precision-Recall曲线的对比:



- 在上图中, (a)和(c)为Roc曲线, (b)和(d)为Precision-Recall曲线。(a)和(b)展示的是分类其在原始测试集(正负样本分布平衡)的结果, (c)(d)是将测试集中负样本的数量增加到原来的10倍后, 分类器的结果, 可以明显的看出, ROC曲线基本保持原貌,

代价敏感错误率与代价曲线1/2

		预测结果	
		正例0	负例1
真实情况	正例0	0	$Cost_{01}$
	负例1	$Cost_{10}$	0

$Cost_{10} = C(+|-)$ 表示实际为反例但预测成正例的代价; $Cost_{01} = C(-|+)$ 表示实际为正例但是预测为反例的代价



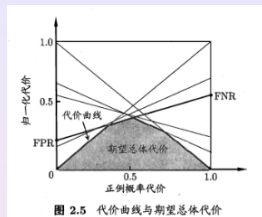
$$\mathbb{E}(f, D_i, cost) = \frac{1}{N} \left(\sum_{x_i \in D^+} \mathbb{I}(f(x_i) \neq y_i) cost_{01} + \sum_{x_i \in D^-} \mathbb{I}(f(x_i) \neq y_i) cost_{10} \right)$$

- 非均衡代价下，ROC曲线不能直接反映出学习器的期望总代价，于是出现了一个新的曲线“代价曲线”
- 分析理解：

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

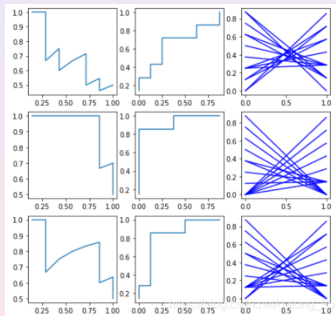
这里 $E(Cost) = Cost_{10} * P(Cost = 1|负例)P(负例) + Cost_{01} * P(Cost = 0|正例)P(正例)$

代价敏感错误率与代价曲线2/2



- 期望总体代价越小，则模型的泛化能力越强；
- 与ROC曲线的区别与练习：ROC主要考量均等代价，代价敏感曲线主要考量非均等代价。两者都是衡量某一学习器在不同场景下的综合表现情况，而不是单一场景。ROC通过阈值变化来体现不同场景，即高阈值表现了重视查准率的场景，低阈值则重视查全率的场景。代价敏感曲线则是通过 P 值，即正例的先验概率即原本正例占比的变化来体现不同场景。代价敏感曲线上方直线是根据不同决策阈值下做出的（含有参数 P ，固定参数），横轴 P 值确定时，即确定了一种情景，直线对应的点体现了 P 这种情景下在该阈值下的代价。因此当画出不同阈值下的直线时，某一 P 总能找到最小代价进而找出对应的决策阈值。这样每一个 P 都能对应一个在学习器下的最下代价，进而下方面积就是概率和代价的积分则为期望：该学习器的整体期望或整体表现。

如下9幅图为3个模型对同一个样本的分类结果。第1行为 $model_1$ 的 $P-R$ 图、 ROC 曲线、代价曲线。第2行为 $model_2$,第3行为 $model_3$ 。3个模型对比,判断哪个模型比较理想?



1. (03-21) 请根据例题中的测试得分绘制ROC曲线, $P-R$ 曲线
2. (03-29) 请根据例题中的测试得分绘制均衡代价曲线(注: $cost_{10} = 1, cost_{01} = 1$)。如果该数据是一个银行追查网络钓鱼攻击的一个学习器的训练数据, 漏网的代价(FN)是误报代价(FP)的5倍 ($cost_{10} = 1, cost_{01} = 5$) 那么请绘制代价曲线, 如果是3倍, 代价曲线会怎样变化, 请比较三种不同代价下的代价曲线有什么变化?
3. (03-29) 如下数据是16次实验中的误分类结果, 请比较测试误差率是否小于30% ($\epsilon \leq 30\%$), 实验中每次测试集 $n = 10$ 的误分类数量(mc) 如下, 凭借这些数据使用 t 检验对改算法的泛化误差给出判断? 请给出检验问题的表述, 计算统计量, 计算统计量的 p -值, 给出结论。

3	4	1	2	0	1	2	3
1	2	4	0	0	3	4	5

比较检验（比较什么？）

- 比较的是学习器的泛化性能，但性能度量观测 $(\epsilon_1, \dots, \epsilon_K)$ 可能和泛化性能 (ϵ) 并不相同；
- 测试集上的性能与测试集本身的选择有很大关系，测试集不同，结果可能也会发生改变，场景不同结果可能也不同；
- 很多机器学习算法具有很大的随机性，同一份测试集多次测试结果也可能不同(假设参数相同)；统计假设检验为进行学习算法的性能比较提供了重要依据。基于假设检验的结果可以推断出:若在测试集上观察到学习算法 L_A 比 L_B 好，则 L_A 的泛化性能是否在统计意义上优于 L_B ，以及这个结论的把握有多大。
- 假设是对学习器泛化错误率分布的某种判断或猜想，用测试错误率估计泛化错误率，以检查学习器性能。

机器学习中的假设检验

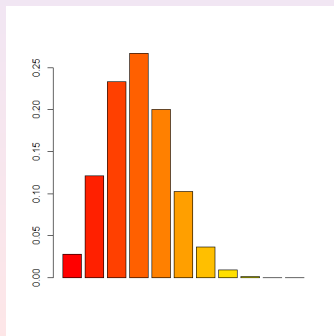
问题1 当问起一个概念学习算法的泛化误差，多数人的回答是 $\epsilon \leq 30\%$ ，如果这在行业中成为一种共识。一项研究中16次算法实验中每次测试集 $n = 10$ 的误分类数量(mc) 如下，凭借这些数据可以对这个泛化误差给出怎样的判断？

3	4	1	2	0	1	2	3
1	2	4	0	0	3	4	5

机器学习中的假设检验

问题1 当问起一个概念学习算法的泛化误差，多数人的回答是 $\epsilon \leq 30\%$ ，如果这在行业中成为一种共识。一项研究中16次算法实验中每次测试集 $n = 10$ 的误分类数量(mc) 如下，凭借这些数据可以对这个泛化误差给出怎样的判断？

3	4	1	2	0	1	2	3
1	2	4	0	0	3	4	5



多次留出法以及交叉验证中的 t 检验

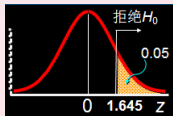
- 设立假设: $H_0 : \epsilon \leq \epsilon_0$ $H_1 : \epsilon > \epsilon_0$;
- (K 次) 训练/测试, 会产生多个测试错误率 $\epsilon_1, \dots, \epsilon_K$, 假设要比较的泛化误差率是 $\hat{\mu}$, 假设样本方差 $\hat{\sigma}^2$ 如下:

$$\hat{\mu}_{\epsilon} = \frac{1}{K} \sum_{k=1}^K \epsilon_k \quad \hat{\sigma}^2 = \frac{1}{K-1} \sum_{k=1}^K (\epsilon_k - \hat{\mu}_{\epsilon})^2.$$

- 将 K 个测试误差看作泛化误差率 ϵ_0 的独立采样, 定义检验统计量

$$T = \frac{\sqrt{K}(\hat{\mu}_{\epsilon} - \epsilon_0)}{\hat{\sigma}}.$$

- T 服从自由度为 $K-1$ 的 t 分布, 检验准则为: p 值为 $P(T > t_{\alpha}(K-1)) < \alpha$ 时拒绝零假设, 认为测试误差率 $> \epsilon_0$;



交叉验证 t -检验

- 基本原理：对两个学习器 A 和 B ，使用 k 折交叉验证法分别得到 k 个测试错误率，如果两个学习器性能相同，则使用相同训练/测试集时测试错误率应该相同. 设 A 的误差率 $\epsilon_1^A, \dots, \epsilon_k^A$ ，设 B 的误差率 $\epsilon_1^B, \dots, \epsilon_k^B$ ，
- 设立假设: $H_0 : \epsilon^A = \epsilon^B$ $H_1 : \epsilon^A \neq \epsilon^B$;
- 首先求两个学习器的 k 个测试错误率的差，即 $\Delta_k = \epsilon_k^A - \epsilon_k^B$ ，若

$$T_{\Delta} = \left| \frac{\sqrt{k} \times \mu_{\Delta}}{\sigma_{\Delta}} \right| < t_{\alpha/2, k-1}$$

, 则认为两个学习器性能相同, 否则, 认为两个学习器性能存在显著性差别, 而且学习错误率较小的那个学习器的性能较优。

•

5 × 2交叉验证t-检验

- 要进行有效的t检验，一个重要前提是测试误差率均为泛化误差的**独立采样**，然而由于样本有限，使用交叉验证等实验估计方法时，不同轮次的训练样本会有一定程序的重叠，这就使得测试误差率实际上并不独立，而且会导致过高估计假设成立的概率
 - 解决方案：采用5 × 2交叉验证法[Dieterich1998],每次2折交叉验证之前用随机数将数据打乱，使得5次交叉验证的数据划分不重复，对两个学习器A和B,第*i*次2折交叉验证将产生两对测试错误率，分别求差，记为 Δ_i^1, Δ_i^2 ，为缓解测试误差率的非独立性，仅计算第1次2折交叉验证的两个结果的平均值 $\mu = 0.5(\Delta_i^1 + \Delta_i^2)$ ，但是对每次2折实验的结果都计算方差 $\sigma_i^2 = (\Delta_i^1 - \mu)^2 + (\Delta_i^2 - \mu)^2$. 用统计量
 - $T = \frac{\mu}{\sqrt{0.2 \sum_{i=1}^5 \sigma_i^2}}$ 服从自由度5的t分布，取双边检验临界值 $t_{\alpha/2,5} = 2.5706 (\alpha = 0.05)$. 进行拒绝性检验
- Machine Learning.
In David Hemmendinger, Anthony Ralston and Edwin Reilly (Eds.), The
Encyclopedia of Computer Science, Fourth Edition, Thomson Computer Press.

Table 1.两学习器分类差别列联表

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

- 若假设A,B学习器性能相同, 则应由 $e_{01} = e_{10}$, 那么 $|e_{01} - e_{10}|$ 应服从正态分布。McNemar 检验考虑变量 $\chi^2 = \frac{(|e_{01} - e_{10}|)^2 - 1}{e_{10} + e_{01}}$. 和 $\chi_{0.975}^2(1) = pchisq(0.975, 1) = 5.023$ 进行比较, 大于它拒绝。
- 【例题】A和B两种算法对同一组数据集进行测试, A算法判出91个正例(65%), B算法判出77名正例(55%), A和B两法一致的正例56名(40%), 问哪种方法学习性能更高?

		算法 正确	A 错误	合计
算法 B	正确	56	35	91
	错误	21	28	49

- $1 - pchisq(10.73, 1) = 0.0010$ p值很小, 于是拒绝零假设, 认为两个算法性能有差异。

Friedman检验和Nemenyi检验

- 有多个数据集多个学习器进行比较时使用，对各个算法在各个数据集上对测试性能排序，对平均序值计算 χ^2 和 F 检验,并进行临界值检验。
- 假定有 B 个数据集上比较 K 个算法，令 R_j 表示第 j 个算法在 D 个数据集的秩和（这里忽略秩打结的情况），根据秩方差分析原理（非参数统计，王星2014），构造统计量

$$Q_{\chi^2} = \frac{12}{BK(K+1)} \sum_{k=1}^K R_j^2 - 3B(K+1)$$

当 $K \times B$ 比较大的时候，上述 Q 统计量服从自由度为 $K - 1$ 的 χ^2 分布，可以分析是不是几个算法性能有差异。

在一组数据集上对多个算法进行比较

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

```
> alcp=c(1,1,1,1,2,2,2,3,2.5,3,3)
> group=c(1,2,3,4,1,2,3,4,1,2,3,4)
> treat=c(1,1,1,1,2,2,2,2,3,3,3,3)
> friedman.test(alcp,treat,group)

Friedman rank sum test

data:  alcp, treat and group
Friedman chi-squared = 7.6, df = 2, p-value = 0.02237
```

- Friedman检验的后续分析1F检验：当要比较的组数很多的时候（ K 比较大）

$$F = \frac{(B-1)Q_{\chi^2}}{B(K-1) - Q_{\chi^2}},$$

F 服从自由度为 $(K-1)(B-1)$ 的 F 分布

- 【例题】 $B=4, K=3$,
 $F = (B-1) * 7.6 / (B * (K-1) - 7.6) = 57$, 显著大于临界值
 $F_{0.05}((K-1), (K-1)(B-1)) = qf(0.95, 2, 2 * 3) =$
[1]5.143253, 于是拒绝零假设，认为几组学习器之间的测试误差性能具有很大差异。

- 当Friedman检验发现多个学习器之间性能存在差异时，想知道哪些算法性能之间差异较大，需要进行平均秩之差的Nemanyi 后续检验。

-

$$CD = q_{\alpha} \sqrt{\frac{K(K+1)}{6B}}$$

- 【例题】 $CD = 2.344 * \text{sqrt}(K * (K + 1) / (6 * B))$, $CD = 1.657$
- 比较各组平均
秩 $\Delta_{1,2} = 1.125 < CD$, $\Delta_{2,3} = 0.75$, $\Delta_{1,3} = 1.775 > CD$
- 于是，有关各组算法之间性能的结论是：在显著性水平 $\alpha = 0.05$ 之下，算法A与算法C 的差异是比较大的，其他算法之间没有发现明显差异。

偏差与方差(1)

$$\begin{aligned}E(f; D) &= \mathbb{E}[(f(x; D) - y_D)^2] \\&= \mathbb{E}_D [(f(x; D) - \mathbb{E}(y|x) + \mathbb{E}(y|x) - y_D)^2] \\&= \mathbb{E}_D [(f(x; D) - \mathbb{E}(y|x))^2 + \mathbb{E} [\mathbb{E}(y|x) - y_D]^2] \\&\quad + \mathbb{E}_D [2(f(x; D) - \mathbb{E}(y|x))(\mathbb{E}(y|x) - y_D)] \\&= \mathbb{E}_D [(f(x; D) - \mathbb{E}(y|x))^2] + \mathbb{E} [\mathbb{E}(y|x) - y_D]^2 \\&= \mathbb{E}_D [(f(x; D) - \mathbb{E}(y|x))^2] + \mathbb{E} [\mathbb{E}(y|x) - y + y - y_D]^2 \\&= \mathbb{E}_D [(f(x; D) - \mathbb{E}(y|x))^2] + \mathbb{E} [(\mathbb{E}(y|x) - y)^2] + \mathbb{E} [(y - y_D)^2] \\&\quad + 2\mathbb{E}_D [(\mathbb{E}(y|x) - y)(y - y_D)] \\&= \mathbb{E}_D [(f(x; D) - \mathbb{E}(y|x))^2] + [(\mathbb{E}(y|x) - y)^2] + \mathbb{E}_D [(y - y_D)^2]\end{aligned}$$

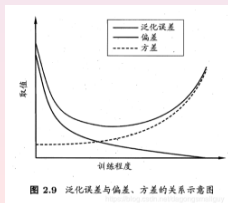
- 于是

$$E(f; D) = \text{bias}^2(x) + \text{var}(x) + \epsilon^2$$

- 泛化误差分解为偏差、方差与噪声之和。

偏差与方差(2)

- 偏差-方差分解是解释算法泛化性能的一种重要工具。偏差-方差分解试图对学习算法的期望泛化错误率进行拆解。
- 泛化误差可分解为：偏差，方差与噪声之和。
- **偏差**度量了学习算法的期望预测与真实结果的偏离程度，即刻画了学习算法本身的拟合能力。
- **方差**度量了同样大小的训练集的变动所导致的学习性能的变化，即刻画了数据扰动所造成的影响。
- **噪声**表达了在当前任务上任何学习算法所能达到的期望泛化误差的下界,即刻画了学习问题本身的难度
- 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。



- 本章主要基于NFL，学习器性能很难判断，经验误差无法代表泛化误差进行训练，否则难免不出现过度拟合；
- 需要合适的评估方法来给出泛化误差的近似值，划分训练、测试集，估计出不同的泛化估计；
- 用哪些性能指标来计算测试误差（错误率、查准率，查全率，F1,ROC,AUC,代价敏感曲线）；
- 计算完误差以后如何比较不同学习器的性能？
- 比较完性能后如何理解性能上的差异。

3.29课后作业

- 在一个钓鱼问题中（数据fishing.xls），感兴趣的是具备怎样的条件可以钓到鱼，我们的数据分成训练集和测试集（由标签变量定义），在训练集上建立了两个模型，分别获得了测试集和训练集的预测评分($E(Y|X)$ 的估计)，要求编写程序，实现以下分析内容，做出数据分析报告：
 - 请对训练数据和测试数据分别对模型1和模型2制作P-R图，ROC曲线图，用经验AUC估算面积，判断整体预测效果，如果按照实验结果给出一个符合实际的预测阈值，请说明这个阈值选取的依据；
 - 在给定的阈值上，输出混淆矩阵。
 - 如果将钓不到鱼的游客判为能钓到鱼，这个损失将会使钓鱼中心的声誉受损更大，将其设为将能钓到鱼的游客判为钓不到鱼的损失的2倍，请根据新的非平衡代价损失制作两个模型的测试数据的代价曲线，并给出比较分析结果。
- 在一个汉字识别问题中（数据手写体数字数据.xls），我们用5折交叉验证法将数据分成5折训练集和测试集，获得了3个0-1分类模型
 - 编写程序计算宏、微查全率和查准率，根据经验选择合适的阈值。
 - 编写程序求解宏期望代价曲线($Macro_{FPR}$, $Macro_{FNR}$)，微期望代价曲线($Micro_{FPR}$, $Micro_{FNR}$)并画图表示,比较三种模型的代价曲线。

3.29课后作业1

- 在一个图像去噪研究中，作者研究了四种算法 $K-SVD$, $K-SVD-N$, NLM 和 $K-SVD-N-NL$ 的实验效果，分别在5组测试图像上获得信噪比评分PSNR,该信噪比数值越高表示方法越有效，请完成以下检验：

		表1 去噪图像的 PSNR 值比较									dB
测试图像	算法	$\sigma=20$	$\sigma=30$	$\sigma=40$	$\sigma=50$	$\sigma=60$	$\sigma=70$	$\sigma=80$	$\sigma=90$	$\sigma=100$	
Lena	K-SVD	30.42	27.77	25.66	23.94	22.51	21.28	20.21	19.25	18.40	
	K-SVD-N	30.69	28.33	26.48	25.16	23.97	23.05	22.17	21.30	20.43	
	NLM	30.61	28.56	27.03	25.72	24.60	23.65	22.85	22.19	21.62	
	K-SVD-N-NL	31.57	29.96	28.57	27.51	26.75	25.91	25.47	24.92	24.36	
Barbara	K-SVD	29.25	27.01	25.04	23.34	21.90	20.70	19.67	18.76	17.96	
	K-SVD-N	29.42	27.10	25.24	23.72	22.45	21.48	20.63	19.97	19.33	
	NLM	29.26	26.98	24.76	23.43	22.39	21.56	20.89	20.34	19.88	
	K-SVD-N-NL	29.74	28.00	26.61	25.37	24.40	23.61	23.00	22.47	22.04	
Boat	K-SVD	29.09	26.89	25.00	23.40	22.04	20.88	19.85	18.95	18.13	
	K-SVD-N	29.37	27.22	25.62	24.23	23.03	22.12	21.20	20.59	19.66	
	NLM	28.50	26.41	24.98	23.93	23.08	22.36	21.75	21.21	20.75	
	K-SVD-N-NL	29.64	28.01	26.84	25.79	25.01	24.36	23.79	23.30	22.82	
Cameraman	K-SVD	26.96	25.82	24.39	22.98	21.70	20.52	19.50	18.59	17.75	
	K-SVD-N	27.52	26.20	24.74	23.45	22.29	21.34	20.46	19.70	19.15	
	NLM	28.75	26.56	24.84	23.56	22.57	21.75	21.05	20.43	19.89	
	K-SVD-N-NL	27.87	26.83	25.86	24.99	24.15	23.38	22.74	22.14	21.61	
Peppers256	K-SVD	29.18	26.94	25.01	23.34	21.92	20.71	19.59	18.64	17.79	
	K-SVD-N	29.49	27.33	25.57	24.03	22.93	21.98	21.01	20.28	19.48	
	NLM	28.95	26.57	24.74	23.19	21.89	20.84	20.01	19.33	18.78	
	K-SVD-N-NL	29.76	28.09	26.74	25.58	24.56	23.87	23.08	22.46	21.85	

- 在Cameraman数据上，分别比较四种算法，请问哪一种算法的信噪比高于20%。
- 在 $\sigma=20$ 上，四种算法是否有显著差异？