

《机器学习》读书笔记

黄奕诚

目录

1	绪论	4
1.1	引言	4
1.2	基本术语	4
1.3	假设空间	5
1.4	归纳偏好	5
1.5	发展历程	6
1.6	应用现状	6
2	模型评估与选择	7
2.1	经验误差与过拟合	7
2.2	评估方法	7
2.2.1	留出法	7
2.2.2	交叉验证法	8
2.2.3	自助法	8
2.2.4	调参与最终模型	9
2.3	性能度量	9
2.3.1	错误率与精度	9
2.3.2	查准率、查全率与 $F1$	10
2.3.3	ROC 与 AUC	12
2.3.4	代价敏感错误率与代价曲线	12
2.4	比较检验	13
2.4.1	假设检验	13
2.4.2	交叉验证 t 检验	14
2.4.3	McNemar 检验	15

2.4.4	Friedman 检验与 Nemenyi 后续检验	15
2.5	偏差与方差	16
3	线性模型	17
3.1	基本形式	17
3.2	线性回归	17
3.3	对数几率回归	19
3.4	线性判别分析	19
3.5	多分类学习	21
3.6	类别不平衡问题	22
4	决策树	23
4.1	基本流程	23
4.2	划分选择	24
4.2.1	信息增益	24
4.2.2	增益率	24
4.2.3	基尼指数	25
4.3	剪枝处理	25
4.3.1	预剪枝	26
4.3.2	后剪枝	26
5	神经网络	27
6	支持向量机	27
7	贝叶斯分类器	27
8	集成学习	27
9	聚类	27
10	降维与度量学习	27
11	特征选择与稀疏学习	27
12	计算学习理论	27

13	半监督学习	27
14	概率图模型	27
15	规则学习	27
16	强化学习	27

1 绪论

1.1 引言

- 机器学习致力于研究如何通过计算的手段，利用经验来改善系统自身的性能。
- 机器学习研究的主要内容：在计算机上从数据中产生“模型”的算法，即“学习算法”。

1.2 基本术语

- 数据集 (data set): 一组记录的集合
- 示例 (instance) / 样本 (sample): 每条记录，即关于一个事件或对象的描述
- 属性 (attribute) / 特征 (feature): 反映事件或对象在某方面的表现或性质的事项
- 属性值 (attribute value): 属性上的取值
- 属性空间 (attribute space) / 样本空间 (sample space) / 输入空间: 属性张成的空间，记为 \mathcal{X}
- 特征向量 (feature vector): 一个示例（在样本空间对应的坐标向量）
- 学习 (learning) / 训练 (training): 从数据中学得模型的过程
- 训练数据 (training data): 训练过程中使用的数据
- 训练样本 (training sample): 训练数据中的每个样本
- 训练集 (training set): 训练样本组成的集合
- 假设 (hypothesis): 对应了关于数据的某种潜在规律的学得模型
- 真实 (ground-truth): 潜在规律自身
- 学习器 (learner): 学习算法在给定数据和参数空间上的实例化
- 标记 (label): 关于示例结果的信息

- 样例 (example): 拥有标记信息的示例
- 标记空间 (label space) / 输出空间: 所有标记的集合, 记为 \mathcal{Y}
- 分类 (classification): 预测的是离散值的学习任务 (二分类 $\mathcal{Y} = \{-1, +1\}$ 或 $\{0, 1\}$; 三分类 $|\mathcal{Y}| > 2$)
- 回归 (regression): 预测的是连续值的学习任务 ($\mathcal{Y} = \mathbb{R}$)
- 测试 (testing): 使用学得模型进行预测的过程
- 测试样本 (testing sample): 被预测的样本
- 无监督学习 (unsupervised learning): 训练数据中没有标记信息的学习任务, 代表是聚类 (clustering)
- 监督学习 (supervised learning): 训练数据中具有标记信息的学习任务, 代表是分类和回归
- 泛化 (generalization) 能力: 学得模型适用于新样本的能力

1.3 假设空间

- “从样例中学习”是一个归纳的过程。
- 可以把学习过程看作一个在所有假设组成的空间中进行搜索的过程, 搜索目标是找到与训练集“匹配” (fit) 的假设。
- 假设空间可以表示为一棵属性值中通配符逐渐被具体数值取代的树。
- 可以用许多策略对假设空间进行搜索, 如自顶向下 (从一般到特殊)、自底向上 (从特殊到一般)。
- 可能有多个假设与训练集一致, 即存在着一个与训练集一致的“假设集合”, 称之为“版本空间” (version space)。

1.4 归纳偏好

- 多个与训练集一致的假设所对应的模型在面临新样本时, 可能产生不同的输出。而对于一个具体的学习算法而言, 必须要产生一个模型。此时学习算法本身的偏好会起到关键的作用。

- 归纳偏好 (inductive bias): 机器学习算法在学习过程中对某种类型假设的偏好。
- 奥卡姆剃刀 (Occam's razor): 若有多个假设与观察一致, 则选最简单的那个。【常用的、自然科学研究中最基本的原则】
- 设 f 为希望学习的真实目标函数, 则基于训练数据 X 的算法 \mathfrak{L}_a 在训练集之外的所有样本上的误差与学习算法无关, 即

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

“没有免费的午餐”定理 (NFL 定理): 所有学习算法的期望性相同。

1.5 发展历程

1. 二十世纪五十年代到七十年代初: “推理期”——赋予机器逻辑推理能力
2. 二十世纪七十年代中期开始: “知识期”
 - a. 机械学习 (信息存储与检索)
 - b. 示教学习 (从指令中学习)
 - c. 类比学习 (通过观察和发现学习)
 - d. 归纳学习 (从样例中学习)
 - 符号主义学习 (决策树、基于逻辑的学习)
 - 连接主义学习 (神经网络)
 - 统计学习 (支持向量机、核方法)
 - 深度学习

1.6 应用现状

- 计算机科学诸多分支学科领域 (如计算机视觉、自然语言处理)
- 交叉学科 (如生物信息学)
- 数据挖掘 (机器学习领域和数据库领域是数据挖掘的两大支撑)
- 人类日常生活 (天气预报、搜索引擎、自动驾驶、政治选举等)
- 促进人们理解 “人类如何学习”

2 模型评估与选择

2.1 经验误差与过拟合

- 设在 m 个样本中有 a 个样本分类错误，则错误率 (error rate) 为 $E = a/m$ ，精度 (accuracy) 为 $1 - a/m$ 。
- 误差 (error)：学习器的实际预测输出与样本的真实输出之间的差异。
训练误差 (training error) / 经验误差 (empirical error)：学习器在训练集上的误差。泛化误差 (generalization error)：学习器在新样本上的误差。想要使泛化误差最小，而新样本未知，所以努力使经验误差最小化。
- 过拟合 (overfitting)：学习器将训练样本自身的一些特点当作为所有潜在样本都会具有的一般性质。【关键障碍、无法彻底避免】
欠拟合 (underfitting)：学习器对训练样本的一般性质尚未学好。【较容易克服】
若 “ $P \neq NP$ ”，过拟合就不可避免。

2.2 评估方法

为了对学习器对泛化误差进行评估，需要使用一个测试集 (testing set) 来测试学习器对新样本的判别能力，然后以测试集上的测试误差 (testing error) 作为泛化误差的近似。

若当前只有一个包含 m 个样例的数据集

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

，则对其进行适当的处理，从中产生训练集 S 和测试集 T 。

2.2.1 留出法

直接将数据集 D 划分为两个互斥的集合，其中一个作为训练集 S ，另一个作为测试集 T ，即 $D = S \cup T, S \cap T = \emptyset$ ，在 S 上训练出模型后，用 T 来评估其测试误差，作为对泛化误差的估计。

- 划分尽可能保持数据分布的一致性，例如在分类任务中至少要保持样本的类别比例相似 (分层采样)。

- 一般采用若干次随机划分、重复进行实验评估后取平均值作为评估结果。
- S 和 D 大小权衡没有完美的解决方案，常见做法是 $2/3 \sim 4/5$ 的训练样本比例。

2.2.2 交叉验证法

将数据集 D 划分为 k 个大小相似的互斥子集，即

$$D = D_1 \cup D_2 \cup \dots \cup D_k, D_i \cap D_j = \emptyset (i \neq j)$$

每个子集 D_i 都尽可能保持数据分布的一致性（分层抽样）。然后从中选取 $k-1$ 个子集为训练集，剩下一个子集为测试集。可进行 k 次训练和测试，最终返回 k 个测试结果的均值。也称为“ k 折交叉验证”（ k -fold cross validation）。

- k 最常用的取值是 10，常用的还有 5、20 等。
- 留一法（Leave-One-Out）不受随机样本划分的影响，评估结果比较准确，但计算开销大。

2.2.3 自助法

以自助采样法（bootstrap sampling）为基础，给定包含 m 个样本的数据集 D ，每次随机从 D 中挑选一个样本，将其拷贝放入 D' ，再将该样本放回初始数据集 D 中。这个过程重复执行 m 次后，得到了包含 m 个样本的数据集 D' 。此时将 D' 用作训练集， $D \setminus D'$ 用作测试集。

- D 有约 36.8% 的样本未出现在采样数据集 D' 中。
- 亦称为“包外估计”（out-of-bag estimate）。
- 自助法在数据集较小、难以有效划分训练/测试集时很有用，且能从初始数据集中产生多个不同的训练集。
- 因为自助法产生的数据集改变了初始数据集的分布，会引入估计偏差。

2.2.4 调参与最终模型

- 常用的调参做法：对每个参数选定一个范围和变化步长，进行计算开销和性能估计之间的折中。
- 在模型选择完成后，学习算法和参数配置已选定，此时用数据集 D 重新训练模型，使用所有 m 个样本，得到最终提交给用户的模型。

2.3 性能度量

给定样例集

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$$

其中 y_i 是示例 \mathbf{x}_i 的真实标记。要评估学习器 f 的性能，即把学习器预测结果 $f(\mathbf{x})$ 与真实标记 y 进行比较。

回归任务中最常用的性能度量：“均方误差” (mean squared error)

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

更一般地，对于数据分布 \mathcal{D} 和概率密度函数 $p(\cdot)$ ，均方误差可描述为

$$E(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} (f(\mathbf{x}) - y)^2 p(\mathbf{x}) d\mathbf{x}$$

对于分类任务——

2.3.1 错误率与精度

- 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

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

- 对于数据分布 \mathcal{D} 和概率密度函数 $p(\cdot)$, 错误率

$$E(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}(f(\mathbf{x}) \neq y) p(\mathbf{x}) d\mathbf{x}$$

精度

$$\text{acc}(f; \mathcal{D}) = \int_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}(f(\mathbf{x}) = y) p(\mathbf{x}) d\mathbf{x} = 1 - E(f; \mathcal{D})$$

2.3.2 查准率、查全率与 $F1$

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 (precision)

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

查全率 (recall)

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

- 平衡点 (Break-Even Point, BEP): $R = P$ 时的取值, 数值越高可以认为学习器越优。

- $F1$ 度量

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

实际上 $F1$ 是 R 和 P 的调和平均

$$\frac{1}{F1} = \frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)$$

- F_β 度量: 考虑 R 与 P 的不同偏好, 设 β 为查全率 R 对查准率 P 的相对重要性, 则

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

实际上 F_β 是加权调和平均

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

- 宏 $F1$: 在各混淆矩阵上分别计算出各自的 (P_i, R_i) , 再计算平均值:

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

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

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

- 微 $F1$: 先将各混淆矩阵的对应元素进行平均得到四个指标, 再基于这些平均值计算 $F1$:

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

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

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

△ 混淆矩阵介绍: 每一列代表了预测类别, 每一列的总数表示预测为该类别的数据的数目; 每一行代表了数据的真实归属类别, 每一行的数据总数表示该类别的数据实例的数目。例如共有 150 个样本数据, 预测为 1、2、3 类各 50 个, 分类结束后得到的混淆矩阵为

		预测		
		类 1	类 2	类 3
实际	类 1	43	2	0
	类 2	5	45	1
	类 3	2	3	49

2.3.3 ROC 与 AUC

ROC 全称是“受试者工作特征” (Receiver Operating Characteristic) 曲线。横轴为“假正例率” (FPR), 纵轴为“真正例率” (TPR)。

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

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

- 现实任务中 ROC 曲线的绘制方法: 给定 m^+ 个正例和 m^- 个反例, 根据学习器预测结果对样例进行排序, 然后把分类阈值设为最大, 此时 FPR 和 TPR 都为 0. 在坐标 $(0,0)$ 处标记一个点, 然后将分类阈值依次设为每个样例的预测值。设当前一个标记点坐标为 (x,y) , 若当前为真正例, 则对应标记点坐标为 $(x, y + \frac{1}{m^+})$; 若当前为假正例, 则对应标记点坐标为 $(x + \frac{1}{m^-}, y)$, 然后用线段连接相邻点即得。
- AUC (Area Under ROC Curve) 即为 ROC 曲线下各部分的面积之和。设 ROC 曲线是由坐标为 $\{(x_i, y_i) | 1 \leq i \leq m\}$ 的点按序连接而成, 则 AUC 可估算为

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

- 给定 m^+ 个正例和 m^- 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

$$\ell_{rank} = \frac{1}{m^+ m^-} \sum_{\mathbf{x}^+ \in D^+} \sum_{\mathbf{x}^- \in D^-} \left(\mathbb{I}(f(\mathbf{x}^+) < f(\mathbf{x}^-)) + \frac{1}{2} \mathbb{I}(f(\mathbf{x}^+) = f(\mathbf{x}^-)) \right)$$

它对应 ROC 曲线之上的面积, 有

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

2.3.4 代价敏感错误率与代价曲线

- 不同类型的错误可能造成不同损失, 所以为错误赋予“非均等代价” (unequal cost)。

- 以二分类为例，可以设定一个“代价矩阵”，如下表所示。

真实类别	预测类别	
	第 0 类	第 1 类
第 0 类	0	$cost_{01}$
第 1 类	$cost_{10}$	0

- “代价敏感” (cost-sensitive) 错误率

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

- 在非均等代价下，“代价曲线” (cost curve) 可以刻画期望总体代价。设 p 是样例为正例的概率。横轴为正例概率代价

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

纵轴为取值为 $[0, 1]$ 的归一化代价

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

- 代价曲线的绘制方法：将 ROC 曲线上的每一点转化为代价平面上的的一条线段，取所有线段的下界，围成的面积即为所有条件下学习器的期望总体代价。

2.4 比较检验

本节默认以错误率 ϵ 为性能度量。

2.4.1 假设检验

- 设一个学习器的泛化错误率为 ϵ ，在 m 个样本中的测试错误率为 $\hat{\epsilon}$ ，则其被测得测试错误率为 $\hat{\epsilon}$ 的概率为

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

它在 $\epsilon = \hat{\epsilon}$ 时最大。

- 二项检验：假设 $\epsilon \leq \epsilon_0$ ，则在 $1 - \alpha$ 的概率内所能观测到的最大错误率为

$$\sum_{i=\epsilon_0 \times m+1}^m \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i} < \alpha$$

$$\bar{\epsilon} = \max \epsilon$$

若测试错误率 $\hat{\epsilon}$ 小于临界值 $\bar{\epsilon}$ ，则能以 $1 - \alpha$ 的置信度认为学习器的泛化错误率不大于 ϵ_0 ，否则假设被拒绝。

- t 检验：若得到了 k 个测试错误率 $\hat{\epsilon}_1, \hat{\epsilon}_2, \dots, \hat{\epsilon}_k$ ，则平均测试错误率 μ 和方差 σ^2 为

$$\mu = \frac{1}{k} \sum_{i=1}^k \hat{\epsilon}_i$$

$$\sigma^2 = \frac{1}{k-1} \sum_{i=1}^k (\hat{\epsilon}_i - \mu)^2$$

它们可看作泛化错误率 ϵ_0 的独立采样，则变量

$$\tau_t = \frac{\sqrt{k}(\mu - \epsilon_0)}{\sigma}$$

服从自由度为 $k-1$ 的 t 分布。若 $|\mu - \epsilon_0|$ 位于 $[t_{-\alpha/2}, t_{\alpha/2}]$ 内，则接受假设 $\mu = \epsilon_0$ ，否则拒绝该假设。

2.4.2 交叉验证 t 检验

- k 折交叉验证“成对 t 检验”：对每对结果求差 $\Delta_i = \epsilon_i^A - \epsilon_i^B$ ，若两个学习器性能相同，则差值均值为 0。做 t 检验，在显著度 α 下，若

$$\tau_t = \left| \frac{\sqrt{k}\mu}{\sigma} \right| < t_{\alpha/2, k-1}$$

则接受假设。其中 $t_{\alpha/2, k-1}$ 指自由度为 $k-1$ 的 t 分布上尾部累积分布为 $\alpha/2$ 的临界值。

- 考虑到交叉验证法等实验估计方法，不同轮次的训练集会有一定程度的重叠，导致测试错误率并不独立。故可采用 5×2 交叉验证法（5 次 2 折交叉验证）。每次 2 折交叉验证之前随机将数据打乱，使得 5 次交

叉验证中的数据划分不重复。设 Δ_i^k 表示第 i 次第 k 上的差值。

$$\mu = 0.5(\Delta_1^1 + \Delta_1^2)$$

$$\sigma_i^2 = \left(\Delta_i^1 - \frac{\Delta_i^1 + \Delta_i^2}{2}\right)^2 + \left(\Delta_i^2 - \frac{\Delta_i^1 + \Delta_i^2}{2}\right)^2$$

变量

$$\tau_t = \frac{\mu}{\sqrt{0.2 \sum_{i=1}^5 \sigma_i^2}}$$

服从自由度为 5 的 t 分布，其双边检验的临界值为 $t_{\alpha/2,5}$ 。

2.4.3 McNemar 检验

对于二分类问题，可统计两个学习器 A 和 B 的分类结果样本数差别，列出“列联表” (contingency table)

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

假设两学习器性能相同，则 $e_{01} = e_{10}$ ，于是 $|e_{01} - e_{10}|$ 服从正态分布，变量

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

服从自由度为 1 的 χ^2 分布，若其小于临界值 χ_α^2 则接受假设，否则拒绝假设，较小者性能更优。

2.4.4 Friedman 检验与 Nemenyi 后续检验

- 在多个数据集上比较算法。
- 算法排序：使用留出法或交叉验证法得到每个算法在每个数据集上的测试结果，然后在每个数据集上根据测试性能由好到坏排序，序值从 1 递增，若相同则平分序值。
- “原始 Friedman 检验”：假定在 N 个数据集上比较 k 个算法，令 r_i 表示第 i 个算法的平均序值，暂不考虑平分序值，则 r_i 均值为 $(k+1)/2$ ，

方差为 $(k^2 - 1)/12$ 。变量

$$\tau_{\chi^2} = \frac{k-1}{k} \cdot \frac{12N}{k^2-1} \sum_{i=1}^k \left(r_i - \frac{k+1}{2}\right)^2 = \frac{12N}{k(k+1)} \left(\sum_{i=1}^k r_i^2 - \frac{k(k+1)^2}{4}\right)$$

当 k 和 N 都较大时服从自由度为 $k-1$ 的 χ^2 分布。

- Friedman 检验：变量

$$\tau_F = \frac{(N-1)\tau_{\chi^2}}{N(k-1) - \tau_{\chi^2}}$$

服从自由度为 $k-1$ 和 $(k-1)(N-1)$ 的 F 分布。

- Nemenyi 检验：若“所有算法的性能相同”这一假设被拒绝，此时计算出平均序值差别的临界值域

$$CD = q_\alpha \sqrt{\frac{k(k+1)}{6N}}$$

若某两个算法的平均序值之差超出了 CD ，则以相应的置信度拒绝“这两个算法性能相同”这一假设。

- Friedman 检验图：横轴为平均序值，纵轴为各个算法，对每个算法以一个圆点表示平均序值，以圆点为中心的横线段表示临界值域的大小。若两个算法的横线段有交叠，则说明它们没有显著区别，否则可以进行显著比较。

2.5 偏差与方差

- 偏差-方差分解：对学习算法的期望泛化错误率进行拆解。
- 学习算法的期望预测

$$\bar{f}(\mathbf{x}) = \mathbb{E}_D[f(\mathbf{x}; D)]$$

- 使用样本数相同的不同训练集产生的方差

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

- 噪声

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

- 偏差（期望输出与真实标记的差别）

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

- 假定 $\mathbb{E}_D[y_D - y] = 0$ ，则可通过多项式展开得到

$$E(f; D) = \mathbb{E}_D[(f(\mathbf{x}; D) - y_D)^2] = bias^2(\mathbf{x}) + var(\mathbf{x}) + \varepsilon^2$$

泛化误差可分解为偏差、方差与噪声之和。

- 偏差：学习算法本身的拟合能力；方差：数据的充分性；噪声：学习问题本身的难度
- 偏差-方差窘境（bias-variance dilemma）：训练不足时偏差主导，训练加深时方差主导，训练充足时容易发生过拟合。

3 线性模型

3.1 基本形式

示例 $\mathbf{x} = (x_1; x_2; \dots; x_d)$ ，其中 x_i 是 \mathbf{x} 在第 i 个属性上的取值，则线性模型

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b = \mathbf{w}^T \mathbf{x} + b$$

3.2 线性回归

- 一元线性回归的目标

$$f(x_i) = wx_i + b, \text{ 使得 } f(x_i) \simeq y_i$$

其中

$$(w^*, b^*) = \arg_{(w, b)} \min \sum_{i=1}^m (y_i - wx_i - b)^2$$

利用“最小二乘法”的最小二乘“参数估计”将 $E_{(w,b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$ 对 w 和 b 分别求导并令为零即可得到 w, b 最优解的闭式解

$$w = \frac{\sum_{i=1}^m y_i(x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m}(\sum_{i=1}^m x_i)^2} \quad b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i)$$

- 多元线性回归的目标

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b, \text{ 使得 } f(\mathbf{x}_i) \simeq y_i$$

设 $\hat{\mathbf{w}} = (\mathbf{w}; b)$, 将数据集 D 表示为一个 $m \times (d+1)$ 大小的矩阵 \mathbf{X}

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} & 1 \\ x_{21} & x_{22} & \cdots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix}$$

其中

$$\hat{\mathbf{w}}^* = \arg_{\hat{\mathbf{w}}} \min (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

令 $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$, 对 $\hat{\mathbf{w}}$ 求导并令为零即 $2\mathbf{X}^T(\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) = 0$ 即可。若 $\mathbf{X}^T\mathbf{X}$ 正定 (满秩), 则求出 $\hat{\mathbf{w}}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$, 此时

$$f(\hat{\mathbf{x}}_i) = \hat{\mathbf{x}}_i^T (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

若 $\mathbf{X}^T\mathbf{X}$ 不满秩, 则可能有多解, 通过引入正则化由归纳偏好决定。

- 广义线性模型: 考虑单调可微函数 $g(\cdot)$

$$y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$$

其中 $g(\cdot)$ 称为“联系函数”。

3.3 对数几率回归

- 单位阶跃函数：对于二分类问题将线性回归模型的实值转化为 0/1 值。其中预测值为临界值零可任意判别

$$y = \begin{cases} 0, & z < 0; \\ 0.5, & z = 0; \\ 1, & z > 0; \end{cases}$$

- 对数几率函数

$$y = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

可化为

$$\ln \frac{y}{1-y} = \mathbf{w}^T \mathbf{x} + b$$

其中 y 可视为 \mathbf{x} 作为正例的可能性， $1-y$ 可视为其作为反例的可能性。

- 对数几率回归的 \mathbf{w}, b 估计方法：极大似然法。【TODO：细节待学完 7.2 节极大似然法及梯度下降法再补充】

3.4 线性判别分析

- 思想：设法将样例投影到一条直线上，使得同类样例的投影尽可能接近、异类样例的投影尽可能远离。对于新样本根据其投影到这条直线的投影点位置进行分类。

△ L2 范数（例如欧氏距离）

$$\| \mathbf{x} \|_2 = \sqrt{\sum_{i=1}^k |x_i|^2}$$

- 目的：使同类样例投影点的协方差 ($\mathbf{w}^T \Sigma_0 \mathbf{w} + \mathbf{w}^T \Sigma_1 \mathbf{w}$) 尽可能小，使类中心之间的距离 ($\| \mathbf{w}^T \boldsymbol{\mu}_0 - \mathbf{w}^T \boldsymbol{\mu}_1 \|_2^2$) 尽可能大。

- 类内散度矩阵

$$\mathbf{S}_w = \Sigma_0 + \Sigma_1 = \sum_{\mathbf{x} \in X_0} (\mathbf{x} - \boldsymbol{\mu}_0)(\mathbf{x} - \boldsymbol{\mu}_0)^T + \sum_{\mathbf{x} \in X_1} (\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{x} - \boldsymbol{\mu}_1)^T$$

- 类间散度矩阵

$$\mathbf{S}_b = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T$$

- 欲最大化目标 (\mathbf{S}_b 与 \mathbf{S}_w 的广义 Rayleigh 商)

$$J = \frac{\|\mathbf{w}^T \boldsymbol{\mu}_0 - \mathbf{w}^T \boldsymbol{\mu}_1\|_2^2}{\mathbf{w}^T \Sigma_0 \mathbf{w} + \mathbf{w}^T \Sigma_1 \mathbf{w}} = \frac{\mathbf{w}^T (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \mathbf{w}}{\mathbf{w}^T (\Sigma_0 + \Sigma_1) \mathbf{w}} = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

- 确定 \mathbf{w} 的方法: J 中的解与 \mathbf{w} 的长度无关, 故令 $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$, 转化为: 已知 $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$ 求 $-\mathbf{w}^T \mathbf{S}_b \mathbf{w}$ 的最小值。由拉格朗日乘子法, 即 $\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$ 。又 $\mathbf{S}_b \mathbf{w}$ 方向恒为 $\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$, 令 $\mathbf{S}_b \mathbf{w} = \lambda(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$, 得

$$\mathbf{w} = \mathbf{S}_w^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$$

对 \mathbf{S}_w 进行奇异值分解 $\mathbf{S}_w = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$, 由 $\mathbf{S}_w^{-1} = \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^{-1}$ 得到 \mathbf{w} 。

△ 上例用拉格朗日乘子法的计算细节: 目标函数为 $f(\mathbf{x}) = -\mathbf{w}^T \mathbf{S}_b \mathbf{w}$, 约束方程 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1 = 0$ 。等价于由方程 $g(\mathbf{x}) = 0$ 确定的 $d-1$ 维曲面上寻找能使 $f(\mathbf{x})$ 最小化的点, 满足

(1) 约束曲面上任意点 \mathbf{x} 的梯度 $\nabla g(\mathbf{x})$ 正交于约束曲面

(2) 在最优点 \mathbf{x}^* , $f(\mathbf{x})$ 在该点的梯度 $\nabla f(\mathbf{x}^*)$ 正交于约束曲面

于是在最优点 \mathbf{x}^* , 梯度 $\nabla g(\mathbf{x})$ 和 $\nabla f(\mathbf{x})$ 的方向必相同或相反, 即存在 $\lambda \neq 0$ 使得

$$\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$$

即

$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

- LDA 推广到多分类任务。假设存在 N 个类, 且第 i 类示例数为 m_i 。

设 μ 是所有示例的均值向量。全局散度矩阵

$$S_t = S_b + S_w = \sum_{i=1}^m (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T$$

类内散度矩阵

$$S_w = \sum_{i=1}^N S_{w_i} = \sum_{i=1}^N \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T$$

类间散度矩阵

$$S_b = S_t - S_w = \sum_{i=1}^N m_i (\mu_i - \mu)(\mu_i - \mu)^T$$

常用实现的优化目标

$$\max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}$$

可以通过广义特征值 $\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$ 求解, \mathbf{W} 的闭式解为 $\mathbf{S}_w^{-1} \mathbf{S}_b$ 的 d' 个最大非零广义特征值对应的特征向量组成的矩阵, 有 $d' \leq N - 1$, 实现了降维。

3.5 多分类学习

- 基本思路: 将多分类任务拆为若干个二分类任务求解, 最经典的拆分策略有三种。
- 一对一 (OvO): 将 N 个类别两两配对, 产生 $N(N-1)/2$ 个二分类任务。最终把预测得最多的类别作为分类结果。
- 一对其余 (OvR): 每次将一个类的样例作为正例, 所有其它类的样例作为反例, 产生 N 个分类任务。最终若仅有一个分类器预测为正类, 则对应的类别标记为最终分类结果, 否则考虑各分类器的置信区间, 选择置信度最大的类别标记作为分类结果。
- 多对多 (MvM): 每次将若干个类作为正类, 若干个其它类作为反类, 可用纠错输出码 (ECOC) 技术。
- ECOC 工作过程:

- (1) 编码（类别划分）：对 N 个类别做 M 次划分，每次划分将一部分类别作为正类、一部分类别作为反类，产生 M 个分类器。
- (2) 解码（距离比较）：用 M 个分类器对测试样本进行预测，这些预测标记组成一个编码，计算其与各个类别各自编码的距离，返回距离最小的类别作为分类结果。

常用的编码矩阵为二元码和三元码（有停用类）。

- 任何两个类别之间的编码距离越远，纠错能力越强，而码长的增加会增大确定最优编码的难度。

3.6 类别不平衡问题

- 类别不平衡（class-imbalance）：分类任务中不同类别的训练样例数目差别很大。以下为类别不平衡学习的策略。
- 再缩放（rescaling）：假设“训练集是真实样本总体的无偏采样”，令

$$\frac{y'}{1-y'} = \frac{y}{1-y} \times \frac{m^-}{m^+}$$

也是代价敏感学习，其中的 m^-/m^+ 可用 $cost^+/cost^-$ 代替。

- 欠采样（undersampling）/下采样（downsampling）：去除一些正例（反例）使得正、反例数目接近。可用 EasyEnsemble 算法将反例划分为若干个集合供不同学习器使用，全局上不会丢失重要信息。
- 过采样（oversampling）/上采样（upsampling）：增加一些正例（反例）使得正、反例数目接近。可用 SMOTE 算法对训练集中的正例进行插值。
- 阈值移动（threshold-moving）：基于原始训练集学习，在用训练好的分类器进行预测时将“再缩放”中的式子嵌入到决策过程中。

4 决策树

4.1 基本流程

- 决策树的性质：包含一个根结点、若干个内部结点、若干个叶结点，叶结点对应于决策结果，其他每个结点对应于一个属性测试，根结点包含样本全集。从根结点到每个叶结点的路径对应了一个判定测试序列。
- 决策树的目的：产生一棵泛化能力强的决策树。
- 决策树的基本流程（分治策略）

Algorithm 1 TreeGenerate(D, A)

Input: 训练集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$;

属性集 $A = \{a_1, a_2, \dots, a_d\}$

Output: 以 node 为根结点的一棵决策树

```

1: 生成结点 node;
2: if  $D$  中样本全属于同一类别  $C$  then
3:   将 node 标记为  $C$  类叶结点;
4:   return
5: end if
6: if  $A = \emptyset$  or  $D$  中样本在  $A$  上取值相同 then
7:   将 node 标记为叶结点，其类别标记为  $D$  中样本数最多的类;
8:   return
9: end if
10: 从  $A$  中选择最优划分属性  $a_*$ ;
11: for  $a_*$  的每一个值  $a_*^v$  do
12:   为 node 生成一个分支;
13:   令  $D_v$  表示  $D$  中在  $a_*$  上取值为  $a_*^v$  的样本子集;
14:   if  $D_v$  为空 then
15:     将分支结点标记为叶结点，其类别标记为  $D$  中样本数最多的类;
16:     return
17:   else
18:     以 TreeGenerate( $D_v, A \setminus \{a_*\}$ ) 为分支结点
19:   end if
20: end for

```

4.2 划分选择

目的：使决策树的分支结点所包含的样本尽可能属于同一类别，结点的纯度越来越高。

4.2.1 信息增益

- ID3 决策树算法以信息增益为准则选择划分属性。
- 信息熵（设 D 中共有 \mathcal{Y} 类样本，第 k 类样本所占的比例为 p_k ，值越小则 D 的纯度越高）

$$\text{Ent}(D) = - \sum_{k=1}^{|\mathcal{Y}|} p_k \log_2 p_k$$

- 设离散属性 a 的取值集合为 $\{a^1, a^2, \dots, a_V\}$ ，用 a 对样本集合 D 进行划分，产生 V 个分支结点，第 v 个分支结点包含了 D 中所有在属性 a 上取值为 a^v 的样本，记为 D^v ，于是信息增益（值越大则纯度提升越大）：

$$\text{Gain}(D, a) = \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

- 最终选择的属性

$$a_* = \arg_{a \in A} \max \text{Gain}(D, a)$$

- 缺点：对可取值数目较多对属性有所偏好。

4.2.2 增益率

- C4.5 决策树算法以增益率为准则选择划分属性。
- 增益率

$$\text{Gain}_{\text{ratio}}(D, a) = \frac{\text{Gain}(D, a)}{\text{IV}(a)}$$

其中 $\text{IV}(a)$ 为属性 a 的“固有值”，可能取值数目越多则通常越大。

$$\text{IV}(a) = - \sum_{v=1}^V \frac{|D^v|}{|D|} \log_2 \frac{|D^v|}{|D|}$$

- 缺点：对可取值数目较少的属性有所偏好，所以选择候选划分属性时使用了一个启发式算法：先从候选划分属性中找出信息增益高于平均水平的属性，再从中选择增益率最高的。

4.2.3 基尼指数

- CART 决策树算法以基尼指数 (Gini index) 为准则选择划分属性。
- 基尼值 (度量数据集 D 的纯度)，反映了从 D 中随机抽取两个样本，其类别标记不一致的概率，值越小则纯度越高。

$$\text{Gini}(D) = \sum_{k=1}^{|\mathcal{Y}|} \sum_{k' \neq k} p_k p_{k'} = 1 - \sum_{k=1}^{|\mathcal{Y}|} p_k^2$$

- 属性 a 的基尼指数

$$\text{Gini}_{\text{index}}(D, a) = \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Gini}(D^v)$$

- 最终选择属性

$$a_* = \arg_{a \in A} \min \text{Gini}_{\text{index}}(D, a)$$

4.3 剪枝处理

- 目的：去掉一些分支来降低过拟合的风险。
- 基本策略之“预剪枝” (prepruning)：在决策树生成过程中对每个结点在划分前先进行估计，若当前结点的划分不能带来决策树泛化性能提升，则停止划分并将当前结点标记为叶结点。
- 基本策略之“后剪枝” (postpruning)：先从训练集生成一棵完整的决策树，然后自底向上对非叶结点考察，若将该结点的子树替换为叶结点能提升决策树的泛化性能，则将该子树替换为叶结点。
- 泛化性能评估方法用 2.2 节的性能评估方法即可，简单地，可以使用留出法并计算错误率 (精度)。

4.3.1 预剪枝

- 降低了过拟合的风险。
- 显著减少决策树训练和测试时间开销。
- 可能因为本身的贪心禁止了那些后续划分使泛化性能显著提升的分支的展开，带来欠拟合的风险。

4.3.2 后剪枝

- 通常比预剪枝决策树保留了更多的分支。
- 欠拟合风险很小，泛化性能往往优于预剪枝决策树。
- 训练时间比未剪枝或预剪枝决策树大得多。

4.4 连续与缺失值

5 神经网络

6 支持向量机

7 贝叶斯分类器

8 集成学习

9 聚类

10 降维与度量学习

11 特征选择与稀疏学习

12 计算学习理论

13 半监督学习

14 概率图模型

15 规则学习

16 强化学习