

模式识别与统计学习

[本页PDF](#)

以下笔记为个人理解，涵盖了这门课的学习重点，若和教材有出入请以教材为主

从EM算法开始属于无监督学习

统计学习概论

常见的定义

统计学习的定义：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和分析

统计学习的对象：数据 - 关于数据的基本假设是**同类数据具有一定的统计规律性**

统计学习的目的：对数据进行预测和分析

统计学习的三要素：方法=模型+策略+算法

期望风险：模型在数据联合分布上的期望损失，衡量模型泛化能力，但实际应用难以求解，使用结构风险代替

经验风险：模型在已知训练集的平均损失

结构风险：经验风险 + 正则化项

泛化误差：衡量模型的泛化能力，等于近似误差 + 估计误差 - 近似误差：模型简单引起的误差，可以看成偏差 - 估计误差：数据太少引起的误差，可以看成方差

欠拟合：高偏差低方差，模型在训练集，验证集，测试集的泛化能力都很差 - 解决办法：增加模型复杂度，增加数据，减少正则化项等

过拟合：低偏差高方差，模型在已知的数据集上预测的很好，但在未知的数据集上预测的很差 - 解决办法：降低模型复杂度，增加数据，增加正则化项，早停等

生成式模型：学习的是数据的联合分布

判别式模型：学习的是数据的条件概率分布或者决策函数

概率模型：学习的是条件概率分布或者联合分布 - 朴素贝叶斯，高斯混合模型

非概率模型：学习的是决策函数 - 有感知机，knn，kmeans，SVM等

注意：逻辑斯特回归既是概率模型也是非概率模型

参数化模型：对数据有基本假设，用带参函数建模，参数量不变 - 有感知机，朴素贝叶斯，逻辑斯特回归等

非参数化模型：对数据没有基本假设，无学习参数 - 有SVM，knn，决策树，Adaboost等

统计学习可以分为：监督学习、无监督学习、强化学习 - 监督学习：从有标注数据中学习预测模型的机器学习问题 - 无监督学习：从无标注数据中学习预测模型的机器学习问题

混淆矩阵

	P	N
T	TP	FP
F	FN	TN

P和N代表样本真实标签，T和F代表样本的预测标签 - TP：T表示样本的预测结果是正确的，P表示样本被预测为正例

要会计算的指标 精确率

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

召回率

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

F1值

$$\frac{2}{F1} = \frac{1}{Precision} + \frac{1}{Recall}$$

感知机

学习的是 R^n 空间的超平面方程

$$w \cdot x + b = 0$$

其中, $w = [w^{(1)}, w^{(2)}, \dots, w^{(n)}], x = [x^{(1)}, x^{(2)}, \dots, x^{(n)}]^T$

点到超平面的几何距离

$$d = \frac{|w \cdot x + b|}{||w||}$$

$$\begin{cases} w \cdot x + b > 0 & \text{在正面} \\ w \cdot x + b = 0 & \text{在超平面上} \\ w \cdot x + b < 0 & \text{在反面} \end{cases}$$

感知机其实就是线性二分类模型（判别式模型）

$$f(x) = sign(w \cdot x + b)$$

$$sign(x) = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases}$$

损失函数

- 只考虑了误分类点， M 是误分类点的集合
- $||w||$ 不被考虑，能正确分类即可，至于真正距离超平面多远，模型并不关心

$$L(w, b) = - \sum_{x_i \in M} y_i (w \cdot x_i + b)$$

随机梯度下降：每次只使用一个样本进行训练，更新模型参数

$$\min_{w,b} L(w,b) = - \sum_{x_i \in M} y_i (w \cdot x_i + b)$$

梯度

$$\nabla_w L(w,b) = - \sum_{x_i \in M} y_i x_i \quad \nabla_b L(w,b) = - \sum_{x_i \in M} y_i$$

若 $y_i(w \cdot x_i + b) \leq 0$ ，则对参数进行更新，其中 η 为学习率

$$w \leftarrow w + \eta y_i x_i \quad b \leftarrow b + \eta y_i$$

对偶形式

将 w 和 b 表示为实例 x_i 和 y_i 的线性组合的形式，求解其系数

$$w = w_0 + \sum_{i=1}^N \alpha_i y_i x_i \quad b = b_0 + \sum_{i=1}^N \alpha_i y_i$$

其中： $\alpha_i = n_i \eta$ ， n_i 为 x_i 被误分类的次数

Gram矩阵 使用Gram矩阵存储内积可以加速计算

$$G = [x_i \cdot x_j]_{N \times N} = \begin{bmatrix} x_1 \cdot x_1 & x_1 \cdot x_2 & \cdots & x_1 \cdot x_N \\ x_2 \cdot x_1 & x_2 \cdot x_2 & \cdots & x_2 \cdot x_N \\ \vdots & \vdots & \ddots & \vdots \\ x_N \cdot x_1 & x_N \cdot x_2 & \cdots & x_N \cdot x_N \end{bmatrix}$$

感知机模型

$$f(x) = \text{sign}(\sum_{j=1}^N \alpha_j y_j x_j \cdot x + b)$$

其中： $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_N)^T$

如果 $y_i(\sum_{j=1}^N \alpha_j y_j x_j \cdot x + b) \leq 0$ ，进行参数更新

$$\alpha_j \leftarrow \alpha_j + \eta \quad b \leftarrow b + \eta y_i$$

我个人的理解：

如果样本 x_i 被错误分类了，对 α_i 进行更新，同时使用更新的 α 进行下次样本判别

knn

多分类回归模型，时间和空间复杂度高，分类边界的特点是**不规则曲线**

决策函数：决定了样本被分成什么类别

- 当 $y_i = c_j$ 时， $I(y_i = c_j)$ 为1，否则为0
- $N_k(x)$ 为数据集中涵盖与 x 最近 k 个点的领域

- 统计这k个样本的类别，选取个数最多的类别作为预测结果

$$y = \operatorname{argmax} \sum_{x_i \in N_k(x)} I(y_i = c_j)$$

距离度量

- L_1 为曼哈顿距离
- L_2 为欧氏距离
- L_∞ 为切比雪夫距离

$$L_k(x, y) = \left(\sum_{i=1}^d |x_i - y_i|^k \right)^{\frac{1}{k}}$$

k值的选取对模型的影响

- k值太小：模型对噪声敏感，整体模型复杂，易过拟合，近似误差减小，估计误差增大
- k值太大：整体模型简单，易欠拟合，估计误差减小，近似误差增大

当 $k = 1$ 时模型的训练误差是0，当 $k > 1$ 的时候不一定

朴素贝叶斯

贝叶斯定理

$$\begin{aligned} p(x|y) &= \frac{p(y|x)p(x)}{p(y)} \\ p(x|y, z) &= \frac{p(x|z)p(y|x, z)}{p(y, z)} \\ &= \frac{p(x, y, z)}{p(y, z)} \\ &= \frac{p(z)p(x|z)p(y|x, z)}{p(z)p(y|z)} \end{aligned}$$

基本思想

- 已知类条件概率密度参数表达式和先验概率
- 利用贝叶斯公式转换成后验概率
- 根据后验概率大小进行决策分类

通过训练数据集学习联合概率 $P(X, Y)$ ，具体可以学 $P(Y = c_k)$ 与 $P(X = x|Y = c_k)$ ，上面这俩就是要学习的参数

先来认识几个概率

$$\underbrace{P(Y = y|X = x)}_{\text{后验概率}} = \frac{\underbrace{P(X = x|Y = y)}_{\text{似然概率}} \underbrace{P(Y = y)}_{\text{先验概率}}}{\underbrace{P(X = x)}_{\text{证据概率}}}$$

条件独立性假设：假设特征之间是相互独立的

$$\begin{aligned}P(X = x|Y = c_k) &= P(X^{(1)} = x^{(1)}, \dots, X^{(n)} = x^{(n)}|Y = c_k) \\&= \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)\end{aligned}$$

具体的推导

$$\begin{aligned}P(Y = c_k|X = x) &= \frac{P(X = x|Y = c_k)P(Y = c_k)}{\sum_k P(X = x|Y = c_k)P(Y = c_k)} \\&= \frac{P(Y = c_k) \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)}{\sum_k P(X = x|Y = c_k)P(Y = c_k)}\end{aligned}$$

分母是常数，可忽略，最后的目标为最大化 y ，即

$$y = f(x) = \arg\max_{c_k} P(Y = c_k) \prod_{j=1}^n P(X^{(j)} = x^{(j)}|Y = c_k)$$

为何最大化后验概率？下文后验概率最大化的含义给出了证明

用人话说就是根据训练集统计出 $P(Y = c_k)$ 和 $P(X = x|Y = c_k)$ ，然后根据测试样本的特征找到对应的 $P(X = x|Y = c_k)$ ，最后算出不同类别的 $P(Y = c_k)$ 和 $P(X = x|Y = c_k)$ 的乘积，选择结果最大的那个作为测试样本的类别

参数估计

极大似然估计

使用样本频率估计概率

$$\begin{aligned}P(Y = c_k) &= \frac{\sum_{i=1}^N I(y_i = c_k)}{N} \\P(X^{(j)} = a_{jl}|Y = c_k) &= \frac{\sum_{i=1}^N I(X_i^{(j)} = a_{jl}, y_j = c_k)}{\sum_{i=1}^N I(y_i = c_k)}\end{aligned}$$

其中： I 是指示函数

贝叶斯估计

可以防止某些特征的取值个数为0导致似然概率为0，加上一个平滑项使最后的乘积结果不为0

$$\begin{aligned}P_\lambda(Y = c_k) &= \frac{\sum_{i=1}^N I(y_i = c_k) + \lambda}{N + K\lambda} \\P_\lambda(X^{(j)} = a_{jl}|Y = c_k) &= \frac{\sum_{i=1}^N I(X_i^{(j)} = a_{jl}, y_j = c_k) + \lambda}{\sum_{i=1}^N I(y_i = c_k) + S_j\lambda}\end{aligned}$$

其中： S_j 为特征属性取值总数， K 为类别的数量

推导朴素贝叶斯中的概率估计公式

贝叶斯估计

对于先验概率

假设进行 N 次实验，先验概率

$$\begin{aligned} P(Y = c_k) &= \frac{1}{K} \\ PK - 1 &= 0 \end{aligned}$$

由频率是概率的极大似然估计

$$P(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k)}{N}$$

可得

$$P(Y = c_k)N - \sum_{i=1}^N I(y_i = c_k) = 0$$

即

$$\lambda(P(Y = c_k)K - 1) + P(Y = c_k)N - \sum_{i=1}^N I(y_i = c_k) = 0$$

可得

$$P(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k) + \lambda}{K\lambda + N}$$

对于似然概率

S_j 是第 j 个特征可能取值的个数，这里假设每个特征的值是等概率分布，因此概率为 $\frac{1}{S_j}$ 。

$$\begin{aligned} P(X^{(j)} = a_{jl} | Y = c_k) &= \frac{1}{S_j} \\ P(X^{(j)} = a_{jl} | Y = c_k) &= \frac{\sum_{i=1}^N I(X^{(j)} = a_{jl}, Y = c_k)}{\sum_{i=1}^N I(Y = c_k)} \end{aligned}$$

同上

$$\lambda(PS_j - 1) + P \sum_{i=1}^N I(Y = c_k) - \sum_{i=1}^N I(X^{(j)} = a_{jl}, Y = c_k) = 0$$

可得

$$P(X^{(j)} = a_{jl} | Y = c_k) = \frac{\lambda + \sum_{i=1}^N I(X^{(j)} = a_{jl}, Y = c_k)}{\sum_{i=1}^N I(Y = c_k) + S_j \cdot \lambda}$$

后验概率最大化的含义

这里来证明一下朴素贝叶斯是如何从损失函数最小化推出后验概率最大化的。

首先我们写出期望风险的公式：

$$\begin{aligned} R_{exp} &= \int \int L(y, f(\vec{x})) P(\vec{x}, y) d\vec{x} dy \\ &= \int_x \int_y L(y, f(\vec{x})) P(y|\vec{x}) dy P(\vec{x}) d\vec{x} \\ &= \mathbb{E}_x \left[\int_y L(y, f(\vec{x})) P(y|\vec{x}) dy \right] \\ &= \mathbb{E}_x \left[\sum_{k=1}^K L(c_k, f(\vec{x})) P(c_k|\vec{x}) \right] \end{aligned}$$

也就是对于每个 x ，我们对 $L(y, f(\vec{x})) P(y|\vec{x})$ 进行最小化：

$$\begin{aligned} f(x) &= \arg \min_{y \in Y} \sum_{k=1}^K L(c_k, y) P(c_k|X = x) \\ &= \arg \min_{y \in Y} \sum_{k=1}^K P(y \neq c_k | X = x) \\ &= \arg \min_{y \in Y} (1 - P(y = c_k | X = x)) \quad \text{所有概率和为1，等价于1减去预测正确的概率} \\ &= \arg \max_{y \in Y} P(y = c_k | X = x) \end{aligned}$$

决策树

理想的决策树：

- 叶结点数最少
- 叶结点深度最小
- 叶结点数最小且叶结点深度最小

防止决策树过拟合：剪枝、强制决策树最大深度、规定叶结点最少样本数

熵：表示随机变量的不确定性度量

$$I(a_i) = p(a_i) \log_2 \frac{1}{p(a_i)}$$

设 X 是一个取有限个值的离散随机变量

$$p(X = x_i) = p_i$$

则随机变量 X 的熵定义为

$$H(X) = - \sum_i p_i \log_2 p_i$$

信息增益

ID3 算法采用的属性选择方式

信息增益其实就是互信息，表示得知特征A的信息而使得不确定性减少的程度， $g(D, A)$ 定义为集合D的经验熵 $H(D)$ 与特征A给定条件下D的经验条件熵 $H(D|A)$ 之差

$$g(D, A) = H(D) - H(D|A)$$

下面这个是大体流程，刚开始看肯定看不懂，做了题之后就明白了

假设拥有训练数据集 D ， $|D|$ 表示其样本容量（样本个数）

设有 K 个类 C_k ， $k = 1, 2, \dots$ ， $|C_k|$ 为属于类 C_k 的样本个数

特征 A 有 n 个不同的取值 a_1, a_2, \dots, a_n ，根据特征 A 的取值，将 D 划分为 n 个子集 D_1, D_2, \dots, D_n

$|D_i|$ 为 D_i 的样本个数，记子集 D_i 中属于类 C_k 的样本集合为 D_{ik} ， $|D_{ik}|$ 为 D_{ik} 的样本个数

1. 数据集 D 的经验熵 $H(D)$

$$H(D) = - \sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$$

2. 特征 A 对数据集 D 的经验条件熵 $H(D|A)$

$$H(D|A) = \sum_{i=1}^N \frac{|D_i|}{|D|} H(D_i) = \sum_{i=1}^N \frac{|D_i|}{|D|} \sum_{k=1}^K \frac{|D_{ik}|}{|D_i|} \log_2 \frac{|D_{ik}|}{|D_i|}$$

3. 信息增益

$$g(D, A) = H(D) - H(D|A)$$

信息增益最大的特征为最优特征！！！！

信息增益比

C4.5算法采用的属性选择方式

特征 A 对于训练数据集 D 的信息增益比定义为信息增益与训练数据集关于特征值 A 的熵之比

$$g_k(D, A) = \frac{g(D, A)}{H_A(D)}$$

其中： n 是特征值 A 的取值个数

$$H_A(D) = - \sum_{i=1}^n \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}$$

信息增益比最大的特征为最优特征！！！！

基尼指数

CART算法采用的属性选择方式

要注意的是：CART算法生成的是二叉树，若一个特征有多种属性，你只能划分是属性A和不是属性A两种情况

- 目标变量离散：生成分类树

- 目标变量连续：生成回归树

使用基尼指数选择最优特征，并决定该特征的最优二值切分点

对于给定的样本集合 D

$$Gini(D) = 1 - \sum_k (\frac{|C_k|}{|D|})^2$$

若样本集合根据特征 A 被划分为 D_1 、 D_2 两个部分，那么在特征 A 条件下，集合 D 的基尼指数定义为：

$$Gini(D, A) = \frac{|D_1|}{|D|}Gini(D_1) + \frac{|D_2|}{|D|}Gini(D_2)$$

基尼指数表示不确定性，基尼指数越大，集合不确定性越大，因此我们要选择基尼指数小的特征进行划分

Logistic回归

Logistic回归是广义线性模型，决策边界是线性的

二项Logistic回归

$$P(Y = 1|X) = \frac{e^{w \cdot x + b}}{1 + e^{w \cdot x + b}}$$
$$P(Y = 0|X) = \frac{1}{1 + e^{w \cdot x + b}}$$

事件发生比：发生与不发生概率之比

$$\frac{P}{1 - P}$$

对数几率

$$\log \frac{P}{1 - P}$$

模型参数估计

使用极大似然估计实现

对于 n 个观测事件 $\{(x_i, y_i)\}_{i=1}^N, x_i \in R^n, y_i \in \{0, 1\}$

已知

$$P(Y = 1|x) = \pi(x) \quad P(Y = 0|x) = 1 - \pi(x)$$

可得到似然函数

$$L(\pi(x)|x, y) = P(X = x, Y = y|\pi(x))$$
$$= \prod_{i=1}^N P(X = x_i)P(Y = y_i|X = x_i, \pi(x))$$
$$\propto \prod_{i=1}^N [\pi(x_i)^{y_i}][1 - \pi(x_i)]^{1-y_i}$$

进行对数化

$$\begin{aligned}\ln L(w) &= \sum_{i=1}^N y_i \ln \pi(x_i) + (1 - y_i) \ln (1 - \pi(x_i)) \\ &= \sum_{i=1}^N y_i \ln \frac{\pi(x_i)}{1 - \pi(x_i)} + \ln (1 - \pi(x_i)) \\ &= \sum_{i=1}^N y_i (w \cdot x_i) - \ln (1 + e^{w \cdot x_i})\end{aligned}$$

对此函数使用梯度下降求极大值，得到 w 估计值

$$\nabla_w (-\ln L(w)) = \sum_{i=1}^N \frac{x_i e^{w \cdot x_i}}{1 + e^{w \cdot x_i}} - y_i x_i$$

当 $y \in \{-1, 1\}$ 时

似然函数为

$$L(w) = - \prod_{i=1}^N \frac{1}{1 + e^{-y_i w x_i}}$$

负对数似然为

$$-\ln L(w) = \sum_{i=1}^N \ln (1 + e^{-y_i w x_i})$$

梯度为

$$\nabla_w (-\ln L(w)) = \sum_{i=1}^N \frac{-y_i x_i e^{-y_i w x_i}}{1 + e^{-y_i w x_i}}$$

支持向量机

定义在特征空间上的间隔最大线性分类器，还包括核函数，使它成为实质上的非线性分类器

不同分类： - 线性可分支持向量机：硬间隔最大化。找到一个超平面，完全正确地将所有样本点分开 - 线性支持向量机：软间隔最大化。找到一个超平面，尽可能正确地将数据点分开，且间隔最大 - 非线性支持向量机：核技巧+软间隔最大化

线性可分支持向量机

线性支持向量机

非线性支持向量机

合页损失函数的推导

在软间隔支持向量机里面，引入松弛变量 ξ_i 。

目标函数为：

$$L = \frac{1}{2} ||w||^2 + C \sum_{i=1}^N \xi_i$$

约束条件：

$$\begin{aligned} y_i(w^T x_i + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned}$$

- 当 $y_i(w^T x_i + b) \geq 1$ 时，说明样本被正确分类且间隔足够大
 - 此时 $1 - y_i(w^T x_i + b) \leq 0$ ，而约束条件为 $\xi_i \geq 1 - y_i(w^T x_i + b)$ ，即 ξ_i 需要大于一个负数且大于 0。
 - 为了让目标函数最小，令 $\xi_i = 0$ 。
- 当 $y_i(w^T x_i + b) < 1$ 时，说明样本间隔不足或者被误分类
 - 此时 $1 - y_i(w^T x_i + b) > 0$ ，同上面的推导，约束条件为 $\xi_i \geq 1 - y_i(w^T x_i + b)$ ，即 ξ_i 需要大于一个正数且大于 0。
 - 为满足约束且使目标函数最小， $\xi_i = 1 - y_i(w^T x_i + b)$ 。

综上所述：

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b))$$

代入目标函数并令 $\lambda = \frac{1}{2C}$ 可得合页损失函数：

$$\sum_{i=1}^N [1 - y_i(w^T x_i + b)]_+ + \lambda ||w||^2$$

其中

$$[z]_+ = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

提升方法

EM算法

EM算法的推导

它的推导和 VAE（变分自编码器）以及 Diffusion Model（扩散模型）很像，都是由于原函数难以求解，去找到它的下界函数，对下界函数求得最优解，达到优化原函数的目的。感兴趣的可以看一下 [VAE](#) 和 [Diffusion Model](#) 的推导。

首先要了解 EM 算法的核心：**EM 算法是计算每个潜在变量 Z 的后验概率，然后通过加权平均，也就是求期望的方式得到一个近似的估计，最后我们要找到一个能让这个近似估计最大的参数 θ 。**

先来认识一下 **KL 散度**：

$$KL(q||p) = \sum_Z q(Z) \log \frac{q(Z)}{p(Z)}$$

描述的是 q 分布和 p 分布之间的距离，因此取值是大于等于 0 的。

首先要明确，EM 算法的核心是在每一步 **最大化 Q 函数**：

$$\begin{aligned} Q(\theta, \theta^{(i)}) &= \mathbb{E}_Z(\log P(Y, Z|\theta)|Y, \theta^{(i)}) \\ &= \sum_Z \log P(Y, Z|\theta) P(Z|Y, \theta^{(i)}) \end{aligned}$$

现在来看如何推导出上面这个 Q 函数。

对于 EM 算法，我们希望最大化观测数据 Y 关于参数 θ 的对数似然函数：

$$\begin{aligned} L(\theta) &= \log P(Y|\theta) = \log \sum_Z P(Y, Z|\theta) \\ &= \log \sum_Z q(Z) \frac{P(Y, Z|\theta)}{q(Z)} \end{aligned}$$

使用 **Jensen 不等式**：

$$\begin{aligned} \log \sum_Z q(Z) \frac{P(Y, Z|\theta)}{q(Z)} &\geq \sum_Z q(Z) \log \frac{P(Y, Z|\theta)}{q(Z)} \\ &= \sum_Z q(Z) \log P(Y, Z|\theta) - \sum_Z q(Z) \log q(Z) \end{aligned}$$

记上式为：

$$\mathcal{L}(q, \theta) = \sum_Z q(Z) \log P(Y, Z|\theta) - \sum_Z q(Z) \log q(Z)$$

因为

$$\log P(Y, Z|\theta) = \log P(Z|Y, \theta) P(Y|\theta)$$

因此有

$$\begin{aligned} \mathcal{L}(q, \theta) &= \sum_Z q(Z) \log P(Y, Z|\theta) - \sum_Z q(Z) \log q(Z) \\ &= \sum_Z q(Z) \log P(Z|Y, \theta) P(Y|\theta) - \sum_Z q(Z) \log q(Z) \\ &= \sum_Z q(Z) (\log P(Z|Y, \theta) + \log P(Y|\theta)) - \sum_Z q(Z) \log q(Z) \\ &= \log P(Y|\theta) \sum_Z q(Z) + \sum_Z q(Z) \log \frac{P(Z|Y, \theta)}{q(Z)} \\ &= \log P(Y|\theta) - KL(q(Z) || P(Z|Y, \theta)) \end{aligned}$$

最终等式

$$\log P(Y|\theta) = KL(q(Z) || P(Z|Y, \theta)) + \mathcal{L}(q, \theta)$$

且

$$\log P(Y|\theta) \geq \mathcal{L}(q, \theta)$$

老师上课的时候没有说 $q(Z)$ 为什么要取 $P(Z|Y, \theta^i)$ ，这里解释一下：

在 E 步的时候，理论上我们可以取任意的 $q(Z)$ 来构造下界 $\mathcal{L}(q, \theta)$ ，但是为了让当前参数的似然函数与下界相等，我们要让

$$q(Z) = P(Z|Y, \theta^i)$$

为什么要这样？这样选择能让在参数 $\theta^{(i)}$ 处，下界紧贴似然函数，即：

$$\log P(Y, Z|\theta^{(i)}) = \mathcal{L}(q, \theta^{(i)})$$

（此时 KL 散度为 0）

当下界紧贴似然函数时，提升下界会对似然函数产生最大程度的提升，收敛速度最快。

如果下界很松（KL 散度很大），即使提升了下界，似然函数可能只有很小的增长。

注意：选择 $q(Z) = P(Z|Y, \theta^{(i)})$ 只是为了在 E 步的当前点 $\theta^{(i)}$ 使得 KL 散度为零、下界紧贴似然函数，从而让后续优化更高效，一旦进入 M 步，参数 θ 开始更新，真实后验 $P(Z|Y, \theta)$ 随之改变，KL 散度立即变为正数，下界与似然函数不再紧贴。

因此

$$\begin{aligned} \log P(Y|\theta) &\geq \mathcal{L}(q, \theta) \\ &= \sum_Z P(Z|Y, \theta^{(i)}) \log P(Y, Z|\theta) + \text{const} \\ &= \underbrace{\mathbb{E}_Z(\log P(Y, Z|\theta)|Y, \theta^{(i)})}_{\text{Q函数}} + \underbrace{\sum_Z P(Z|Y, \theta^{(i)}) \log P(Z|Y, \theta^{(i)})}_{\text{常数}} \end{aligned}$$

单调性的证明

在 E 步的时候，取

$$q(Z) = P(Z|Y, \theta)$$

做完 M 步之后，新下界大于旧下界

$$\mathcal{L}(q, \theta^{(i+1)}) \geq \mathcal{L}(q, \theta^{(i)})$$

旧似然等于旧下界

$$\log P(Y, Z|\theta^i) = \mathcal{L}(q, \theta^{(i)})$$

新似然做完 M 步后 KL 散度大于 0，因此

$$\log P(Y, Z|\theta^{i+1}) \geq \mathcal{L}(q, \theta^{(i+1)})$$

最后有

$$\log P(Y, Z|\theta^{i+1}) \geq \mathcal{L}(q, \theta^{(i+1)}) \geq \mathcal{L}(q, \theta^{(i)}) = \log P(Y, Z|\theta^i)$$

保证了 EM 算法的单调性。

高斯混合模型

高斯分布

$$\phi(y|\theta_k) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)$$

用多个高斯分布的加权组合描述复杂数据的分布

$$\begin{aligned} P(y|\theta) &= \sum_{k=1}^K \alpha_k \phi(y|\theta_k) \\ \text{s.t. } \quad &\alpha_k \geq 0 \ (k = 1, \dots, K), \quad \sum_{k=1}^K \alpha_k = 1 \end{aligned}$$

Q 函数的定义

$$Q(\theta, \theta^{(i)}) = \mathbb{E}_{\gamma}[\log P(y, \gamma|\theta)|y, \theta^{(i)}]$$

观测数据是由高斯混合模型产生的，模型参数为

$$\theta = (\alpha_k, \mu_k, \sigma_k^2)$$

隐变量定义为

$$\gamma_{jk} = \begin{cases} 1 & \text{第 } j \text{ 个观测样本来自第 } k \text{ 个高斯模型} \\ 0 & \end{cases}$$

现在来推导完全数据的似然函数，因为后续会使用这个代入 Q 函数

$$\begin{aligned} P(y, \gamma|\theta) &= \prod_{j=1}^N p(y_j, \gamma_{j1}, \gamma_{j2}, \dots, \gamma_{jK}|\theta) \\ &= \prod_{j=1}^N \prod_{k=1}^K [\alpha_k \phi(y_j|\theta_k)]^{\gamma_{jk}} \end{aligned}$$

令

$$n_k = \sum_{j=1}^N \gamma_{jk}$$

有

$$\prod_{j=1}^N \prod_{k=1}^K [\alpha_k \phi(y_i|\theta_k)]^{\gamma_{jk}} = \prod_{k=1}^K \alpha_k^{n_k} \prod_{j=1}^N [\phi(y_j|\theta_k)]^{\gamma_{jk}}$$

代入到 Q 函数中

$$\begin{aligned}
 Q(\theta, \theta^{(i)}) &= \mathbb{E}_{\gamma} \left\{ \sum_{k=1}^N \{n_k \log \alpha_k + \sum_{j=1}^N \gamma_{jk} [\log \phi(y_j | \theta_k)]\} | y, \theta^{(i)} \right\} \\
 &= \mathbb{E}_{\gamma} \left\{ \sum_{k=1}^N \{n_k \log \alpha_k + \sum_{j=1}^N \gamma_{jk} [(\log \frac{1}{\sqrt{2\pi}} - \log \sigma_k - \frac{(y_j - \mu_k)^2}{2\sigma_k^2})]\} | y, \theta^{(i)} \right\}
 \end{aligned}$$

我们可以很轻松地发现，真正需要计算的只有 $\mathbb{E}(\gamma_{jk} | y, \theta^{(i)})$ ，因为其它的都可以通过期望的线性性质直接拆开得到结果

$$\begin{aligned}
 \mathbb{E}(\gamma_{jk} | y, \theta^{(i)}) &= P(\gamma_{jk} = 1 | y, \theta^{(i)}) \\
 &= \frac{P(\gamma_{jk} = 1 | \theta^{(i)}) P(y_j | \gamma_{jk} = 1, \theta^{(i)})}{P(y | \theta^{(i)})} \\
 &= \frac{\alpha_k \phi(y_j | \theta_k)}{\sum_{k=1}^K \alpha_k \phi(y_j | \theta_k)} \\
 &= \hat{\gamma}_{jk}
 \end{aligned}$$

并且我们有

$$\begin{aligned}
 n_k &= \sum_{j=1}^N \gamma_{jk} \\
 \mathbb{E}[n_k | y, \theta^{(i)}] &= \sum_{j=1}^N \mathbb{E}[\gamma_{jk} | y, \theta^{(i)}] = \sum_{j=1}^N \hat{\gamma}_{jk}
 \end{aligned}$$

最后得到 Q 函数的最终形式

$$\begin{aligned}
 Q(\theta, \theta^{(i)}) &= \sum_{k=1}^N \left\{ \sum_{j=1}^N (\mathbb{E} \gamma_{jk}) \log \alpha_k + \sum_{j=1}^N (\mathbb{E} \gamma_{jk}) \left[\left(\log \frac{1}{\sqrt{2\pi}} - \log \sigma_k - \frac{(y_j - \mu_k)^2}{2\sigma_k^2} \right) \right] | y, \theta^{(i)} \right\} \\
 &= \sum_{k=1}^N \left\{ \sum_{j=1}^N \hat{\gamma}_{jk} \log \alpha_k + \sum_{j=1}^N \hat{\gamma}_{jk} \left[\left(\log \frac{1}{\sqrt{2\pi}} - \log \sigma_k - \frac{(y_j - \mu_k)^2}{2\sigma_k^2} \right) \right] | y, \theta^{(i)} \right\}
 \end{aligned}$$

解释一下上方公式难以理解的点

- $\mathbb{E}(\gamma_{jk} | y, \theta^{(i)})$ 的值其实就是 $P(\gamma_{jk} = 1 | y, \theta^{(i)})$ ，因为前面定义过隐变量的值不是 1 就是 0，这里可以很容易看出
- $\hat{\gamma}_{jk}$ 其实就是样本点 y_j 对高斯分布 k 的隶属度
- $P(y | \theta^{(i)})$ 不是条件概率，而是在参数为 $\theta^{(i)}$ 的情况下的似然函数，而 y 可以来自不同的高斯分布，因此写成求和的形式
- $P(\gamma_{jk} = 1 | \theta^{(i)})$ 就是在参数为 $\theta^{(i)}$ 的情况下数据点来自第 k 个分量的概率，也就是权重 α_k

最后对模型参数 α, μ, σ 进行更新

$$\begin{aligned}\hat{\alpha}_k &= \frac{n_k}{N} = \frac{\sum_{j=1}^N \hat{\gamma}_{jk}}{N} \\ \hat{\mu}_k &= \frac{\sum_{j=1}^N \hat{\gamma}_{jk} y_j}{\sum_{j=1}^N \hat{\gamma}_{jk}} \\ \hat{\sigma}_k &= \frac{\sum_{j=1}^N \hat{\gamma}_{jk} (y_j - \mu_k)^2}{\sum_{j=1}^N \hat{\gamma}_{jk}}\end{aligned}$$

聚类方法

核心是相似度或者距离，因为是通过这两个指标进行聚类的

闵可夫斯基距离

$$d_{ij} = \left(\sum_{k=1}^m |x_{ki} - x_{kj}|^p \right)^{\frac{1}{p}}$$

距离越大，相似度越小

- $p = 1$: 曼哈顿距离
- $p = 2$: 欧氏距离
- $p = \infty$: 切比雪夫距离

相关系数

$$r_{ij} = \frac{\sum_{k=1}^m (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)}{[\sum_{k=1}^m (x_{ki} - \bar{x}_i)^2 \sum_{k=1}^m (x_{kj} - \bar{x}_j)^2]^{\frac{1}{2}}}$$

其中： m 为样本维度，且

$$\bar{x}_i = \frac{1}{m} \sum_{k=1}^m x_{ki} \quad \bar{x}_j = \frac{1}{m} \sum_{k=1}^m x_{kj}$$

相关系数的绝对值越接近1，样本越相似；越接近0，越不相似

夹角余弦

$$S_{ij} = -$$

SVD

PCA

PPT

统计学习概论

1. 什么是统计学习

统计学习是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和分析

2. 什么是过拟合

过拟合指的是模型在已知的数据上预测得很好，但是在未知的数据上预测得很差

1. 增加更多的训练样本总是可以避免过拟合。[错误]

计算Precision和Recall，体温超过38度为阳性患者

No.	1	2	3	4	5	6	7	8	9	10	11	12
Ground Truth	P	P	F	F	F	P	P	F	F	F	F	F
Tempreture	40°C	39°C	38.7°C	38.6°C	38.3°C	38.1°C	37.8°C	37.6°C	37.4°C	37.2°C	37°C	36.6°C

image-20251216104253184

	P	N
T	TP: 3	FP: 1
F	FN: 3	TN: 5

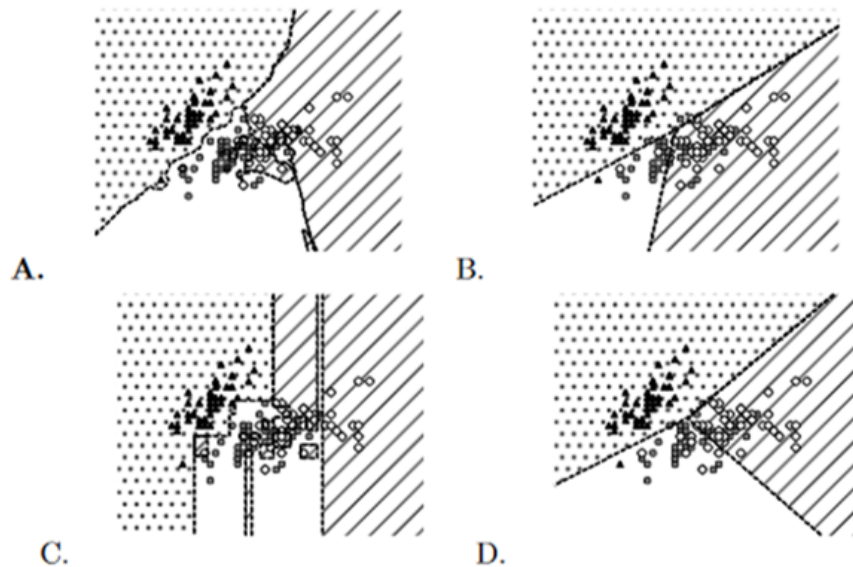
$Precision = \frac{TP}{TP + FP} = 0.75$

$Recall = \frac{TP}{TP + FN} = 0.5$

knn

1. 1-NN分类器的训练误差为0。[正确]
2. 当K=N时，K近邻算法产生的决策边界比1近邻算法更复杂。[错误]
3. 选择不同的距离范数不会影响1-NN的决策边界。[错误]
4. 最近邻算法是一种参数化方法。[错误]

(1 point) Which of the following decision boundaries is most likely to be generated by a k-NN?



3

image-20251216110606101

A, 解析:

- 只有A的分类边界是不规则的

朴素贝叶斯

1. 感知机的训练错误率为0。[错误]
2. 最近邻分类器的训练错误率为0。[正确]
3. 朴素贝叶斯分类器的训练错误率为0。[错误]

决策树

1. ID3方法能够保证得到最优的决策树。[错误]

下列哪些策略可能减小决策树的过拟合问题:

- A. 剪枝
- B. 强制叶节点最小样本数
- C. 强制决策树的最大深度
- D. 确保每个叶节点的样本都属于同一类

ABC, 解析:

- D: 这样子会使得模型为了使叶子节点的样本只有一类, 会使得树的深度骤增, 使模型过拟合

支持向量机

1. 支持向量机计算 $P(y|x)$ 。[错误]
2. 在支持向量机中, 非支持向量的 α_i 的值为0。[正确]

3. 在支持向量机中，正例对应的拉格朗日乘子之和等于负例对应的拉格朗日乘子之和。[正确]

4. 线性可分SVM的损失函数只考虑误分类的样本。[错误]

5. 非线性SVM在训练时不需要设置任何超参数。[错误]

6. SVM在处理非线性问题时使用的技术称为：

A. 神经网络 B. 随机森林

C. 核技巧 D. 梯度提升

C

2. 当你想要提高SVM模型的泛化能力，你应该：

A. 增加正则化参数 C B. 减少正则化参数 C

C. 增加核函数参数 γ D. 减少核函数参数 γ

BD

提升方法

1. 下列哪个说法是正确的：

A. Adaboost算法可直接用于1NN分类器

B. Adaboost算法运用软间隔线性支持向量机做弱分类器可能得到非线性的分类边界

C. Adaboost算法对测试样本的弱分类器输出运用投票法进行决策

D. Adaboost算法适合任何的一组弱分类器

B，解析：

- A：1NN分类器是强分类器，不适合用在Adaboost里面
- C：是加权投票，不是平均投票
- D：不适合，弱分类器要比随即猜测略好，误差小于0.5

2. 在 Adaboost 中，我们从训练样本上的高斯权重分布开始。[错误]

EM算法

EM算法的E步中，我们通常计算什么？（）

A. 参数的后验概率 B. 缺失数据的期望值

C. 观测数据的期望值 D. 参数的最大似然估计

B，解析：

- 应该是完全数据似然函数在隐变量后验概率分布下的期望

聚类

1. 层次聚类中，距离度量不包括以下哪种（）

A. 欧氏距离 B. 曼哈顿距离

C. 余弦相似度 D. 标准差

D

2. 层次聚类算法的类型包括 ()

A. 凝聚的 B. 分裂的 C. 以上都是 D. 以上都不是

C

3. 层次聚类的结果通常表示为 ()

A. 一个聚类数 B. 一个聚类划分
C. 一个树状图（Dendrogram） D. 一个概率模型

C

- 1. 层次聚类方法需要预定义聚类数量。[错误]
- 2. 单链接聚合聚类算法基于两个聚类中点之间的最大距离来合并这两个聚类。[错误]
- 3. K-means算法是一种无监督学习算法。[正确]
- 4. K-means算法是KNN算法的特例。[错误]
- 5. K-means算法总是能找到全局最优解，与初始中心选择无关。[错误]
- 6. K-means算法的K值必须由用户预先指定。[正确]
- 7. K-means算法的时间复杂度是O(nkm)，其中n是数据点的数量，k是簇的数量，m是样本维数。[正确]
- 8. K-means算法可以通过肘部法则来确定K值。[正确]

三要素

方法	模型	策略	算法	损失函数
感知机	二分类超平面	极小化误分类点到超平面的距离	随机梯度下降	误分类点到超平面的距离
knn	\	\	\	\
朴素贝叶斯	特征与类别的联合概率	极大似然估计，极大后验概率估计	EM算法	对数似然损失
决策树	分类树、回归树	正则化的极大似然估计	特征选择、生成、剪枝	对数似然损失
逻辑斯特回归	对数线性模型	正则化的极大似然估计	梯度下降，拟牛顿法	逻辑斯蒂损失
支持向量机	分离超平面	最小间隔最大化，极小化带正则的合页损失	SMO算法	合页损失
提升方法	弱分类器的线性组合	极小化加法模型的指数损失	前向分布加法算法	指数损失
EM算法	含隐变量的概率模型	极大似然估计	迭代算法	对数似然损失

方法	模型	策略	算法	损失函数
层次聚类	聚类树	类内样本距离最小	启发式算法	\
k均值聚类	k中心聚类	样本到类中心的距离之和最小	迭代算法	\
高斯混合模型	高斯混合模型	似然函数最大	EM算法	\
PCA	低维正交空间	方差最大	SVD或者特征值分解	\

我觉得易错的点

- 1. 感知机对偶是用一个 α 来存储每个样本的误分类次数乘上学习率， α 的维度等于样本个数
- 2. knn是懒惰学习，是无参数模型，缺点是时间和空间复杂度高，它的模型三要素是距离度量、k值选择、分类决策规则
- 3. adaboost的训练轮数需要预先指定，训练几轮就有几个弱学习器
- 4. CART算法是二元划分
- 5. 用最小训练误差划分决策树会导致过拟合，不能正确反映划分后两个子集各自的数据分布，生成的决策树也易受噪声影响
- 6. 分裂聚类的时间复杂度是 $O(n^2mk)$ ，层次聚类是层次化的，kmeans是非层次化的
- 7. 逻辑斯蒂回归是线性分类模型，是广义线性模型，但输入输出不在线性关系
- 8. 决策树模型使用的是 `if-then` 规则，特征是互斥且完备，直接应用是 `CLS` 算法
- 9. 奇异值分解的基本定理是奇异值分解对任意实矩阵存在
- 10. 监督学习的基本策略是经验风险最小化和结构风险最小化
- 11. SVM的惩罚系数C太大会导致模型过拟合，C太小会导致模型欠拟合
- 12. SVM是一个求解凸二次规划的问题

考试题型

- 五选择
- 六判断
- 三简答
- 三到四个计算题
- 一道综合题（除了第1问都很难）

重点

::: tabs @tab:active 概论 统计学习定义 统计学习对象 统计学习目的 统计学习三要素 每个模型的三要素是什么 比较用三要素 模型分类概率模型和非概率模型(p11-12) 生成模型和判别模型(生成模型关注联合分布 p28有)。损失函数与风险函数 期望损失 损

失函数的期望 期望风险 经验风险 结构风险 监督学习定义p6 无监督学习定义p8 过拟合定义 解决办法 混淆矩阵 召回率 精确率 F1值

@tab 监督学习 监督学习 七个方法 感知机收敛定理。模型咋写 损失函数是什么 如何推导 knn决策规则 距离度量k值选择造成的影响 不考kd树 朴素贝叶斯对什么做条件独立性假设 最大似然估计和拉普拉斯平滑 决策树 id3 CART 熵以2为底 决策树剪枝的作用 不考回归树 逻辑斯谛回归要会写最大似然函数 会对其求导 支持向量机原问题 对偶问题 约束 KKT条件 原问题如何变成对偶问题 软间隔svm 合页损失 C变大变小的影响 根据kexi取值范围判断是不是支持向量 核方法 正定核充要矩阵:GRAM矩阵半正定 提示方法概念 弱分类器单独分类正确率要超过0.5 adaboost指数损失

@tab 无监督学习 EM算法要会推导。收敛到局部最大值 鞍点 GMM概念 EM算法两步 层次聚类 单链接还是多链接 不考中心距离 平均距离 例14.1 层次聚类图 kmeans 例14.2 会收敛 不是全局最优 SVD定义 怎么求UsigmaV 例15.5 奇异值分解不唯一 定理15.3 截断奇异值近似误差 F范数 pca性质 y_i 的特征值和 方差贡献率是什么 相关系数是什么 用拉格朗日乘子法证明要找方差最大的作为变换向量 要会求相关矩阵 规范化变量情况下的性质