

8.2 对于 0/1 损失函数来说, 指数损失函数并非仅有的一致替代函数. 考虑式(8.5), 试证明: 任意损失函数 $\ell(-f(x)H(x))$, 若对于 $H(x)$ 在区间 $[-\infty, \delta]$ ($\delta > 0$) 上单调递减, 则 ℓ 是 0/1 损失函数的一致替代函数.

$\therefore \ell(-f(x)H(x))$ 关于 $f(x)H(x)$ 在 $(-\infty, \delta]$ ($\delta > 0$) 上单调,

$\therefore \ell$ 的最小值所对应的模型标位置在 δ 右侧, 即最小值点 $x^* = f(x^*)H(x^*) > \delta > 0$

即 $f(x^*)H(x^*)$ 为正数, $\text{sign}(H(x^*)) = f(x^*) = \arg \max_{y \in \{-1, 1\}} P(f(x) = y | x)$

即当指数损失函数最小化, 则分类错误率也将最小化

8.8 MultiBoosting 算法 [Webb, 2000] 将 AdaBoost 作为 Bagging 的基学习器, Iterative Bagging 算法 [Breiman, 2001b] 则是将 Bagging 作为 AdaBoost 的基学习器. 试比较二者的优缺点.

Multi Boosting 先降低模型的偏差再降低方差. 由于集合了 Bagging、AdaBoost, 可以有效降低偏差和方差. 但训练成本和预测成本都会显著增加.

Iterative Bagging 先降低方差再降低偏差. 相比于 Bagging 偏差下降, 但方差上升. 训练成本和预测成本也会显著增加.

- 给定任意的两个相同长度向量 x, y , 其余弦距离为 $1 - \frac{x^T y}{|x||y|}$, 证明余弦距离不满足传递性, 而余弦夹角 $\arccos\left(\frac{x^T y}{|x||y|}\right)$ 满足

(1) 取 $x = (1, 0)^T$, $y = (1, 1)^T$, $z = (0, 1)^T$

使用余弦距离, 有 $\text{dist}(x, z) = 1 - \frac{x^T z}{|x||z|} = 1$

$$\text{dist}(x, y) = 1 - \frac{x^T y}{|x||y|} = 1 - \frac{1}{\sqrt{2}}$$

$$\text{dist}(y, z) = 1 - \frac{y^T z}{|y||z|} = 1 - \frac{1}{\sqrt{2}}$$

$$\text{dist}(x, y) + \text{dist}(y, z) = 2 - \sqrt{2} < 1 = \text{dist}(x, z)$$

\therefore 余弦距离不满足传递性.

(2) 不妨设 x, y, z 均为单位向量, 即 $|x| = |y| = |z| = 1$

要证 $\arccos x^T z \leq \arccos x^T y + \arccos y^T z$

① $\arccos x^T y + \arccos y^T z > \pi$ 则不等式显然成立

② $0 \leq \arccos x^T y + \arccos y^T z \leq \pi$

又 $\arccos x \in [0, \pi]$, $\forall x \in \mathbb{R}$, $\cos x$ 在 $[0, \pi]$ 上单调递减

记 $x^T z = a$, $x^T y = b$, $y^T z = c$

故即证 $\cos(\arccos a) \geq \cos(\arccos b + \arccos c)$

$$a \geq \cos(\arccos b) \cos(\arccos c) - \sin(\arccos b) \sin(\arccos c)$$

$$\sin(\arccos b) \sin(\arccos c) \geq bc - a$$

$$(1 - \cos^2(\arccos b))(1 - \cos^2(\arccos c)) \geq (bc - a)^2$$

$$(1 - b^2)(1 - c^2) \geq b^2 c^2 - 2abc + a^2$$

$$1 - b^2 - c^2 + b^2 c^2 \geq b^2 c^2 - 2abc + a^2$$

即证

$$a^2 + b^2 + c^2 - 2abc - 1 \leq 0$$

设 $m = (x, y, z)$

设 $A = m^T m = \begin{pmatrix} 1 & x^T y & x^T z \\ y^T x & 1 & y^T z \\ z^T x & z^T y & 1 \end{pmatrix}$

则 $|A| = |m|^2 = 1 + 2abc - a^2 - b^2 - c^2 \geq 0$, 证毕.

• 证明k-means算法的收敛性

对于每个数据点 x_n , 引入一组对应的二值指示变量 $r_{nk} \in \{0, 1\}$,

其中 $k=1, \dots, K$, 表示数据点 x_n 属于 K 个聚类中心中的哪一个.

定义目标函数 $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$,

表示每个数据点与它被分配到的向量 μ_k 之间的距离平方和.

目标是找到 $\{r_{nk}\}$ 和 $\{\mu_k\}$ 的值, 使得 J 达最小值. 可以使用一种迭代的方法.

1. 初始化 μ_k .
2. 关于 r_{nk} 最小化 J , 保持 μ_k 固定.
3. 关于 μ_k 最小化 J , 保持 r_{nk} 固定.
4. 不断重复步骤2和3, 直至收敛.

可见更新 r_{nk} 和更新 μ_k 的两个阶段分别对应于EM算法中的E和M.

首先考虑 r_{nk} , J 是 r_{nk} 的一个线性函数, 因此最优化过程可以得到解析解.

将数据点的聚类设置为最近的聚类中心 $r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{else} \end{cases}$

然后考虑 r_{nk} 固定时, 关于 μ_k 的优化. J 是 μ_k 的一个二次函数, 令它关于 μ_k 的导数为0, 即可达到最小值, 即

$$\begin{aligned} \frac{\partial J}{\partial \mu_k} &= \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \\ \mu_k &= \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}, \quad \text{即令 } \mu_k \text{ 为类别 } k \text{ 的所有数据点的均值.} \end{aligned}$$

重新为数据点分配聚类为步骤2以及重新计算聚类均值的步骤重复进行, 直至聚类的分配不改变或迭代次数超过了某个最大值.

由于每个阶段都减小了目标函数 J 的值, 因此算法的收敛性得到了保证.

但算法可能收敛到 J 的一个局部最小值而不是全局最小值.

- 在k-means算法中替换欧式距离为其他任意的度量, 请问“聚类簇”中心如何计算?

$$J = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}(x, \mu_i), \quad \text{若 } J \text{ 可导, 令 } \frac{\partial J}{\partial \mu_i} = 0 \quad (i=1, \dots, k)$$

1. 随机选取 k 个样本作为初始均值向量 $\{\mu_1, \mu_2, \dots, \mu_k\}$
2. 求 $\frac{\partial J}{\partial \mu_i} = 0$, 更新 μ_i ($i=1, \dots, k$), 若能求出解析解, 则使用解析解作为更新值, 否则使用梯度下降法.