

# 统计学习第三章

2022年8月10日 星期三 12:50

## K近邻法

训练数据集:  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

输出:  $y = \arg \max_{C_j} \sum_{x_i \in N_k(x)} I(y_i = C_j), \dots$  ( $I$  为指示函数,  $y_i = C_j$  时  $I$  为 1, 否则为 0)

设特征空间  $X$  为  $n$  维实数向量空间  $\mathbb{R}^n$ ,  $x_i, x_j \in X$

$$x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(n)})^T$$

$$x_j = (x_j^{(1)}, x_j^{(2)}, \dots, x_j^{(n)})^T$$

其中  $x_i, x_j$  的  $L_p$  距离定义为:

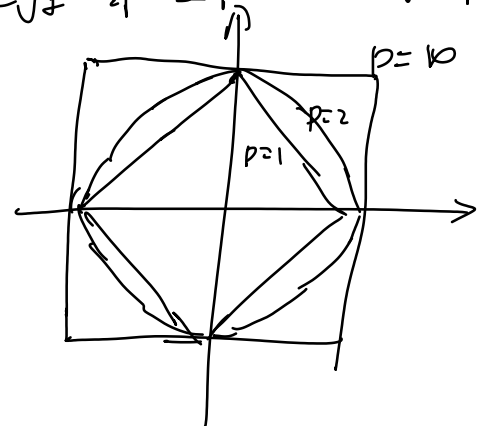
$$L_p(x_i, x_j) = \left( \sum_{l=1}^n |x_i^{(l)} - x_j^{(l)}|^p \right)^{\frac{1}{p}} \quad p \geq 1$$

$$p=2: \text{欧氏距离: } L_2(x_i, x_j) = \left( \sum_{l=1}^n |x_i^{(l)} - x_j^{(l)}|^2 \right)^{\frac{1}{2}}$$

$$p=1: \text{曼哈顿距离: } L_1(x_i, x_j) = \sum_{l=1}^n |x_i^{(l)} - x_j^{(l)}|$$

$$p=\infty: \text{各坐标距离的 max: } L_\infty(x_i, x_j) = \max_l |x_i^{(l)} - x_j^{(l)}|$$

eg1: 取  $L_p = 1$  时  $p=1, p=2, p=\infty$  时的图形:



「解释」:

- $p=1$  指特征向量中每一个值加起来和为 1
- $p=2$  指图形上一点到原点距离为 1
- $p=\infty$  指正方形上一点最大值 (取了绝对值之后)

eg2:  $x_1 = (1, 1)^T, x_2 = (5, 1)^T, x_3 = (4, 4)^T$  在  $p$  取不同值时  $L_p$  距离下  $x_1$  的最近邻点.

1°  $p=1$  时,  $L_1(x_1, x_2) = 4, L_1(x_1, x_3) = 6$

2°  $p=2$  时,  $L_2(x_1, x_2) = 4, L_2(x_1, x_3) = 3\sqrt{2}$

3°  $p=\infty$  时,  $L_\infty(x_1, x_2) = 4, L_\infty(x_1, x_3) = 3$

$p=1, 2$  时, 最近邻点为  $x_2$ ,  $L=\infty$  时, 最近邻点为  $x_3$

注:  $k$  值较小, 相当于用较小邻域

$k$  值较大, 相当于用较大邻域.

## 分类决策规则

$$f: \mathbb{R}^n \rightarrow \{C_1, C_2, \dots, C_K\}$$

$$\text{误分类概率: } P(Y \neq f(X)) = 1 - P(Y = f(X))$$

$$\frac{1}{K} \sum_{x_i \in N_k(x)} I(y_i \neq C_j) = 1 - \frac{1}{K} \sum_{x_i \in N_k(x)} I(y_i = C_j)$$

( $I$  为指示函数, 是为 1, 不是为 0)

## K近邻法的实现: kd 树.

构造 kd 树: 构造根结点, 将实例保存在相应结点.

算法: 输入:  $T = \{x_1, x_2, \dots, x_n\}$ , 其中  $x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(K)})^T$   
输出: kd 树.

1) 开始: 构造根结点, 对应于包含  $T$  的  $K$  维空间的超矩形区域  
选择  $x^{(1)}$  为坐标轴, 以  $T$  中所有实例  $x^{(1)}$  坐标的中位数为切分点

2) 重复: 对深度为  $j$  的结点, 选择  $x^{(l)}$  为切分的坐标轴,  
 $l = j \pmod K + 1$

3) 直到两个子区域没有实例存在时停止, 从而形成 kd 树的区域划分

eg: 给定二维数据集:

$$T = \{(2, 3)^T, (5, 4)^T, (9, 6)^T, (4, 7)^T, (8, 1)^T, (7, 2)^T\}$$

构造一个平衡 kd 树.

