

KNN (K邻近法)

1. K邻近算法

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

其中, $x_i \in X \subset R^n$ 是实例的特征向量, $y_i \in Y = \{c_1, c_2, \dots, c_K\}$ 为实例的类别, $i = 1, 2, \dots, N$; 实例特征向量 x ;

(1) 根据给定的度量距离, 在训练集 T 中找出与 x 最邻近的 k 个点, 涵盖这 k 个点的 x 的邻域记作 $N_k(x)$

(2) 在 $N_k(x)$ 中根据分类决策规则(如多数表决)决定 x 的类别 y

$$y = \underset{c_j}{\operatorname{argmax}} \sum_{x_i \in N_k(x)} I(y_i = c_j), \quad i = 1, 2, \dots, N; \quad j = 1, 2, \dots, K$$

2. K邻近模型

2.1 距离度量

设特征空间 X 是 n 维实数向量空间 R^n , $x_i, x_j \in X$, x_i, x_j 的 L_p 距离定义为

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

$p=1$ 时称为曼哈顿距离 (Manhattan Distance)

$p=2$ 时称为欧氏距离 (Euclidean Distance)

$$p=\infty \text{ 时, } L_\infty(x_i, x_j) = \max_l |x_i^{(l)} - x_j^{(l)}|$$

2.2. K值的选择

若 k 值较小, 易发生过拟合; 若 k 值较大, 可能造成较大近似误差

一般使用交叉验证法选取 k 的值

2.3 分类决策规则

一般选取“投票法” (majority voting rule)

解读: 多数投票法等价于经验损失最小化 (0-1 损失)

$$\hat{p}(Y \neq f(x)) = \frac{1}{k} \sum_{x_i \in N_k(x)} I(y_i \neq c_j) = 1 - \sum_{x_i \in N_k(x)} I(y_i = c_j)$$

最大化 $\sum_{x_i \in N_k(x)} I(y_i = c_j)$ 等价于最小化 $\frac{1}{k} \sum_{x_i \in N_k(x)} I(y_i \neq c_j)$