



机器学习 - KNN学习笔记 /

一、模型

三要素: ① 距离度量 ② K-值的选择 ③ 分类决策规则

利用训练数据集对特征向量空间进行划分

二、距离度量

① L_p 距离: $L_p(x_i, x_j) = (\sum_{k=1}^n |x_i^{(k)} - x_j^{(k)}|^p)^{\frac{1}{p}}$

② 欧氏距离: $L_2(x_i, x_j) = (\sum_{k=1}^n |x_i^{(k)} - x_j^{(k)}|^2)^{\frac{1}{2}}$

③ 曼哈顿距离 $L_1(x_i, x_j) = \sum_{k=1}^n |x_i^{(k)} - x_j^{(k)}|$

④ 切比雪夫距离 $L_\infty(x_i, x_j) = \max_k |x_i^{(k)} - x_j^{(k)}|$

例: $x_1 = (1, 1)^T$ $x_2 = (15, 1)^T$ $x_3 = (14, 4)^T$

当 $p=1$ 时: $L_1(x_1, x_2) = 4$ $L_1(x_1, x_3) = 6$ 故 x_2 为最近邻点

当 $p=2$ 时 $L_2(x_1, x_2) = 4$ $L_2(x_1, x_3) = 4.24$ 故 x_2 为最近邻点

当 $p=3$ 时 $L_3(x_1, x_2) = 4$ $L_3(x_1, x_3) = 3.78$ 故 x_3 为最近邻点

当 $p=4$ 时 $L_4(x_1, x_2) = 4$ $L_4(x_1, x_3) = 3.5$ 故 x_3 为最近邻点

当 $p > 4$ 时 $L_p(x_1, x_2) = 4$ $L_p(x_1, x_3) < 4$ 故 x_3 为最近邻点



当 $p=\infty$ 时 $L_\infty(x_1, x_2)=4$ $L_\infty(x_1, x_2)=3$ 的 x_2 为最近邻点

综上: 较小的 k 值: 学习近似误差 \downarrow 但估计误差 \uparrow 敏感性 \uparrow 模型复杂

较大 k 值: 估计误差 \downarrow 近似误差 \uparrow

k 的取值可以交叉验证选取, 低于样本量的平方根

三. 分类决策规则

多数表决规则: 由输入实例的 k 个邻近的训练实例中的多数决定输入实例的类。

分类函数: $f: \mathbb{R}^n \rightarrow \{c_1, c_2, \dots, c_k\}$

0-1 损失函数: $L(Y, f(x)) = \begin{cases} 1, & Y \neq f(x) \\ 0, & Y = f(x) \end{cases}$

误分类概率: $P(Y \neq f(x)) = 1 - P(Y = f(x))$

$$x \in X, \frac{1}{N} \sum_{i \in N(x)} I(y_i \neq c_i) = 1 - \frac{1}{N} \sum_{i \in N(x)} I(y_i = c_i)$$

$$\arg \max_{c_i} \sum_{i \in N(x)} I(y_i = c_i)$$

KNN 应对分类回归问题