

定义: 若样本点未知, 则通过找与它距离最近的  $K$  个邻居样本点, 通过 投票法 / 均值法 获取样本的 类别 / 值, 完成对样本的 分类 / 预测.

$K$  近邻搜索作用 { 分类  
预测

关键点:

$K$  近邻搜索

如何确定  $K$   
①

如何定义距离从而找出近邻. ②

① 如何确定合适的  $K$ ?

重要性:  $K$  过小  $\Rightarrow$  过拟合

$K$  过大  $\Rightarrow$  欠拟合

Method: A: 将  $K$  设置得大一些, 但设置近邻投票权重 (与距离成反比)

B:  $m$  折交叉验证, 将  $K$  取不同值. 于每种值

下进行  $m$  折交叉验证得平均值, 最后选出平均误差最小的  $K$  的平均值

$\triangle$  此交叉验证方法在 ① 决策树 & 随机森林 ② 岭回归与 LASSO 回归中均有使用.

② 如何确定距离?

Method: A: 欧式距离

$A = (x_1, x_2, \dots, x_n)$

$B = (y_1, y_2, \dots, y_n)$

$$d_{A,B} = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

B: 曼哈顿距离.

$$d_{A,B} = |y_1 - x_1| + |y_2 - x_2| + \dots + |y_n - x_n|$$

C: 余弦相似度.

$$\begin{aligned} \text{Similarity}_{A,B} = \cos\theta &= \frac{\vec{A} \cdot \vec{B}}{||\vec{A}|| ||\vec{B}||} \\ &= \frac{x_1 y_1 + x_2 y_2 + \dots + x_n y_n}{\sqrt{x_1^2 + x_2^2 + \dots + x_n^2} \cdot \sqrt{y_1^2 + \dots + y_n^2}} \end{aligned}$$

$\cos\theta \uparrow$ ,  $\theta \downarrow$ , 则相似度越大.

D: 闵可夫斯基距离公式.

$$d_{A,B} = \sqrt[p]{\sum_{i=1}^n (x_i - y_i)^p}$$

当  $p=1$  时, 为曼哈顿距离;

$p=2$  时, 为欧式距离;