

K近邻

K近邻

1.模型

K近邻模型由3要素构成：距离度量，k指选择和分类决策规则

距离度量：常见的有闵可夫斯基距离，欧式距离，曼哈顿距离

k值：直接关系到模型的复杂度，k越小，预测结果越精确，模型越复杂，对近邻的实例点越敏感

k太大则模型太简单，会忽略训练实例中的大量有用信息

分类决策规则：多数表决规则。其实就是分类损失函数为0-1函数时的损失函数最小化原则

算法1（k邻近法）

输入：训练数据集

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

其中 $x_i \in \mathcal{X} \subseteq \mathbb{R}^n$ 为样本的特征向量， $y_i \in \mathcal{Y} = \{c_1, c_2, \dots, c_K\}$ 为实例的类别， $i=1, 2, \dots, N$ ；样本特征向量 x （新样本）；

输出：样本 x 所属的类 y 。

(1) 根据给定的距离度量，在训练集 T 中找出与 x 最相邻的 k 个点，涵盖这 k 个点的邻域记作 $N_k(x)$ ；

(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 ：

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

式中 I 为指示函数，即当 $y_i = c_j$ 时为1，否则为0。

2.kd树

可参照<https://zhuanlan.zhihu.com/p/22557068>理解

每次寻找

算法 3.2（构造平衡 kd 树）

输入： k 维空间数据集 $T = \{x_1, x_2, \dots, x_N\}$ ，

其中 $x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(k)})^T$ ， $i=1, 2, \dots, N$ ；

输出： kd 树。

(1) 开始：构造根结点，根结点对应于包含 T 的 k 维空间的超矩形区域。

选择 $x^{(l)}$ 为坐标轴，以 T 中所有实例的 $x^{(l)}$ 坐标的中位数为切分点，将根结点对应的超矩形区域切分为两个子区域。切分由通过切分点并与坐标轴 $x^{(l)}$ 垂直的超平面实现。

由根结点生成深度为 1 的左、右子结点：左子结点对应坐标 $x^{(l)}$ 小于切分点的子区域，右子结点对应于坐标 $x^{(l)}$ 大于切分点的子区域。

将落在切分超平面上的实例点保存在根结点。

(2) 重复：对深度为 j 的结点，选择 $x^{(l)}$ 为切分的坐标轴， $l = j(\bmod k) + 1$ ，以该结点的区域中所有实例的 $x^{(l)}$ 坐标的中位数为切分点，将该结点对应的超矩形区域切分为两个子区域。切分由通过切分点并与坐标轴 $x^{(l)}$ 垂直的超平面实现。

由该结点生成深度为 $j+1$ 的左、右子结点：左子结点对应坐标 $x^{(l)}$ 小于切分点的子区域，右子结点对应坐标 $x^{(l)}$ 大于切分点的子区域。

将落在切分超平面上的实例点保存在该结点。

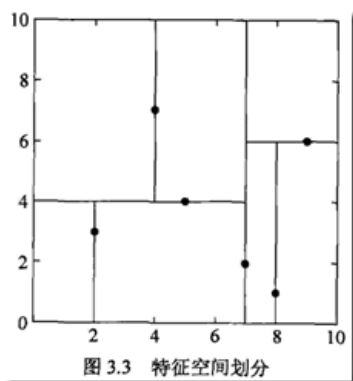
(3) 直到两个子区域没有实例存在时停止。从而形成 kd 树的区域划分。 ■

例子：

例 3.2 给定一个二维空间的数据集：

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

构造一个平衡 kd 树^③。

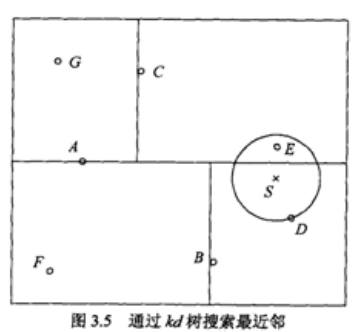


搜索kd树:

对于搜索k近邻, 就是k次搜索最近邻, 因此只需要关注最近邻搜索

例子:

例 3.3 给定一个如图 3.5 所示的 kd 树, 根结点为 A , 其子结点为 B, C 等. 树上共存储 7 个实例点; 另有一个输入目标实例点 S , 求 S 的最近邻.



假设要搜索的点为 S , 首先找到包含点 S 的那个叶节点 (这里是 D), 以 S 为圆心, SD 为半径画圆。如果在 D 节点的所有实例

点钟有距离 S 更近的点, 那么该点一定在这个圆内。以 D 节点为近似最近邻, 然后回溯到父节点 B , 在 B 的另一个子节点 F 内搜索最近邻。 F 的区域与圆不相交, 不可能有最近邻点。继续向上回溯到节点 A , 在 A 的另一个子节点 C 内搜索最近邻。 C 区域与圆相交, 切点 E 比 D 距离 S 更近, E 成为新的最近邻。最后得到点 E 是 S 的最近邻。