

# Chap03-kNN

## I. 算法 (k-nearest neighbor, k-NN)

输入: 训练数据集  $T = \{(x_1, y_1), \dots, (x_N, y_N)\}$   
 $x_i \in X \subseteq \mathbb{R}^n$ ,  $y_i \in Y = \{c_1, c_2, \dots, c_K\}$

输出:  $x$  所属的类  $y$

- (1) 根据给定的距离度量, 在训练集  $T$  中找出与  $x$  最近邻的  $k$  个点, 涵盖这  $k$  个点的  $x$  邻域记为  $N_k(x)$
- (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  $y$ :

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

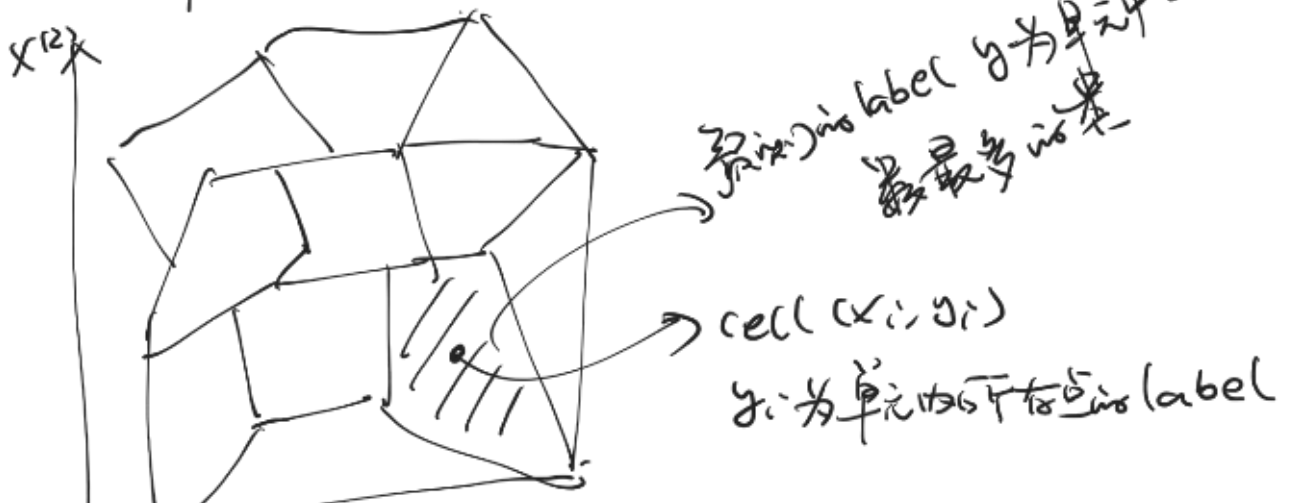
$$I(y_i = c_j) = \begin{cases} 1 & y_i = c_j \\ 0 & \text{else} \end{cases}$$

$k=1$  时, 最近邻算法.

kNN 没有显式的学习过程?

## II. k-NN 模型

### 1. 模型





## 2. 距离度量 $L_p$ -norms or Minkowski-distance

$$\begin{aligned} & X_i^{(1)}, X_j^{(1)} \in \mathbb{R}^n \\ L_p(X_i, X_j) &= \left( \sum_{l=1}^n |X_i^{(l)} - X_j^{(l)}|^p \right)^{1/p} \\ \left\{ \begin{aligned} L_2 &= \left( \sum_{l=1}^n (X_i^{(l)} - X_j^{(l)})^2 \right)^{1/2} \\ L_1 &= \sum_{l=1}^n |X_i^{(l)} - X_j^{(l)}| \\ L_\infty &= \max_l |X_i^{(l)} - X_j^{(l)}| \end{aligned} \right. \end{aligned}$$

由不同的距离度量所确定的最近邻点不同.

### 3. $\lambda$ 值的选择

$K \downarrow \rightarrow \text{训练误差} \downarrow \rightarrow \text{预测误差} \uparrow$   
 $K \uparrow \rightarrow \text{训练误差} \uparrow \rightarrow \text{预测误差} \downarrow$

( $\Leftarrow$ ) 模型复杂度:  $k$  值减小  $\Rightarrow$  容易 over-fitting  
 $k$  值增加  $\Rightarrow$  模型过于简单

在应用中,  $k$  值一般取一个比较小的数值, 然后通过交叉验证法来选取最优的  $k$  值。

4. 决策的决策规则)

多数表决:  $N_k(x)$  中的多数类  $\Rightarrow x$  属于  $y$

0-1 loss function

$N_k(x)$  中的实例被误分类的概率为

$$\begin{aligned} p &= \frac{1}{K} \sum_{x_i \in N_k(x)} I(y_i \neq C_{\hat{y}}) \\ &= \frac{1}{K} \sum_{x_i \in N_k(x)} [1 - I(y_i = C_{\hat{y}})] \\ &= 1 - \frac{1}{K} \sum_{x_i \in N_k(x)} I(y_i = C_{\hat{y}}) \end{aligned}$$

to min  
argmax

## IV. kd tree

kNN 关键是如何进行最近邻搜索.

Linear Scan =  $N + N - 1 + \dots + 1 \propto O(N^2)$

1. 构造平衡 kd 树 作为代表

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

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

Output: kd-tree

SPLIT( $T$ , Node):

if  $T.shape[0] = 1$   
Node.val =  $T$   
return

$\hat{j} = \text{Node} \rightarrow \text{depth}$

$l = \hat{j} \bmod k + 1$

...

Use partition to split  $x$  into  
 $(L, X_{\text{median}}^{(k)}, R)$   $\rightarrow O(n)$

$\text{Node} \rightarrow \text{Val} = X_{\text{median}}^{(k)}$

if  $T.\text{shape}[n] = 2$  Use larger one as median  
 $\text{Node} \rightarrow \text{left} = LNode$

$LNode \rightarrow \text{Val} = L$

$LNode \rightarrow \text{depth} = \delta + 1$

return

$\text{Node} \rightarrow \text{left} = LNode$

$\text{Node} \rightarrow \text{right} = RNode$

$LNode \rightarrow \text{depth} = \delta + 1$

$RNode \rightarrow \text{depth} = \delta + 1$

$\text{SPLIT}(L, LNode)$

$\text{SPLIT}(R, RNode)$

$\text{BUILD-KD-TREE}(T, \text{root}) =$   
 $\text{SPLIT}(T, \text{root})$   $O(n \log n)$

2. KD 树的最近邻搜索:

<https://zhuanlan.zhihu.com/p/23966698>

$\text{Find-Bottom}(X, \text{root}) =$

$\text{temp} = \text{root}, \text{dim} = X.\text{size}$

```

while temp not None:
    p = temp
    key = temp.depth (mod dim) + 1
    if X[key] < temp.val[key]:
        temp = temp.left
    else:
        temp = temp.right
return p

```

$k\_list = [ ]$ ,  $search\_list = [ ]$

$(c\_nearest\_search(X, lowest\_node, k)) =$

```

current_node = lowest_node
if current_node not in search_list:
    search_list.append(current_node)
if k_list.size < k:
    k_list.append(current_node)
elif D(X, current_node.val) < Max(k_list):
    replace(k_list, current_node)

```

while current\_node not root:

```

temp = current_node.parent
if temp not in search_list:
    search_list.append(temp)
    if k_list.size < k:
        k_list.append(temp)
    elif D(X, temp.val) < Max(k_list):
        replace(k_list, temp)

```

```

        replace (k-list, temp)
        key = temp → depth (mod d) + 1
        sd = abs(temp → val(key) - x[key])
        if sd < max(k-list) or
            k-list < k =
            if current_node = temp → left
                and temp → right not None:
                low = Find-Bottom(x, temp, r)
            elif current_node = temp → right
                and temp → left not None:
                low = Find-Bottom(x, temp, left)
            k-nearest(x, low, k)
        else: current = current → parent
    else = current = current → parent

```

上次修改：02:19