

9.4.1 k均值算法.

$$D = \{d_1, \dots, d_m\} \quad C = \{C_1, \dots, C_k\} \text{ 最小化}$$

$$E = \sum_{i=1}^k \sum_{d \in C_i} \|d - \mu_i\|_2^2 \quad \mu_i = \frac{1}{|C_i|} \sum_{d \in C_i} d$$

输入: 样本集  $D = \{d_1, d_2, \dots, d_m\}$

聚类数  $k$ .

过程:

从  $D$  中随机选择  $k$  个样本作为初始均值向量  $\{\mu_1, \mu_2, \dots, \mu_k\}$ .

repeat:

  令  $C_i = \emptyset \quad (1 \leq i \leq k)$

  for  $j = 1, 2, \dots, m$  do

    计算样本  $d_j$  与均值向量  $\mu_i \ (1 \leq i \leq k)$  距离:  $d_{ji} = \|d_j - \mu_i\|_2$

    根据距离最近的均值向量确定  $d_j$  的簇标记:  $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ji}$

    将样本  $d_j$  列入相应簇:  $C_{\lambda_j} = C_{\lambda_j} \cup \{d_j\}$

  end for

  for  $i = 1, 2, \dots, k$  do

    计算新均值向量:  $\mu'_i = \frac{1}{|C_i|} \sum_{d \in C_i} d$

    if  $\mu'_i \neq \mu_i$  then

$\mu_i = \mu'_i$

  until 当前均值向量均未更改.

输出: 簇划分  $C = \{C_1, C_2, \dots, C_k\}$ .









