

PART5 聚类

2020-05-25

河北师范大学 软件学院

基本内容:

1. 什么是聚类? 什么是分类? 二者的区别与联系。
2. (1)以连续实值特征向量描述的样本点之间的距离度量方式

欧式距离、切比雪夫距离、曼哈顿距离、马氏距离

例: 若多维特征空间特征分布的协方差矩阵为 Σ , 则该空间任意两点之间

x_i, x_j 的马氏距离为:

$$d_M(x_i, x_j) = \left[(x_i - x_j)^T \Sigma^{-1} (x_i - x_j) \right]^{0.5}$$

对于 1 维特征空间, 特征分布标准差 σ , 则该空间任意两点之间 x_i, x_j 的马氏距离为:

$$d_M(x_i, x_j) = \left| \frac{x_i - x_j}{\sigma} \right|$$

(2)样本点与集合(例: 聚类簇)之间距离

(3)集合与集合之间(例: 两个聚类簇之间)的距离

最小距离、最远距离、平均距离

3. **动态聚类**。掌握 K-Means Clustering 算法(目标函数? 哪些因素影响聚类性能? 实现步骤?)
4. **系统聚类**。以**聚合式系统聚类**为例, 掌握系统聚类(实现步骤? 哪些因素影响聚类性能?)。
5. **密度聚类**。以 DBSCAN 算法为例, 理解密度聚类实现的基本流程, 掌握有关概念, 哪些因素会影响聚类的效果。

练习:

1. 给定数据集 $D = \{x_i, i = 1, \dots, m\}$, 其中 $x_i \in R^d$.若采用 K-均值聚类 算法将该数据集 D 划分为**K簇** $\{C_1, \dots, C_K\}$, 请完成如下工作:

- (1) 写出 K-Means 算法所对应的准则函数，并给出必要的参数说明；
- (2) 对 K-Means 算法的实现过程进行描述；
- (3) 指出影响数据集 D 划分结果的可能因素。

解：

(1)

聚类准则--“总的簇内误差平方和”最小

将样本集 D 划分成 k 个簇: $D = C_1 \cup \dots \cup C_k$

总的簇内误差平方和 $E(\mu_1, \dots, \mu_k, C_1, \dots, C_k) = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$

其中 C_i --第 i 个簇(聚类), $i = 1, \dots, k$

$N_i = |C_i|$ --第 i 个簇的样本数目

μ_i --第 i 个聚类簇中心, $\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$

(2)

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;

聚类簇数 k .

过程:

```

1: 从  $D$  中随机选择  $k$  个样本作为初始均值向量  $\{\mu_1, \mu_2, \dots, \mu_k\}$ 
2: repeat
3:   令  $C_i = \emptyset$  ( $1 \leq i \leq k$ )
4:   for  $j = 1, 2, \dots, m$  do
5:     计算样本  $x_j$  与各均值向量  $\mu_i$  ( $1 \leq i \leq k$ ) 的距离:  $d_{ji} = \|x_j - \mu_i\|_2$ ;
6:     根据距离最近的均值向量确定  $x_j$  的簇标记:  $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ji}$ ;
7:     将样本  $x_j$  划入相应的簇:  $C_{\lambda_j} = C_{\lambda_j} \cup \{x_j\}$ ;
8:   end for
9:   for  $i = 1, 2, \dots, k$  do
10:    计算新均值向量:  $\mu'_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$ ;
11:    if  $\mu'_i \neq \mu_i$  then
12:      将当前均值向量  $\mu_i$  更新为  $\mu'_i$ 
13:    else
14:      保持当前均值向量不变
15:    end if
16:  end for
17: until 当前均值向量均未更新
输出: 簇划分  $C = \{C_1, C_2, \dots, C_k\}$ 

```

(3) K值的大小、样本集是否进行规范化预处理、聚类中心的初始化方式

2. (1)若要采用合并式层次聚类将样本集 $D = \{x_i, i = 1, \dots, m\}$ 划分为 k 个聚类簇，其

中 $x_i \in R^p$. 请对该聚类算法的实现流程予以描述.

(2)上述聚类过程中, 需要进行不同聚类簇之间的距离计算, 请分别采用最近距离、最远距离, 估算任意两个聚类簇 C_j, C_l 之间的距离.

解:

(1)

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;
 聚类簇距离度量函数 d ;
 聚类簇数 k .

过程:

```

1: for  $j = 1, 2, \dots, m$  do
2:    $C_j = \{x_j\}$ 
3: end for
4: for  $i = 1, 2, \dots, m$  do
5:   for  $j = 1, 2, \dots, m$  do
6:      $M(i, j) = d(C_i, C_j)$ ;
7:      $M(j, i) = M(i, j)$ 
8:   end for
9: end for
10: 设置当前聚类簇个数:  $q = m$ 
11: while  $q > k$  do
12:   找出距离最近的两个聚类簇  $C_{i^*}$  和  $C_{j^*}$ ;
13:   合并  $C_{i^*}$  和  $C_{j^*}$ :  $C_{i^*} = C_{i^*} \cup C_{j^*}$ ;
14:   for  $j = j^* + 1, j^* + 2, \dots, q$  do
15:     将聚类簇  $C_j$  重编号为  $C_{j-1}$ 
16:   end for
17:   删除距离矩阵  $M$  的第  $j^*$  行与第  $j^*$  列;
18:   for  $j = 1, 2, \dots, q - 1$  do
19:      $M(i^*, j) = d(C_{i^*}, C_j)$ ;
20:      $M(j, i^*) = M(i^*, j)$ 
21:   end for
22:    $q = q - 1$ 
23: end while

```

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

(2)

最近距离 $d_{\min}(C_j, C_l) = \min_{\substack{x \in C_j \\ z \in C_l}} \text{dist}(x, z)$

最远距离 $d_{\max}(C_j, C_l) = \max_{\substack{x \in C_j \\ z \in C_l}} \text{dist}(x, z)$

3. DBSCAN是一种基于密度的聚类算法, 若要基于该算法, 对观测点集 $D = \{x_i, i = 1, \dots, N\}$ 进行聚类, 需要提供两个全局参数 $(\varepsilon, \text{MinPts})$, 两个参数的意

义是什么？如何确定核心对象？什么是密度直达？什么是密度可达？什么是密度相连？

解：

(1) 两个**全局邻域参数**($\varepsilon, MinPts$)

ε --邻域最大半径

$MinPts$ --给定样本的 ε -**邻域**内最小样本数.

其中：

ε -**邻域** 对于 $\forall x_j \in D$, x_j 的 ε -邻域为 $N_\varepsilon(x_j) = \{x_i \in D \mid dist(x_i, x_j) \leq \varepsilon\}$

(2) **核心对象**(*core object*)

若 $x_j \in D$ 并且 $|N_\varepsilon(x_j)| \geq MinPts$, 则称 x_j 为一个**核心对象**.

(3) **密度直达**(*directly density-reacheable*)

若 $x_j \in N_\varepsilon(x_i)$, 并且 x_i 为一个核心对象, 则称 x_j 为由 x_i **密度直达**.

(4) **密度可达**(*density-reacheable*)

对于 x_i, x_j , 若存在样本序列 p_1, p_2, \dots, p_n , 其中 $p_1 = x_i, p_n = x_j$, 且 p_{i+1} 由 p_i 密度直达, 则称 x_j 由 x_i **密度可达**.

(5) **密度相连**(*density-connected*)

对于 x_i 与 x_j , 若存在样本 x_k , 使得 x_i 与 x_j 均由 x_k **密度可达**, 则称 x_i 与 x_j **密度相连**