

§9 聚类

1. 聚类评价与性能度量

无监督训练样本的标记信息是未知的。

有效性指标 (validity index)

聚类性能度量

- 外部指标: 将聚类结果与某个参考模型比较
- 内部指标: 直接考察聚类结果而不利用任何参考模型

外部: Jaccard 系数, FM 指数, Rand 指数

内部: DB 指数, Dunn 指数

2. 距离

$$\text{欧式 distink}(x_i, x_j) = \left(\sum_{u=1}^n |x_{iu} - x_{ju}|^p \right)^{\frac{1}{p}}$$

$p=2$ 欧式 $\text{distink}(x_i, x_j) = \|x_i - x_j\|_2 = \sqrt{\sum_{u=1}^n |x_{iu} - x_{ju}|^2}$

distman 曼哈顿 $= \|x_i - x_j\|_1 = \sum_{u=1}^n (x_{iu} - x_{ju})$

3. 原型聚类

(1) k-means $E = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|_2^2$ 最小化误差

思考: 对于给定的类别数 k , 首先给定初始划分, 通过迭代改变样本所属的类别关系, 使每次处理后得到比上次好

步骤: 选择初始中心 k 个, $a_1, a_2, \dots, a_k \rightarrow$ 对于每个样本 x_i , 将其标记为

距离类别中心 a_j 最近的类别 $j = \arg\min_k \sqrt{\sum_{u=1}^n (x_{iu} - a_{ju})^2} \rightarrow$ 更新 a_j 为均值

\rightarrow 重复直到中止。

中止: 迭代次数, 样本到中心距离平方和, 簇中心点变化率

优点: 说明简单, 较快的收敛性和高效率, 高维分布效果不错

缺点: k 未知, 对初始值敏感, outlier 影响大, 不适合发现非凸形状的簇

改进: k-medoids, kmeans++ kmeans||

(2) 迭代聚类 (LVE)

利用样本之任意位置寻找质心

随机选择 p → 计算样本 x_j 与 p 距离 $d_{ij} = \|x_j - p\|_2$ → 找出与 x_j 距离最近的原
点向量 p^* → 更新向量 $p' \rightarrow$ 迭代中止

(3) GMM 常采用 EM 迭代

(4) 密度聚类

思想: 只要样本点密度大于等于某个阈值, 即将该样本加入到最近簇中

适用于任意形状, 对噪声 data 不敏感, 但计算量大

DBSCAN: 用一个点之 ϵ 邻域内的邻居点数量衡量该点所在空间的密度
不需要给定 k

流程: 若点 x 的 ϵ 邻域包含多于 ϵ 个对象, 则创建一个 x 作为新 cluster

→ 寻找并合并核心对象直接密度可达对象, 无新对象可合并时结束

优点: 无需给定 k , 可任意形状 cluster, 可成 outlier 不敏感, 只需 2 个参数

几乎不依赖节点之遍历顺序

缺点: 依赖噪声识别 (噪声) 不适合数据集中密度差异很小的情况

(5) 层次聚类

凝聚: AGNES 自底向上, 两两间距离由距离最近之数据点之相似/度
确定, 合并过程反复至所有对象满足预数目

分裂: DIANA 自顶向下, 按照比例细分为越来越小的簇, 直至 cluster 数中止

优点: 合并/分裂点选择不易, 不可撤回操作, 大数据不适合