

非监督式机器学习

聚类(Clustering)

--算法评价

河北师范大学软件学院

2018.04.03-04.12

评价的意义

- (1) 避免所发现的数据结构源自噪声干扰
- (2) 不同聚类算法的比较
- (3) 两个聚类集合 (**two sets of clusters**) 的比较
- (4) 两个聚类的比较

--> { **聚类趋向：验证给定数据集是否具有聚类结构；**
发现数据中真实的结构

评价的几个角度

➤ 明确给定数据集合中“聚类的趋势”

如：区分给定数据集内是否存在非随机性“结构”

➤ “外部评价”--将聚类分析的结果与给定的结果
(带有类别标签的专门数据)比较

➤ “内部评价”--评估聚类分析的结果是否与数据
结构相符，而无需参考外部信息-只借助数据本身

➤ 比较不同聚类算法的分析结果，以确定哪种聚类
算法更好

➤ 确定正确的“聚类数目”

评价的几种类型 (*types of validation measures*)

(1) 外部评价 (*external validation*)

需要关于研究对象相关领域的先验知识

如：一个预定义的划分

不足 强化了研究者的主观猜测；
 会忽略某些与之前认识不符的现象；
 导致错过发现新规律新模式的机会

(2) 内部评价 (*internal validation*)

基于数据本身内在的信息，量化分析

外部评价的一些常见指标

给定数据集 $\mathbf{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, $\mathbf{x}_i = [\mathbf{x}_{i1} \quad \dots \quad \mathbf{x}_{id}]^T \in \mathbf{R}^d$

若 $\begin{cases} \text{由聚类给出的簇划分结果} & \mathbf{C} = \{\mathbf{C}_1, \dots, \mathbf{C}_k\} \\ \text{参考模型给出的簇划分结果} & \mathbf{C}^* = \{\mathbf{C}_1^*, \dots, \mathbf{C}_s^*\} \end{cases}$

数据集 \mathbf{D} 内各样本相应簇标记值集合 $\begin{cases} \lambda = \{\lambda_1, \dots, \lambda_m\} \\ \lambda^* = \{\lambda_1^*, \dots, \lambda_m^*\} \end{cases}$

数据集 \mathbf{D} 内样本两两配对，定义：



河北师范大学软件学院
Software College of Hebei Normal University

$$a = |\mathbf{SS}| \quad \mathbf{SS} = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i = \lambda_j, \lambda_i^* = \lambda_j^*, i < j\}$$

$$b = |\mathbf{SD}| \quad \mathbf{SD} = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i = \lambda_j, \lambda_i^* \neq \lambda_j^*, i < j\}$$

$$c = |\mathbf{DS}| \quad \mathbf{DS} = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i \neq \lambda_j, \lambda_i^* = \lambda_j^*, i < j\}$$

$$d = |\mathbf{DD}| \quad \mathbf{DD} = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i \neq \lambda_j, \lambda_i^* \neq \lambda_j^*, i < j\}$$

显然 $a + b + c + d = C_m^2 = \frac{m(m-1)}{2}$



河北师范大学软件学院
Software College of Hebei Normal University

基于上述定义，给出用于聚类性能度量的常见**外部指标**

[1] *Jaccard* 系数 (*Jaccard Coefficient*, 简称 *JC*)

$$JC = \frac{a}{a + b + c} \in [0, 1]$$

[2] *FM* 指数 (*Fowlkes and Mallows Index*, 简称 *FMI*)

$$FMI = \sqrt{\frac{a}{a + b} \cdot \frac{a}{a + c}} \in [0, 1]$$

[3] *Rand* 指数 (*Rand Index*, 简称 *RI*)

$$RI = \frac{a + d}{a + b + c + d} = \frac{2(a + d)}{m(m-1)} \in [0, 1]$$

内部评价的一些常见评价指标

给定数据集 $D = \{x_1, \dots, x_m\}$, $x_i = [x_{i1} \ \dots \ x_{id}]^T \in R^d$

若由聚类给出的簇划分结果 $\mathbf{C} = \{C_1, \dots, C_k\}$

并且 $\begin{cases} \text{dist}(\cdot, \cdot) -- \text{两样本点之间距离} \\ \mu = \frac{1}{|C|} \sum_{x \in C} x -- \text{任意簇 } C \in \mathbf{C} \text{ 的中心点.} \end{cases}$

$\forall C \in \mathbf{C}$, 簇 C 内样本间的平均距离



河北师范大学软件学院
Software College of Hebei Normal University

$$avg(C) = \frac{2}{|C|(|C|-1)} \sum_{1 \leq i < j \leq |C|} dist(x_i, x_j)$$

$\forall C \in \mathbf{C}$, 簇 C 内样本间的最远距离

$$diam(C) = \max_{1 \leq i < j \leq |C|} dist(x_i, x_j)$$

簇 C_i, C_j 样本间最近距离 $d_{\min}(C_i, C_j) = \min_{x_i \in C_i, x_j \in C_j} dist(x_i, x_j)$

簇 C_i, C_j 中心点之间距离 $d_{\text{cen}}(C_i, C_j) = dist(\mu_i, \mu_j)$

基于上述定义，给出用于聚类性能度量的常见**内部指标**

[1] **DB**系数 (*Davies – Bouldin Index*, 简称**DBI**)

$$DBI = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left(\frac{\text{avg}(C_i) + \text{avg}(C_j)}{d_{\text{cen}}(C_i, C_j)} \right)$$

DBI值越小越好.

[2] **Dunn**指数 (*Dunn Index*, 简称**DI**)

minimal intercluster distance / maximal intracluster distance.

$$DI = \min_{1 \leq i \leq k} \min_{j \neq i} \left(\frac{d_{\min}(C_i, C_j)}{\max_{1 \leq l \leq k} \text{diam}(C_l)} \right) \quad DI \in [0, \infty)$$

DI值越大越好.



河北师范大学软件学院
Software College of Hebei Normal University