

控制与决策

Control and Decision

基于过滤模型的聚类算法

邱保志, 张瑞霖, 李向丽

引用本文:

邱保志, 张瑞霖, 李向丽. 基于过滤模型的聚类算法[J]. 控制与决策, 2020, 35(5): 1091–1101.

在线阅读 View online: <https://doi.org/10.13195/j.kzyjc.2018.1089>

您可能感兴趣的其他文章

Articles you may be interested in

基于分量属性近邻传播的多元时间序列数据聚类方法

Multivariate time series clustering based on affinity propagation of component attributes

控制与决策. 2018, 33(4): 649–656 <https://doi.org/10.13195/j.kzyjc.2017.0150>

维度概率摘要模型及其层次聚类算法

Hierarchical clustering algorithm with dimensions probability summary model

控制与决策. 2017, 32(8): 1421–1426 <https://doi.org/10.13195/j.kzyjc.2016.0984>

基于平均差异度优选初始聚类中心的改进K-均值聚类算法

Improved K-means clustering algorithm optimizing initial clustering centers based on average difference degree

控制与决策. 2017, 32(4): 759–762 <https://doi.org/10.13195/j.kzyjc.2016.0274>

参数自适应的可变类FLICM灰度图像分割算法

Self-adaptive FLICM algorithm for gray image segmentation with unknown number of clusters

控制与决策. 2017, 32(2): 262–268 <https://doi.org/10.13195/j.kzyjc.2016.0050>

基于自适应学习的演化聚类算法

Evolving clustering method based on self-adaptive learning

控制与决策. 2016(3): 423–428 <https://doi.org/10.13195/j.kzyjc.2014.1945>

基于维度最大熵数据流聚类的异常检测方法

Data stream clustering algorithm based on the maximum entropy of data dimension and its applications for anomaly detection

控制与决策. 2016(2): 343–348 <https://doi.org/10.13195/j.kzyjc.2014.1783>

基于模糊测度和证据理论的模糊聚类集成方法

Fuzzy clustering ensemble based on fuzzy measure and DS evidence theory

控制与决策. 2015, 30(5): 823–830 <https://doi.org/10.13195/j.kzyjc.2014.0358>

一种改进隶属度函数的FCM聚类算法

An FCM clustering algorithm with improved membership function

控制与决策. 2015, 30(12): 2270–2274 <https://doi.org/10.13195/j.kzyjc.2014.1716>

基于过滤模型的聚类算法

邱保志[†], 张瑞霖, 李向丽

(郑州大学 信息工程学院, 郑州 450001)

摘要: 合理的聚类原型是正确聚类的前提. 针对现有聚类算法原型选取不合理、计算聚类个数存在偏差等问题, 提出基于过滤模型的聚类算法(CA-FM). 算法以提出的过滤模型去除干扰聚类过程的边界和噪声对象, 依据核心对象之间的近邻关系生成邻接矩阵, 通过遍历矩阵计算聚类个数; 然后, 按密度因子将数据对象排序, 从中选出聚类原型; 最后, 将其余对象按照距高密度对象的最小距离划分到相应的簇中, 形成最终聚类. 在人工合成数据集、UCI 数据集以及人脸识别数据集上的实验结果验证了算法的有效性, 与同类算法相比, CA-FM 算法具有较高的聚类精度.

关键词: 聚类算法; 过滤模型; 偏差因子; 聚类原型; 局部密度; 密度因子

中图分类号: TP273

文献标志码: A

Clustering algorithm based on filter model

QIU Bao-zhi[†], ZHANG Rui-lin, LI Xiang-li

(School of Information Engineering, Zhengzhou University, Zhengzhou 450001, China)

Abstract: Reasonable clustering prototype is the premise of correct clustering. Most of the existing clustering algorithms have some shortcomings such as the unreasonable selection of clustering prototypes and calculation deviation of cluster numbers. A clustering algorithm based on filter model (CA-FM) is proposed. The algorithm uses the proposed filtering model to remove the boundary and noise objects which interfere with the clustering process. The adjacency matrix is generated according to the neighbor relationships among the core objects, and the number of clusters is calculated by traversing the matrix. Then, the objects are sorted according to the density factor, and clustering prototypes are selected from them. Finally, the remaining objects are assigned into corresponding clusters according to the minimum distance from the high density objects. The effectiveness of the proposed algorithm is demonstrated by experiments on synthetic datasets, UCI datasets and Olivetti face dataset. Compared with similar algorithms, the CA-FM has a higher clustering accuracy.

Keywords: clustering algorithm; filter model; deviation factor; clustering prototype; local density; density factor

0 引言

聚类是一种学习范式^[1], 旨在发现数据的内部结构, 它在数据探索和知识发现中扮演着重要角色. 目前已提出大量的聚类算法, 并在图像分割、生物学、电子商务、互联网等领域得到广泛应用^[2-7], 如均值聚类算法(K -means)^[8]、模糊均值聚类算法(fuzzy C -means clustering algorithm, FCM)^[9]、基于密度和噪音的空间聚类算法(density-based spatial clustering of applications with noise, DBSCAN)^[10]、密度峰值聚类算法(clustering by fast search and find of density peaks, DPC)^[11]、无参数拉普拉斯中心性聚类算法(parameter-free Laplacian centrality peaks clustering, LPC)^[12]、

优化密度峰值聚类算法(comparative density peaks clustering, CDP)^[13]等. 许多算法是通过寻找聚类骨架完成聚类, 而寻找聚类骨架的关键在于确定聚类原型集. 选取正确的聚类原型可以提高聚类精度, 如何合理地确定聚类原型集已成为聚类算法亟待解决的问题.

以 K -means、FCM 为代表的基于划分的聚类算法将随机选取的 k 个对象作为聚类的初始原型, 按照相似性原则将数据对象分配给相应的原型形成一个簇, 通过反复计算每个簇的原型和再分配, 直至目标函数收敛. 这一机制决定了这一类算法不能有效地处理非球形簇, 且聚类精度不高.

收稿日期: 2018-08-08; 修回日期: 2018-11-27.

基金项目: 河南省基础与前沿技术研究项目(152300410191).

责任编辑: 阳春华.

[†]通讯作者. E-mail: iebzqiu@zzu.edu.cn.

以DBSCAN为代表的基于密度的聚类算法将核心点作为聚类原型,寻找与聚类原型密度可达的对象,形成聚类.它可以发现任意形状的聚类,对噪声具有很好的鲁棒性,但由于算法采用了固定大小的邻域计算密度,使得这一类算法的聚类结果对输入参数敏感,且不能有效处理高维和多密度数据集.

DPC、LPC、CDP等算法通过计算对象密度,选取决策图中密度峰值对象作为聚类原型;依据距峰值最小距离原则,将其余数据对象划分到相应聚类.虽然算法结构简单,易于理解,但DPC、CDP算法的密度度量方式依赖于手动输入的截断距离参数,不合理的参数设置会引起对象划分错误的连锁反应.LPC算法借鉴谱聚类思想,将数据集视为无向图,采用拉普拉斯中心性^[14]表征数据对象的密度,由于LPC算法在提取每一维度的特征值形成拉普拉斯矩阵的计算量较大,对于高维数据集,算法运行时间会指数倍增加,无法适应高维数据聚类要求.

为了解决上述问题,本文以提出的偏差因子建立非核心对象过滤模型,用于过滤掉那些影响聚类原型选取的噪声和聚类边界对象;然后基于提出的原型选取机制自动确定聚类原型;最后,将其余对象分配到各个原型所属的簇中,形成聚类.本文的创新点如下:1)建立非核心对象过滤模型;2)提出一种局部密度计算方法;3)提出一种聚类原型自动选取的机制.

1 基于过滤模型的聚类算法

1.1 相关定义

数据采样可分为静态采样和动态采样.超球采样、立方采样和网格采样^[6]属于静态采样,采样思想是以固定大小的邻域或超立方体中包含的数据对象个数来衡量一个对象的密度^[15].若对象分布不均或对于分布稀疏的高维空间,难以设置合适的邻域半径或边长,造成对象的局部分布特征度量不准确,导致密度度量失衡.

相对于静态采样,动态采样可以更好地反映对象的分布特征, k 近邻采样是动态采样,它提取 k 个最近的数据对象形成动态采样空间,可以更好地表征数据的局部密度分布.设数据集 D 含有 m 个属性, k 近邻的定义如下.

定义1 设 $x \in D$, x 的 k 近邻^[6]是距离 x 最近的 k 个对象的集合,用 $\text{knn}(x)$ 表示,即

$$\begin{aligned} \text{knn}(x) = \{ & x_1, x_2, \dots, x_k | x_i \in D, 1 \leq i \leq k, \\ & \forall y \in D - \{x_1, x_2, \dots, x_k\}, \\ & \text{dist}(x, y) \geq \text{dist}(x, x_i) \}. \end{aligned} \quad (1)$$

定义2 设 $x \in D$, x 的 k 近邻距离^[6] $\text{knn}_d(x)$ 定义为 x 到 k 近邻对象的距离之和,即

$$\text{knn}_d(x) = \sum_{y \in \text{knn}(x)} \text{dist}(x, y). \quad (2)$$

其中: dist 为欧氏距离, $\text{knn}_d(x)$ 反映了 x 周围的分布情况,其值越小,说明 x 周围分布越稠密,反之越稀疏,但对象的 k 近邻距离对参数 k 较为敏感,若 k 值选取过小,则采样不充分,无法表征对象的真实分布; k 值选取过大,任意两个对象的 k 近邻距离几乎一致.为了合理地度量对象的周围分布,本文提出 k 近邻共享度的定义.

定义3 设 $x \in D$,对象 x 的 k 近邻共享度是对象 x 的 k 近邻与每个近邻对象的 k 近邻拥有相同对象个数之和,记为 $\text{knns}(x)$,即

$$\text{knns}(x) = \sum_{y \in \text{knn}(x)} |\text{knn}(x) \cap \text{knn}(y)|, \quad (3)$$

其中 $\text{knns}(x)$ 的取值范围为 $[0, (k-1)^2]$. $\text{knns}(x)$ 值越大,说明 x 周围数据对象分布越稠密; $\text{knns}(n)$ 值越小,说明 x 的近邻与 x 的分布差异越大.

核心对象周围分布较为稠密,其 k 近邻距离较小, k 近邻共享度较大;而非核心对象周围分布较为稀疏,其 k 近邻共享度较小.为了增大核心对象与非核心对象的差异,结合 k 近邻距离与 k 近邻共享度,本文给出了局部密度的定义.

定义4 设 $x \in D$,对象 x 的局部密度 $\text{den}(x)$ 定义如下:

$$\text{den}(x) = \frac{\text{knns}(x)}{\text{knn}_d(x)}. \quad (4)$$

DPC算法认为聚类中心具有较高的密度且距高密度对象的最小距离较远^[11].为了方便聚类中心的选取,这里使用密度因子的概念放大聚类中心与其他对象的特征差异.

定义5 设 $x \in D$,数据对象 x 的密度因子 $R(x)$ 定义为

$$R(x) = \text{den}(x)\delta(x). \quad (5)$$

其中: $\delta(x)$ 表示 x 与高密度数据对象之间的最小距离,即

$$\delta(x) = \begin{cases} \min\{\text{dist}(x, y) | \text{den}(x) < \text{den}(y), y \in D\}, \\ \exists y \in D, \text{den}(x) < \text{den}(y); \\ \max\{\text{dist}(x, y) | y \in D\}, \text{ otherwise.} \end{cases} \quad (6)$$

一个对象的密度因子越大,表明对象成为聚类中心的可能性越大.对象按其密度因子降序排序后,聚类中心一定位于序列的前半部分.

1.2 核心对象的获取

在基于密度的聚类中,数据对象通常划分为核心对象与非核心对象^[16],核心对象构成聚类的骨架,而非核心对象(如噪音、边界)可能会干扰聚类过程^[17].为此,首先建立非核心对象过滤模型(FM)去除边界和噪声对象,得到聚类的核心对象集,然后在核心对象集上获得聚类原型.

定义6 设 $x \in D, y \in \text{knn}(x)$, 对象 x 的 k nn 质心 $\text{knn_centroid}(x)$ 定义为

$$\text{knn_centroid}(x) = \left(\sum_{y \in \text{knn}(x)} \frac{y_1}{k}, \sum_{y \in \text{knn}(x)} \frac{y_2}{k}, \dots, \sum_{y \in \text{knn}(x)} \frac{y_i}{k}, \dots, \sum_{y \in \text{knn}(x)} \frac{y_m}{k} \right), \quad (7)$$

其中 y_i 表示对象 y 在第 i 维上的取值,即 $y = (y_1, y_2, \dots, y_i, \dots, y_m)$.

对象 x 的 k nn 质心反映了 x 周围分布情况. 如果 x 是核心对象,则其 k 近邻相对均匀地分布在 x 的周围;若 x 是非核心对象,则其 k 近邻分布具有较大的偏向性,其 k 近邻聚集在某一方向. 为了反映这一特征,本文将对象 x 与 k nn 质心之间的距离定义为 x 的偏差因子,依据偏差因子过滤非核心对象.

定义7 设 $x \in D$, x 的偏差因子 $d_f(x)$ 定义为 x 到其 k 近邻质心的距离,即

$$d_f(x) = \text{dist}(\text{knn_centroid}(x), x). \quad (8)$$

核心对象的周围分布较为均匀,其偏差因子较小;非核心对象周围分布具有较强的偏向性且稀疏,其偏差因子较大. 将数据对象按偏差因子降序排列,得到序列 d_f_desc , 则聚类的非核心对象和核心对象分别位于该序列的前半部分和后半部分.

设 α 为过滤因子(偏差因子降序排列后的百分位数). 对于数据集 D , 过滤模型 FM 将其分为非核心对象集 non_core_set 与核心对象集 core_set , 定义如下:

$$\begin{cases} \text{core_set} = \\ \{x|x \in D, \text{ and } d_f(x) \leq d_f_desc(\lceil |D| * \alpha \rceil)\}; \\ \text{non_core_set} = \\ \{x|x \in D \text{ and } d_f(x) > d_f_desc(\lceil |D| * \alpha \rceil)\}. \end{cases} \quad (9)$$

其中: d_f_desc 为偏差因子的降序排序, $\lceil \cdot \rceil$ 表示向上取整.

以合成数据集 Syn 为例,该数据集共有 5324 个对象,含有 5 个大小不同的簇,簇周围存在噪声,垂直方向的两个椭圆簇之间存在桥接现象. 图 1 显示了聚类核心对象的获取过程,图 1(b) 表示使用过滤模型得到的边界和噪声对象. 图 1(c) 表示过滤后得到的核心对象.

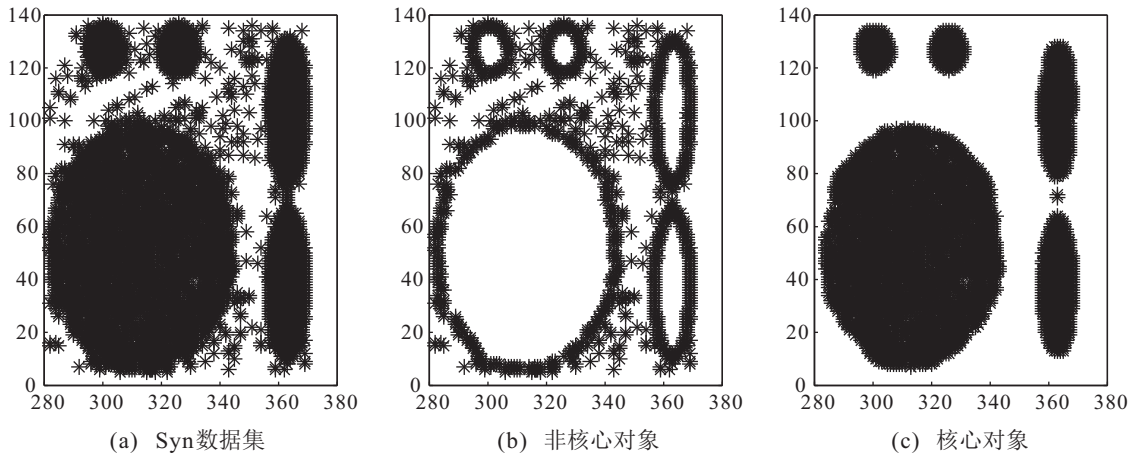


图 1 核心点获取过程 ($\alpha = 19.84$)

1.3 聚类原型选取

自动且准确地获取聚类原型将提高聚类算法的适用性与效率,Chen 等^[18]采用证据积累的思想,迭代运行 FCM 算法,形成累积邻接矩阵,通过图切分的方式获取原型个数,由于 FCM 算法将最小误差平方和作为目标函数,导致部分隶属度矩阵存在偏差,使获取的原型个数并不稳定. Zhang 等^[19]将不同 k 值对应的最小误差平方和作为决策图,将决策图中拐点

所对应的 k 值作为原型个数,对于非球形簇,此算法得到的聚类原型与实际原型存在一定的偏差. DC-MDACC^[20]算法利用残差分析与线性回归,将未在置信区间中的对象视为聚类原型,但原型个数极易受到置信因子参数的影响. DPC 与 LPC 算法在决策图中采用人工方式筛选聚类原型,即算法在选取过程中需要人工参与.

由于过滤了干扰聚类过程的非核心对象,核心对

象集 $core_set$ 可以较好地表征数据集的分布特征. 本文利用非核心对象过滤模型获取核心对象集; 建立核心对象之间的近邻关系, 形成邻接矩阵; 采用遍历算法计算子图个数, 并结合密度因子得到聚类原型.

定义 8 设 $x, y \in core_set$, 邻接矩阵 $Conn_M$ 中元素定义如下:

$$Conn_M(x, y) = \begin{cases} 1, & x \in knn(y); \\ 0, & x \notin knn(y). \end{cases} \quad (10)$$

其中 $Conn_M$ 为一个非对称的 $n \times n$ 矩阵. 通过建立近邻关系, 将数据集的核心部分转换为独立的连通子图, 每个子图代表原始数据集中的一个聚类, 形成邻接矩阵 $Conn_M$, 然后采用广度优先搜索的方式, 得到邻接矩阵中子图的个数 I , 即为聚类原型个数.

由式 (5) 可知, 密度因子 R 越大的数据对象, 作为聚类中心的可能性越大. 依据密度因子进行降序排序后得到序列 R_desc , 聚类原型一定在该序列的前半部分, 则序列 R_desc 的前 I 个数据对象即为聚类原型, 即 $Center_set = R_desc[1 : I]$.

1.4 基于过滤模型的聚类算法

基于过滤模型的聚类 (clustering algorithm based on filter model, CA-FM) 算法包含计算、核心点获取、原型获取和分派 4 个步骤. 计算步骤包含计算每个对象的局部密度、偏差因子、 α 距离、密度因子; 核心点获取步骤是依据过滤模型, 过滤那些影响聚类原型提取的非核心对象, 从而得到核心对象集; 原型获取步骤是根据核心集中对象的近邻关系构建邻接矩阵, 采用图搜索方式计算子图个数, 并使用密度因子获取聚类原型; 分派步骤将剩余数据对象划分至相应的簇中形成最终聚类. 详细步骤如下.

算法 1 CA-FM 算法.

输入: 数据集 $data$ 、近邻参数 k 、过滤因子 α ;

输出: 数据集的聚类标签 $Label$.

step 1: 计算. 根据式 (4)、(5)、(6)、(8) 计算数据对象的局部密度 den 、密度因子 R 、距离 δ 、偏差因子 d_f .

step 2: 核心对象获取.

step 2.1: 依据偏差因子 d_f 将对象降序排序, 得到 d_f_desc 序列;

step 2.2: 选取 d_f_desc 集合中位置 $[|D| * \alpha]$ 之后的对象作为核心对象, 形成核心对象集 $core_set$.

step 3: 原型获取.

step 3.1: 根据式 (10), 确定核心对象之间的近邻关系, 构建邻接矩阵 $Conn_M$;

step 3.2: 采用遍历算法搜索邻接矩阵, 计算邻接矩阵中连通子图个数, 得到聚类中心个数 I ;

step 3.3: 选取密度因子降序排序后前 I 个对象作为聚类中心.

step 4: 分派. 将剩余数据对象划分至距离最近的高密度对象所在的簇中, 返回聚类标签 $Label$.

2 实验结果与分析

实验环境: 内存为 4.00 GB, 操作系统为 Microsoft Windows 7, 编译环境为 Matlab R2014a.

数据集包括合成数据集、UCI 数据集和人脸识别数据集^[11]. 详细信息见表 1, 编号 1~编号 6 为合成数据集, 用来检验算法在不同数据分布形态下的聚类效果; 编号 7~编号 16 是 UCI 数据集, 用来检测算法在真实数据下的聚类效果; 编号 17 为人脸识别数据集, 用来检测算法在高维数据下的检测效果. 算法使用标准化互信息 (normalized mutual information, NMI)^[21]、准确率 (accuracy, ACC)^[21]、纯度 (Purify)^[22]、兰德指数 (rand index, RI)^[23]、FM 指数 (fowlkes and mallows index, FMI)^[24]、杰卡德相似系数 (jaccard similarity coefficient, JC)^[24] 聚类评价指标从多角度衡量聚类质量.

表 1 数据集的基本信息

NO.	datasets	data Sources	m	class	instance
1	Compound ^[25]	synthesis	2	6	399
2	R15 ^[25]	synthesis	2	15	600
3	Spiral ^[25]	synthesis	2	3	312
4	Aggregation ^[25]	synthesis	2	7	788
5	Jain ^[25]	synthesis	2	2	373
6	4k2-far ^[25]	synthesis	2	4	400
7	Wine ^[26]	UCI	13	3	178
8	Sonar ^[27]	UCI	60	2	208
9	Soybean ^[27]	UCI	35	4	47
10	Zoo ^[27]	UCI	16	7	101
11	Parkinson ^[2]	UCI	22	2	195
12	Glass ^[26]	UCI	9	7	214
13	Iris ^[26]	UCI	4	3	150
14	Tic-Tac-Toe ^[26]	UCI	10	2	958
15	Mushroom ^[21]	UCI	22	2	8 124
16	German Credit ^[20]	UCI	20	2	1 000
17	Face dataset ^[11]	ORL	10 304	10	100

在对比实验中, 将 K -means、FCM、CDP 算法的参数: 原型个数设置为正确的聚类个数, 各运行 10 次, 取各聚类指标的均值作为最终的聚类效果, 其他算法则使用最优的聚类效果作为最终的结果.

2.1 人工合成数据集

实验选取的人工合成数据集包含了以下形态: 多密度、流型螺旋、多形状、微型簇、簇间半包含以及簇间嵌套, 并含有随机噪声或桥接噪声. Compound 数据集共有 6 个密度分布不均匀的聚类, 且聚类之间

存在嵌套分布,用以测试算法能否准确识别多密度聚类和嵌套的聚类. R15数据集包含15个微型簇,用来检测算法是否能完整识别数据集中所有聚类. Spiral、

Jain数据集为流型簇,簇之间为半包含关系,用来检测算法是否可以识别任意形状的聚类. Aggregation数据集共有7个形状不同的聚类,簇之间存在桥接噪

表2 人工合成数据集上的聚类结果比较

数据集	算法	参数	ACC/%	NMI	Purify	JC	RI	FMI
Compound	<i>K</i> -means	$k = 6$	62.6566	0.7149	0.8321	0.4549	0.8406	0.6337
	DBSCAN	MinPts = 4, EPS = 2.2	78.4461	0.8053	0.7845	0.5552	0.9341	0.7324
	FCM	$k = 6$	65.6642	0.7103	0.8321	0.4627	0.8427	0.6404
	CDP	$k = 6, d_c = 0.4$	65.9148	0.7407	0.8321	0.4808	0.8442	0.6535
	DPC	$d_c = 1.25$	63.1579	0.7589	0.9321	0.4529	0.8266	0.6251
	LPC		77.4436	0.7679	0.7744	0.4968	0.9304	0.6700
	CA-FM	$k = 10, \alpha = 0.0275$	83.2080	0.8594	0.8321	0.6015	0.8977	0.7748
Spiral	<i>K</i> -means	$k = 3$	34.6154	0.00005	0.3494	0.1960	0.5540	0.3278
	DBSCAN	MinPts = 10, EPS = 1	100	1.0000	1.0000	1.0000	1.0000	1.0000
	FCM	$k = 3$	33.9744	0.0002	0.3429	0.1955	0.5541	0.3272
	CDP	$k = 3, d_c = 0.24$	100	1.0000	1.0000	1.0000	1.0000	1.0000
	DPC	$d_c = 1.7443$	100	1.0000	1.0000	1.0000	1.0000	1.0000
	LPC		100	1.0000	1.0000	1.0000	1.0000	1.0000
	CA-FM	$k = 2, \alpha = 0.05$	100	1.0000	1.0000	1.0000	1.0000	1.0000
Aggregation	<i>K</i> -means	$k = 7$	73.3503	0.8036	0.8883	0.5676	0.8958	0.7321
	DBSCAN	MinPts = 4, EPS = 0.83	82.7411	0.8894	0.8274	1.0000	1.0000	1.0000
	FCM	$k = 7$	79.6954	0.8427	0.9315	0.6433	0.9187	0.7926
	CDP	$k = 7, d_c = 0.03$	87.6904	0.8756	0.8858	0.7625	0.9374	0.8673
	DPC	$d_c = 1$	94.0355	0.9705	0.9403	0.9591	0.9911	0.9793
	LPC		98.7310	0.9700	0.9873	0.9599	0.9912	0.9796
	CA-FM	$k = 5, \alpha = 0.3807$	99.6193	0.9896	0.9962	0.9898	0.9850	0.9949
R15	<i>K</i> -means	$k = 15$	79.5000	0.8989	0.7950	0.6075	0.9606	0.7704
	DBSCAN	MinPts = 5, EPS = 0.32	78.1667	0.9121	0.7850	0.5927	0.9627	0.7642
	FCM	$k = 15$	99.6667	0.9942	0.9967	0.9866	0.9991	0.9932
	CDP	$k = 15, d_c = 0.24$	100	1.0000	1.0000	1.0000	1.0000	1.0000
	DPC	$d_c = 0.95$	99.5000	0.9922	0.9950	0.9801	0.9987	0.9900
	LPC		92.1667	0.9410	0.9217	0.7915	0.9851	0.8846
	CA-FM	$k = 10, \alpha = 0.3333$	100	1.0000	1.0000	1.0000	1.0000	1.0000
Jain	<i>K</i> -means	$k = 2$	78.5523	0.3690	0.7855	0.5348	0.6621	0.7005
	DBSCAN	MinPts = 4, EPS = 2.8	73.9946	0.0001	0.7399	0.5000	0.4530	0.7071
	FCM	$k = 2$	77.4799	0.3555	0.7748	0.5218	0.6501	0.6894
	CDP	$k = 2, d_c = 0.01$	86.0590	0.5052	0.8606	0.6497	0.7594	0.7904
	DPC	$d_c = 0.95$	86.0590	0.5052	0.8606	0.6497	0.7594	0.7904
	LPC		89.2761	0.5709	0.8928	0.7135	0.8080	0.8348
	CA-FM	$k = 9, \alpha = 0.0268$	92.2252	0.6447	0.9223	0.7804	0.8120	0.8779
4k2-far	<i>K</i> -means	$k = 4$	100	1.0000	1.0000	1.0000	1.0000	1.0000
	DBSCAN	MinPts = 4, EPS = 0.5	99.7500	0.9941	1.0000	0.9940	0.9982	0.9970
	FCM	$k = 4$	100	1.0000	1.0000	1.0000	1.0000	1.0000
	CDP	$k = 3, d_c = 0.01$	87.7500	0.8430	0.8775	0.7287	0.9085	0.8441
	DPC	$d_c = 0.2168$	83.2500	0.9081	1.0000	0.7749	0.9312	0.8803
	LPC		87.2500	0.8745	0.8725	1.0000	1.0000	1.0000
	CA-FM	$k = 10, \alpha = 0.15$	100	1.0000	1.0000	1.0000	1.0000	1.0000

声,用来检测算法能否处理带有桥接干扰的聚类. 4k2-far数据集含有4个簇和随机噪声,用来测试算法能否在噪声干扰下准确聚类.

表2给出了各算法的聚类结果评价指标值,表中加粗的数据是聚类指标最好的情况. FCM、 K -means算法均采用最小误差平方和作为迭代的目标函数,因此二者均不能有效地处理非球形簇,如Spiral、Jain数据集.但FCM算法引入了模糊理论^[26],其聚类效果好于 K -means的聚类结果. DBSCAN算法虽然可以处理非球形簇,但DBSCAN算法采用固定邻域来度量密度,算法无法有效地处理多密度数据集,所以在Jain数据集上聚类效果较差.

DPC算法的聚类效果整体优于 K -means、FCM、DBSCAN聚类算法,但DPC算法的中心选取过程依赖于数据对象在决策图中的分布,不同的截断参数会产生不同形状的决策图,可能会造成聚类中心与其他对象之间的特征差距变小,难以选取正确的聚类中心,最终影响聚类结果,如算法在4k2-far数据集上聚类效果不佳.

LPC算法聚类效果好于DPC算法,虽然LPC不需要显式的输入参数,但仍需要根据实际的聚类个数和决策图来选择聚类中心.

CDP算法依据距高密度点最小距离与距低密度

点最小距离的差,将决策图中聚类中心与其他点进一步分离,便于选出聚类中心,但算法需要输入聚类个数,并且由于采用了截断距离 d_c 进行密度计算,导致其处理多密度数据的能力不强,如在Compound、Jain数据集上聚类效果不理想.

CA-FM算法在6个人工合成数据集中的4个均达到最佳聚类效果,在Compound数据集与Aggregation数据集的聚类效果大部分达到最佳,其余的聚类指标取值与最佳的聚类效果相差很小,说明CA-FM算法在处理多种数据分布的聚类时是有效的.

2.2 UCI数据集

实验选用了来自医疗、生物工程、地质勘探、化学等领域的UCI数据集,用来检验在真实、高维数据下的聚类效果. 其中Iris、Soybean、Zoo、Mushroom数据集来自生物工程领域,包含了不同动植物的各种特征; Parkinson数据集来自医疗领域,记录了病患的多种生理指标; Glass数据集来自化学领域,记录了玻璃的不同化学成分; Sonar数据集来自地质勘探领域,包含了不同物体的声呐强度; Tic-Tac-Toe数据集来自游戏博弈领域; German Credit来自金融信贷领域,记录了用户的信用情况. 表3给出了各算法在UCI数据集上的聚类结果评价指标值.

表3 UCI数据集上的聚类结果比较

数据集	算法	参数	ACC/%	NMI	Purify	JC	RI	FMI
Parkinson	K -means	$k = 2$	72.3077	0.0000	0.7538	0.5792	0.5975	0.7444
	DBSCAN	MinPts = 5, EPS = 10	50.2564	0.0904	0.7446	0.3320	0.4779	0.5093
	FCM	$k = 2$	71.7949	0.1037	0.7538	0.4810	0.5929	0.6516
	CDP	$k = 2, d_c = 0.34$	65.6410	0.0140	0.7538	0.4552	0.5466	0.6260
	DPC	$d_c = 9.8975$	73.3333	0.0533	0.7538	0.5321	0.6069	0.6948
	LPC		74.8718	0.0049	0.7538	0.6208	0.6218	0.7861
	CA-FM	$k = 15, \alpha = 0.1026$	75.3846	0.0001	0.7538	0.7000	0.9404	0.7971
Sonar	K -means	$k = 2$	54.3264	0.0068	0.5433	0.3374	0.5013	0.5046
	DBSCAN	MinPts = 1, EPS = 0.5	35.0962	0.2075	0.4875	0.2922	0.5033	0.4546
	FCM	$k = 2$	55.2885	0.0088	0.5529	0.3358	0.5032	0.5028
	CDP	$k = 2, d_c = 0.44$	50.0000	0.0517	0.5337	0.4812	0.4976	0.4818
	DPC	$d_c = 0.8186$	50.9615	0.0011	0.5337	0.3450	0.4978	0.5133
	LPC		52.4038	0.0000	0.5337	0.3469	0.4987	0.5154
	CA-FM	$k = 8, \alpha = 0.2403$	58.1731	0.0228	0.5817	0.3473	0.5110	0.5156
Soybean	K -means	$k = 4$	48.9362	0.5293	0.5745	0.3472	0.7391	0.5167
	DBSCAN	MinPts = 8, EPS = 4	78.7234	0.8377	0.7872	1.0000	1.0000	1.0000
	FCM	$k = 4$	72.3404	0.7158	0.7872	0.4888	0.8316	0.6568
	CDP	$k = 4, d_c = 0.34$	63.8293	0.6376	0.8085	0.3739	0.8002	0.5500
	DPC	$d_c = 0.667$	65.9574	0.6656	0.6596	0.4286	0.7817	0.6020
	LPC		70.2128	0.7786	0.8936	0.5213	0.8649	0.6947
	CA-FM	$k = 5, \alpha = 0.4255$	78.7234	0.8377	0.7872	1.0000	1.0000	1.0000

表3 (续)

数据集	算法	参数	ACC/%	NMI	Purify	JC	RI	FMI
Zoo	<i>K</i> -means	$k = 4$	74.468 1	0.759 0	0.787 2	0.564 4	0.827 9	0.734 7
	DBSCAN	MinPts = 4, EPS = 2	72.277 2	0.601 8	0.722 8	0.774 8	0.943 8	0.873 3
	FCM	$k = 4$	65.346 5	0.715 4	0.821 8	0.507 2	0.864 2	0.677 9
	CDP	$k = 4, d_c = 0.34$	62.376 2	0.708 0	0.821 8	0.405 8	0.981 5	0.586 5
	DPC	$d_c = 0.12$	77.227 7	0.808 6	0.910 9	0.578 0	0.882 2	0.733 8
	LPC		75.247 5	0.725 6	0.831 7	0.582 0	0.859 6	0.741 4
	CA-FM	$k = 10, \alpha = 0.396 0$	78.217 8	0.737 0	0.782 2	0.774 8	0.993 8	0.873 3
Wine	<i>K</i> -means	$k = 4$	58.427 0	0.380 4	0.704 7	0.344 9	0.703 2	0.516 0
	DBSCAN	MinPts = 2, EPS = 1.3	38.202 2	0.026 8	0.398 9	0.486 4	0.732 4	0.688 8
	FCM	$k = 4$	65.730 3	0.407 3	0.657 3	0.695 7	0.903 4	0.820 6
	CDP	$k = 4, d_c = 0.14$	56.176 8	0.313 8	0.606 7	0.330 0	0.621 2	0.498 8
	DPC	$d_c = 17.147 2$	52.808 9	0.393 9	0.646 1	0.398 3	0.610 4	0.588 9
	LPC		66.853 9	0.401 7	0.674 2	0.473 9	0.700 9	0.655 4
	CA-FM	$k = 9, \alpha = 0.337 0$	70.786 5	0.419 3	0.707 9	0.411 7	0.719 0	0.583 2
Tic-Tac-Toe	<i>K</i> -means	$k = 2$	50.939 5	0.000 0	0.653 4	0.352 9	0.499 7	0.522 2
	DBSCAN	MinPts = 4, EPS = 9	64.509 4	0.000 0	0.653 4	0.528 5	0.541 6	0.717 0
	FCM	$k = 2$	51.878 9	0.001 0	0.653 4	0.353 4	0.500 2	0.522 8
	CDP	$k = 2, d_c = 1.1$	64.509 4	0.005 2	0.653 4	0.528 5	0.541 6	0.707 1
	DPC	$d_c = 10$	59.185 8	0.004 1	0.653 4	0.494 8	0.516 4	0.681 3
	LPC		57.933 2	0.055 2	0.018 6	0.467 5	0.512 1	0.648 6
	CA-FM	$k = 16, \alpha = 0.100 2$	70.041 8	0.069 4	0.700 4	0.580 8	0.579 9	0.651 9
Iris	<i>K</i> -means	$k = 3$	89.333 3	0.758 2	0.893 3	0.695 9	0.879 7	0.820 8
	DBSCAN	MinPts = 10, EPS = 0.5	68.666 7	0.604 4	0.686 7	0.537 5	0.771 9	0.705 4
	FCM	$k = 3$	89.333 3	0.749 6	0.893 3	0.694 3	0.879 7	0.819 7
	CDP	$k = 3, d_c = 0.24$	72.000 0	0.635 7	0.720 0	0.549 3	0.782 0	0.714 5
	DPC	$d_c = 0.316 2$	90.666 7	0.805 7	0.906 7	0.724 8	0.892 3	0.840 7
	LPC		69.333 3	0.709 8	0.693 3	0.584 3	0.777 7	0.756 7
	CA-FM	$k = 10, \alpha = 0.333 3$	96.000 0	0.870 5	0.960 0	0.857 8	0.949 5	0.923 4
Glass	<i>K</i> -means	$k = 7$	50.000 0	0.362 1	0.500 0	0.339 3	0.575 0	0.552 1
	DBSCAN	MinPts = 5, EPS = 0.7	45.327 1	0.306 8	0.495 3	0.362 5	0.679 3	0.559 4
	FCM	$k = 7$	49.532 7	0.336 9	0.630 8	0.234 0	0.726 6	0.385 4
	CDP	$k = 7, d_c = 0.29$	48.130 8	0.373 1	0.579 4	0.255 6	0.684 4	0.407 2
	DPC	$d_c = 0.378 9$	48.130 8	0.370 8	0.495 3	0.333 1	0.543 7	0.553 8
	LPC		48.130 8	0.287 0	0.490 7	0.260 7	0.621 9	0.421 7
	CA-FM	$k = 6, \alpha = 0.140 1$	53.271 0	0.383 9	0.532 7	0.384 7	0.700 1	0.582 1
Mushroom	<i>K</i> -means	$k = 2$	64.598 7	0.071 2	0.646 0	0.403 5	0.542 6	0.576 4
	DBSCAN	MinPts = 20, EPS = 30	64.598 7	0.071 2	0.646 0	0.403 5	0.542 6	0.576 4
	FCM	$k = 2$	64.598 7	0.071 2	0.646 0	0.403 5	0.542 6	0.576 4
	CDP	$k = 2, d_c = 0.79$	64.598 7	0.071 2	0.646 0	0.403 5	0.542 6	0.576 4
	DPC	$d_c = 1.93$	60.462 8	0.083 1	0.646 0	0.371 6	0.532 8	0.581 3
	LPC		64.327 9	0.068 3	0.642 3	0.402 0	0.541 0	0.574 8
	CA-FM	$k = 100, \alpha = 0.384 5$	68.746 9	0.110 4	0.687 5	0.414 8	0.570 2	0.586 8
German Credit	<i>K</i> -means	$k = 2$	67.100 0	0.012 0	0.700 0	0.490 6	0.558 0	0.661 8
	DBSCAN	MinPts = 8, EPS = 15.9	57.500 0	0.004 7	0.700 0	0.420 3	0.514 5	0.592 1
	FCM	$k = 2$	67.000 0	0.012 9	0.700 0	0.486 5	0.557 4	0.657 6
	CDP	$k = 2, d_c = 1.59$	60.600 0	0.007 6	0.700 0	0.450 1	0.522 0	0.576 8
	DPC	$d_c = 7.700 0$	67.300 0	0.007 4	0.700 0	0.531 2	0.568 0	0.705 1
	LPC		68.400 0	0.001 5	0.700 0	0.558 0	0.567 3	0.706 8
	CA-FM	$k = 15, \alpha = 0.2$	78.000 0	0.235 6	0.700 0	0.500 0	0.510 2	0.707 1

真实数据集的维度普遍较高,其数据对象的分布并不严格符合统计规律,且高维空间中的数据形状难以确定,并非传统的球形簇,因此适用于球形簇的FCM、 K -means算法的聚类效果不佳,如 K -means在Soybean上效果较差,FCM在Zoo上效果较差.由于高维空间上数据分布稀疏,DBSCAN算法的半径参数难以确定,算法在Iris、Wine、German Credit数据集上效果较差.

由于采用了在决策图中人工选择聚类原型的方式,DPC、LPC、CDP算法在大部分数据集上均可有效地聚类.LPC算法采用拉普拉斯中心性作为密度度量,其度量原理依据统计分布规律,对于高维数据集,密度度量会产生偏差,进而影响聚类结果,如算法在Iris、Glass数据集上效果不佳.DPC、CDP算法由于截断距离 d_c 的设置不合理,导致部分对象在决策图中的位置重合或者过于接近,干扰了聚类中心的选取,如CDP算法在Zoo、German Credit数据集上效果一般,DPC算法在Wine、Mushroom数据集上效果一般.CA-FM算法在各个数据集上的聚类指标大部分达到最佳,其余指标均与最佳指标相差不大.说明CA-FM算法在处理真实数据的聚类时是有效的.

2.3 高维数据集

本文使用人脸识别数据集检测算法在高维数据下的聚类效果. ORL人脸数据库(olivetti research laboratory)来自剑桥Olivetti实验室.共有40个不同年龄、不同性别和不同种族的对象.每个人有10幅图像,图像尺寸是 92×112 .从人脸数据集中随机抽取10

个人物的图像,构成测试数据集Face dataset,共计100张图片,其中人脸部分表情和细节均有变化,例如笑与不笑、眼睛睁着或闭着、戴或不戴眼镜等.数据集预处理方式为:将每张图像的像素矩阵转化成一维矩阵,即每张图像的维度为10304,形成高维数据集.

算法的聚类结果如图2所示,每一行代表识别的不同聚类,每个聚类中亮度最高的图像表示算法识别的聚类原型.详细聚类评价指标值如表4所示,可以看出,CA-FM算法对高维数据聚类是有效的.



图2 算法在人脸识别数据集上的聚类结果

表4 人脸识别数据集上的聚类结果比较

数据集	算法	参数	ACC/%	NMI	Purify	JC	RI	FMI
Face dataset	K -means	$k = 10$	65.0000	0.7983	0.6500	0.5049	0.9491	0.6822
	DBSCAN	MinPts = 1, EPS = 6	10.0000	0.0000	0.1000	0.5000	0.9909	0.7071
	FCM	$k = 10$	41.0000	0.5122	0.4100	0.2060	0.7766	0.3866
	CDP	$k = 10, d_c = 0.34$	88.0000	0.9582	0.9000	0.7887	0.9766	0.8851
	DPC	$d_c = 3.0288e^3$	81.0000	0.9188	0.8100	0.8000	0.9800	0.8926
	LPC		36.0000	0.6178	0.4300	0.2445	0.7859	0.4498
	CA-FM	$k = 5, \alpha = 0$	100.0000	1.0000	1.0000	1.0000	1.0000	1.0000

2.4 算法分析

2.4.1 时间复杂度分析

CA-FM算法的计算开销主要有:计算偏差因子、计算局部密度、计算 α 距离、建立邻接矩阵、计算子图个数和对象划分.

算法采用KD树建立数据集的 k 近邻关系^[2],则

计算偏差因子的时间复杂度为 $O(kn \log_2 n)$,计算数据对象距高密度点最小距离的时间复杂度为 $O(n \log_2 n + n)$,计算局部密度的时间复杂度为 $O(n \log_2 n + n)$,建立邻接矩阵的时间复杂度为 $O(kn)$,计算子图个数即对邻接矩阵构成的有向图进行搜索,时间复杂度为 $O(n + e)$,其中 e 为具有近邻

关系的边数,最大为 nk . 在得到聚类中心时,算法只一次遍历便可完成对象划分,时间复杂度为 $O(n)$. 算法的时间复杂度之和为 $O(kn \log_2 n + 2n \log_2 n + 3n + 2kn)$,则算法的时间复杂度为 $O(kn \log_2 n)$. 表5给出了CA-FM算法与对比算法的时间复杂度,其中Item表示算法迭代次数, t 表示CDP算法中被迪杰斯特拉算法处理的对象个数, K 表示数据集的真实聚类个数. 可以看出,CA-FM算法的时间复杂度优于DBSCAN、DPC、LPC和CDP算法.

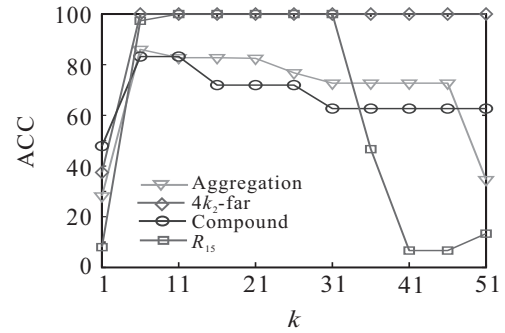
表5 算法的时间复杂度分析

算法	时间复杂度
K -means	$O(\text{Item} \times nK)^{[13]}$
AP	$O(n^2 \log_2 n)^{[13]}$
FCM	$O(\text{Item} \times nK)^{[9]}$
DBSCAN	$O(n^2)^{[10]}$
DP	$O(n^2)^{[11]}$
LPC	$O(n^2)^{[12]}$
CDP	$O(n^2 \log_2 n + tn^2 + n \log_2 K)^{[13]}$
CA-FM	$O(kn \log_2 n)$

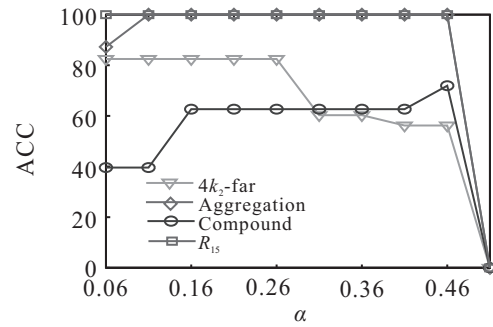
2.4.2 参数敏感性分析

CA-FM算法有两个参数:近邻对象数 k 、过滤因子 α . 本文选取4个二维数据集与4个高维数据集对两个参数进行敏感性分析. 图3(a)和图3(c)给出了参数 k 与聚类精度之间的关系,可以看出,随着 k 增大,对象的采样空间不断扩展,局部密度可有效表征数据的真实分布情况,其聚类精度越来越高. 当 k 继续增大时,会导致数据对象的近邻关系延伸至其他簇中,造成多个簇合并从而导致聚类精度下降. 一般情况下,当 k 的取值为5~31时,算法能达到理想的聚类效果.

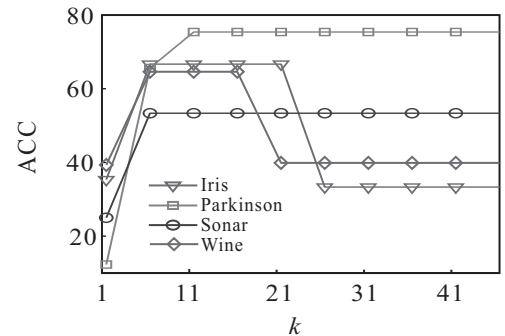
图3(b)和图3(d)给出了参数与聚类精度之间的关系. 参数 α 为过滤因子,其作用是过滤掉干扰聚类的对象. 对于二维数据集,当 α 取值偏小时,无法有效地去除全部的非核心对象,所以无法得到准确的聚类骨架,这时聚类精度偏低;对于高维数据集,其数据集的结构并不遵循传统的统计分布, α 取值较小时不会对聚类精度产生影响,但随着参数 α 增大,部分核心对象被过滤掉,从而得不到完整的聚类骨架,造成聚类精度下降. 特殊情况下,当 $\alpha = 1$ 时,表示全部数据对象均被视为非核心对象,聚类精度最差. 一般情况, α 的取值为0.1~0.41时,算法能达到理想的聚类效果.



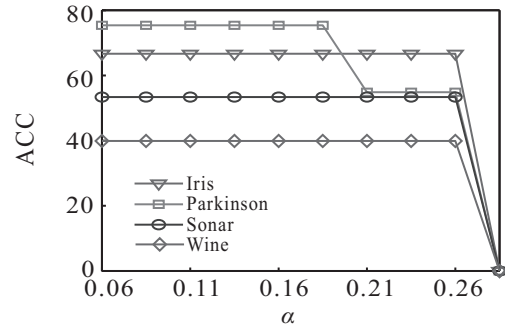
(a) 近邻参数 k 与精度的关系(二维数据集)



(b) 近过滤因子 α 与精度的关系(二维数据集)



(c) 近邻参数 k 与精度的关系(高维数据集)



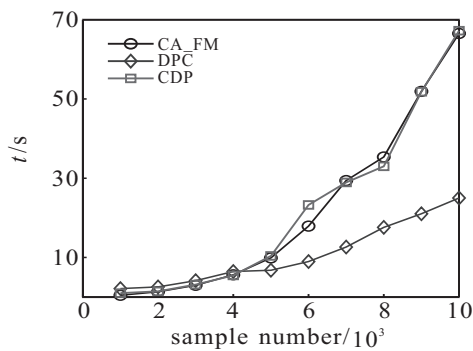
(d) 近过滤因子 α 与精度的关系(高维数据集)

图3 聚类精度与参数 k 、 α 的关系

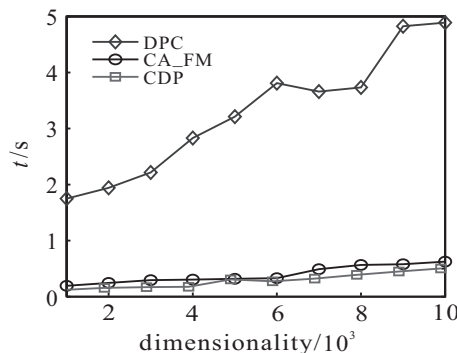
2.4.3 伸缩性分析

本文将合成数据集Flame^[6]扩展为大样本量数据集,将UCI数据集Soybean扩充为高维度数据集,用来检验算法的运行效率,并与DPC、CDP算法进行比较. 由于LPC算法需要计算每一个维度的特征值,其时间消耗特别大,数据集维度1000时,运行时间已经超过所显示范围,因此实验中没有给出与LPC算法的曲线. 图4(a)、图4(b)分别给出了样本量、维度与运

行时间的关系. 由图4可看出, CA-FM算法的运行时间介于DPC和CDP算法之间. 图4(a)中, 为了在固定范围内得到直观的趋势对比, 将DPC算法所需的复杂参数DATA提前进行预处理, 因此, DPC算法的运行时间最少. 由图4(b)可知, CA-FM算法的时间消耗与CDP算法大致相当, 而DPC算法由于需要人工选取决策图中的聚类原型, 其运行时间与选取过程的耗时有关, 因此曲线成不规则的上升趋势. 本算法的运行时间与样本量呈类正比函数关系, 且对维度变化的敏感性不大, 运行时间变化不大.



(a) 运行时间与样本量的关系



(b) 运行时间与维度的关系

图4 算法的运行时间与样本量、维度的关系

3 结论

本文在定义偏差因子的基础上建立了非核心对象过滤模型, 消除了噪声和边界点对原型提取的干扰; 以提出的密度因子和邻接矩阵, 解决了原型选取不合理的问题. CA-FM算法可以对含有噪声的多密度和高维数据集进行有效的聚类, 并具有较高的精度, 同时为研究聚类原型提供了一种框架.

参考文献(References)

[1] Kacprzyk J, Pedrycz W. Springer handbook of computational intelligence[M]. Berlin: Springer Publishing Company, 2015: 578-600.
 [2] Li X L, Han Q, Qiu B Z. A clustering algorithm using skewness-based boundary detection[J]. Neurocomputing, 2018, 275: 618-626.

[3] Mondal S A. An improved approximation algorithm for hierarchical clustering[J]. Pattern Recognition Letters, 2018, 104: 23-28.
 [4] Chen H Z, Wang W W, Feng X C, et al. Discriminative and coherent subspace clustering[J]. Neurocomputing, 2018, 284: 177-186.
 [5] 李海林, 王成, 邓晓懿. 基于分量属性近邻传播的多元时间序列数据聚类方法[J]. 控制与决策, 2018, 33(4): 649-656.
 (Li H L, Wang C, Deng X Y. Multivariate time series clustering based on affinity propagation of component attributes[J]. Control and Decision, 2018, 33(4): 649-656.)
 [6] 李向丽, 曹晓锋, 邱保志. 基于矩阵模型的高维聚类边界模式发现[J]. 自动化学报, 2017, 43(11): 1962-1972.
 (Li X L, Cao X F, Qiu B Z. Clustering boundary pattern discovery for high dimensional space base on matrix model[J]. Acta Automatica Sinica, 2017, 43(11): 1962-1972.)
 [7] 张秦, 方志耕, 蔡佳佳, 等. 基于多元异构不确定性案例学习的广义区间灰数熵权聚类模型[J]. 控制与决策, 2018, 33(8): 1481-1488.
 (Zhang Q, Fang Z G, Cai J J, et al. Generalized interval grey entropy-weight clustering model based on multiple heterogeneous uncertainty cases study[J]. Control and Decision, 2018, 33(8): 1481-1488.)
 [8] Macqueen J. Some methods for classification and analysis of multivariate observations[C]. Proc of 5th Berkeley Symposium on Mathematical Statistics and Probability. Berkeley: California Press, 1967: 281-297.
 [9] Bezdek J C, Robert E, Full W. FCM: The fuzzy c -means clustering algorithm[J]. Computers & Geosciences, 1984, 10(2/3): 191-203.
 [10] Ester M, Kriegel H P, Xu X W, et al. A density-based algorithm for discovering clusters in large spatial databases with noise[J]. Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD-96). Portland: Association for the Advancement of Artificial Intelligence, 1996: 226-231.
 [11] Rodriguez A, Laio A. Clustering by fast search and find of density peaks[J]. Science, 2014, 344(6191): 1492-1496.
 [12] Yang X H, Zhu Q P, Huang Y J, et al. Parameter-free laplacian centrality peaks clustering[J]. Pattern Recognition Letters, 2017, 100: 167-173.
 [13] Li Z J, Tang Y C. Comparative density peaks clustering[J]. Expert Systems with Applications, 2018, 95: 236-247.
 [14] Qi X, Fuller E, Wu Q, et al. Laplacian centrality: A new centrality measure for weighted networks[J]. Information Sciences, 2012, 194: 240-253.

- [15] Kumar K M, Reddy A R M. A fast dbscan clustering algorithm by accelerating neighbor searching using groups method[J]. *Pattern Recognition*, 2016, 58: 39-48.
- [16] Li X L, Han Q, Qiu B Z. A clustering algorithm with affine space-based boundary detection[J]. *Applied Intelligence*, 2018, 48(2): 432-444.
- [17] Qiu B Z, Cao X F. Clustering boundary detection for high dimensional space based on space inversion and Hopkins statistics[J]. *Knowledge-Based Systems*, 2016, 98: 216-225.
- [18] Chen H P, Shen X J, Lv Y D, et al. A novel automatic fuzzy clustering algorithm based on soft partition and membership information[J]. *Neurocomputing*, 2017, 236: 104-112.
- [19] Zhang Y, Madziuk J, Chai H Q, et al. Curvature-based method for determining the number of clusters[J]. *Information Sciences*, 2017, 415: 414-428.
- [20] 陈晋音, 何辉豪. 基于密度的聚类中心自动确定的混合属性数据聚类算法研究[J]. *自动化学报*, 2015, 41(10): 1798-1813.
(Chen J Y, He H H. Research on density-based clustering algorithm for mixed data with determine cluster centers auto matically[J]. *Acta Automatica Sinica*, 2015, 41(10): 1798-1813.)
- [21] Ding S F, Du M J, Sun T F, et al. An entropy-based density peaks clustering algorithm for mixed type data employing fuzzy neighborhood[J]. *Knowledge-Based Systems*, 2017, 133: 294-313.
- [22] Aliguliyev R M. Performance evaluation of density-based clustering methods[J]. *Information Sciences*, 2009, 179(20): 3583-3602.
- [23] 乔颖, 王士同, 杭文龙. 大规模数据集引力同步聚类[J]. *控制与决策*, 2017, 32(6): 1075-1083.
(Qiao Y, Wang S T, Hang W L. Clustering by gravitational synchronization on large scale dataset[J]. *Control and Decision*, 2017, 32(6): 1075-1083.)
- [24] 徐明亮, 王士同, 杭文龙. 一种基于同类约束的半监督近邻反射传播聚类方法[J]. *自动化学报*, 2016, 42(2): 255-269.
(Xu M L, Wang S T, Hang W L. A semi-supervised affinity propagation clustering method with homogeneity constraint[J]. *Acta Automatica Sinica*, 2016, 42(2): 255-269.)
- [25] Chen M, Li L J, Wang B, et al. Effectively clustering by finding density backbone based-on knn[J]. *Pattern Recognition*, 2016, 60: 486-498.
- [26] Campo D N, Stegmayer G, Milone D H. A new index for clustering validation with overlapped clusters[J]. *Expert Systems with Applications*, 2016, 64: 549-556.
- [27] Du M, Ding S, Xue Y. A novel density peaks clustering algorithm for mixed data[J]. *Pattern Recognition Letters*, 2017, 97: 46-53.

作者简介

邱保志(1962—), 男, 教授, 博士, 从事机器学习与数据挖掘等研究, E-mail: iebzqiu@zzu.edu.cn;

张瑞霖(1995—), 男, 硕士生, 从事机器学习与数据挖掘的研究, E-mail: zzurlz@163.com;

李向丽(1965—), 女, 教授, 博士, 从事计算机网络等研究, E-mail: iexlli@zzu.edu.cn.

(责任编辑: 孙艺红)