

# 一种基于密度和网格的簇心可确定聚类算法

何熊熊<sup>†</sup>, 管俊轶, 叶宣佐, 詹亦钊

(浙江工业大学 信息工程学院, 杭州 310023)

**摘要:** 以网格化数据集来减少聚类过程中的计算复杂度, 提出一种基于密度和网格的簇心可确定聚类算法. 首先网格化数据集空间, 以落在单位网格对象里的数据点数表示该网格对象的密度值, 以该网格到更高密度网格对象的最近距离作为该网格的距离值; 然后根据簇心网格对象同时拥有较高的密度和较大的距离值的特征, 确定簇心网格对象, 再通过一种基于密度的划分方式完成聚类; 最后, 在多个数据集上对所提出算法与一些现有聚类算法进行聚类准确性与执行时间的对比实验, 验证了所提出算法具有较高的聚类准确性和较快的执行速度.

**关键词:** 数据挖掘; 数据聚类; 网格; 密度

中图分类号: TP18

文献标志码: A

## A density-based and grid-based cluster centers determination clustering algorithm

HE Xiong-xiong<sup>†</sup>, GUAN Jun-yi, YE Xuan-zuo, ZHAN Yi-zhao

(College of Information Engineering, Zhejiang University of Technology, Hangzhou 310023, China)

**Abstract:** A density and grid based cluster centers determination clustering algorithm is proposed. The computational complexity of the clustering process is reduced by using the gridding dataset. Firstly, the dataset space is divided into grids with the same size, and the number of data objects that are contained in grid is defined as the value of the grid density. The nearest distance from one grid to another with higher density is defined as the value of grid distance. The cluster center grids can be found since these grids always have high density value and large distance value. Then, a density-based division approach is used to accomplish the task of clustering. Finally, a comprehensive comparison is carried out to examine the clustering accuracy and execution time between the proposed clustering algorithm and some classical algorithms. Experiment results show that the proposed algorithm can lead to a higher accuracy with less execution time.

**Keywords:** data mining; data clustering; grid; density

## 0 引言

数据挖掘 (DM)<sup>[1]</sup> 是指从大量数据中发现未知的、有价值的模式或规律等知识的复杂过程. 聚类分析<sup>[2]</sup> 作为数据挖掘技术中一个重要组成部分, 它的目的是将物理对象或抽象对象的集合中相似的对象聚在同一个簇中. 通俗地讲, 簇就是相似对象的集合, 同簇中的对象相似度较高, 不同簇之间的对象相似度较低. 聚类分析是一种常见的数据分析工具, 在统计学、机器学习、数据挖掘、生物学、市场营销等领域都有着广泛的应用前景<sup>[3-4]</sup>. 最早的聚类算法是由 MacQueen 在 1967 年提出的  $K$ -means<sup>[5]</sup> 聚类算法, 经过数十年国内外专家学者的深入研究, 在聚类算法

上已取得了相当丰硕的成果<sup>[6]</sup>. 如基于划分的聚类算法<sup>[7]</sup> 有  $K$ -means、 $K$ -medoids<sup>[8]</sup>、CLARANS<sup>[9]</sup> 等; 基于密度的聚类算法有 DBSCAN<sup>[10]</sup>、GDBSCAN<sup>[11]</sup>、OPTICS<sup>[12]</sup> 等; 基于层次的聚类算法有 CURE<sup>[13]</sup>、CHAMELEON<sup>[14]</sup>、BIR-CH<sup>[15]</sup> 等; 基于网格的聚类算法有 STING<sup>[16]</sup>、Wave-Cluster<sup>[17]</sup>、CLIQUE<sup>[18]</sup> 等. 然而传统的聚类算法都存在着各自的不足, 如  $K$ -means、 $K$ -medoids、Fuzzy  $K$ -means、Spectral Clustering<sup>[19-20]</sup> 聚类算法, 无法确定簇心个数;  $K$ -means、 $K$ -medoids、Fuzzy  $K$ -means<sup>[21]</sup>、BIRCH 聚类算法不能用来处理任意形状的数据集; 尽管 DBSCAN、Spectral Clustering 聚类算法适用于任意形状的数据集, 但它们的聚类质

收稿日期: 2016-02-29; 修回日期: 2016-07-05.

基金项目: 国家自然科学基金项目 (61473262).

作者简介: 何熊熊 (1965—), 男, 教授, 博士生导师, 从事重复学习控制、网络控制系统等研究; 管俊轶 (1991—), 男, 硕士生, 从事数据挖掘的研究.

<sup>†</sup>通讯作者. E-mail: hxx@zjut.edu.cn

量过于依赖参数的设定;而基于网格的聚类算法如 STING、CLIQUE 常被用作处理大规模的数据集,可减少计算复杂度,但会降低聚类质量。

Rodriguez 等<sup>[22]</sup>在 Science 期刊上提出了一种可以处理任何形状数据集的算法.该算法假设聚类中心具有较高的密度  $\rho$ ,且它与更高密度点具有较大的距离  $\delta$ .与传统聚类算法相比,该算法具有较高的聚类质量,但仍存在着下列问题:1)数据对象的  $\rho$ 、 $\delta$  两值依赖于截断距离参数  $dc$  的取值;2)算法需要计算所有数据对象之间的两两距离,计算量过于庞大。

针对上述问题,本文提出一种基于密度和网格的簇心可确定聚类算法(DGCCD).该算法首先网格化数据集空间,形成网格对象集;其次,分别定义网格对象的密度值  $\rho$  和距离值  $\delta$ ;再次,将同时具有较高密度值且较大距离值的网格对象选作聚类中心;然后根据一种基于密度的划分方法,完成网格对象集的聚类;最后对每个数据对象的类标进行划分,完成数据集的聚类。

### 1 DGCCD 算法的具体步骤

DGCCD 算法需要将数据空间的每一维都均匀划分成相同的段数,记为  $f_G$ ;形成若干个等大的网格对象,剔除网格对象中数据对象个数为 0 的网格对象,剩余网格对象为网格对象集,记为  $G$ ;网格对象集中网格对象数记为  $N_G$ .实验结果表明,网格对象集中网格对象数量  $N_G$  在大于或等于数据集中数据量  $n$  的 1/6 的情况下,算法具有较好的聚类质量. DGCCD 算法的具体步骤如下。

Step 1: 将含有  $n$  个数据对象的数据集空间中每一维均匀划分成相同段数,记为  $f$  (赋予  $f$  初始值 2),形成网格;

Step 2: 剔除网格对象中数据对象数量为 0 的网格,剩余非空网格对象数记为  $N$ ;

Step 3: 如果  $N < n/6$ ,则令  $f = f + 1$ ,返回 Step 1;

Step 4: 确定网格对象集  $G$ ,划分段数  $f_G$ ,网格对象数  $N_G$ .

本文将以上过程称为 DGCCD 算法的网格化过程。

针对样本数据集  $D_1$ ,其二维分布如图 1(a) 所示.将数据集  $D_1$  通过上述 DGCCD 算法网格化后的网格对象集分布如图 1(b) 所示。

网格对象集一旦确定,本文接下来需通过计算求得网格对象的密度值  $\rho$  和距离值  $\delta$ ,定义分别如下。

定义 1 以落在网格对象  $i$  中数据对象的数量作为网格对象密度值,记为  $\rho_i$ .

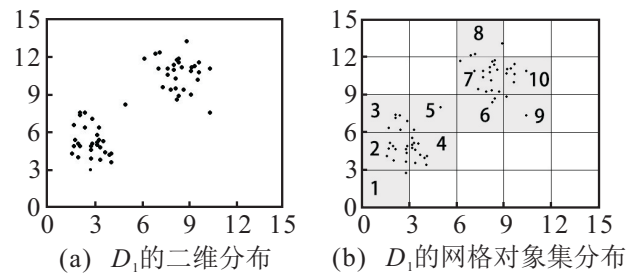


图 1  $D_1$  数据集的二维分布及其网格对象集的二维分布

定义 2 以网格对象  $i$  到更高密度网格对象  $j$  的最近距离作为网格对象的距离值,记为  $\delta_i$ ,定义如下:

$$\delta_i = \min_{j:\rho_j > \rho_i} (d_{ij}), \tag{1}$$

其中  $d_{ij}$  为网格对象  $i$  中心位置到网格对象  $j$  中心位置之间的欧氏距离。

令网格对象集中最高密度的网格对象为  $i_m$ ,其距离值  $\delta_{i_m} = \max(\delta_j)$ ,其中  $j$  为除  $i_m$  以外的所有网格对象。

在 Rodriguez 和 Laio 在 Science 期刊上提出的聚类中心具有较高的密度  $\rho$  且与更高密度点具有较大的距离  $\delta$  的相同假设下, DGCCD 算法网格化后的网格对象集中,处在簇心位置的网格对象会同时具有较高的密度  $\rho$  和较大距离值  $\delta$ .如图 2 所示,将编号为“4”、“7”的网格对象作为簇心网格对象,它们都同时具有较高的密度  $\rho$  和较大的距离  $\delta$ .在确定网格对象的密度值  $\rho$  和距离值  $\delta$  之后,作出网格对象集对应的  $\rho$  和  $\delta$  分布图.然后根据处在簇心位置的网格对象会同时具有较高的密度  $\rho$  和较大距离值  $\delta$  的特征,可通过人为监督的方式选取出簇心网格对象。

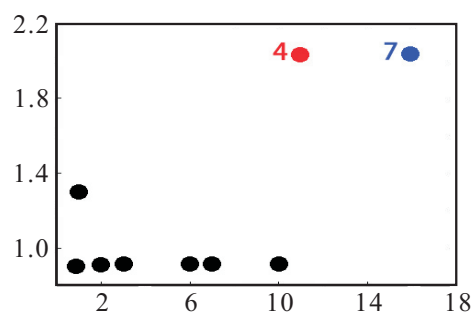


图 2  $D_1$  的网格对象集的  $\rho$  和  $\delta$  分布

在簇心网格对象确定之后,首先给每个簇心网格对象赋予不同的类标,然后采用一种基于密度的划分方式,即每个网格对象的类标跟随它的最近且密度高于它的网格对象类标,对剩余网格对象进行划分,完成网格对象集的聚类.网格对象集聚类完成后,每个数据对象类标只需跟随其落入的网格对象类标进行划分,从而完成整个数据集的聚类。

图 3 为网格对象划分的图示,网格中的数字表示网格对象的密度值.设网格“8”类标为 A,对于网格

“6”而言,网格“8”是与它最近且密度值大于它的网格对象,故网格“6”的类标则跟随网格“8”变为A.对于网格“2”和网格“5”,网格“6”是距离它们最近且密度值大于它们的网格对象,故类标跟随网格“6”变为A.

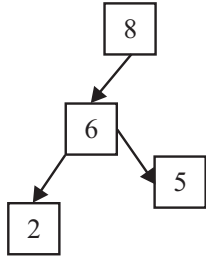


图3 网格对象的划分规则

算法完成划分后,首先选取边缘的网格对象,计算边缘网格对象中的每个数据点与所有非边缘数据点彼此间的距离,以距离该数据点最近的非边缘数据点的类标作为该边缘数据点的类标,从而完成边缘数据点的划分. DataSet1 经算法边缘处理前后的对比如图4所示.图4(a)由于网格化,在类与类之间的边缘划分较为粗糙,经过本文算法的边缘处理,边缘粗糙的问题基本得到解决,如图4(b)所示.

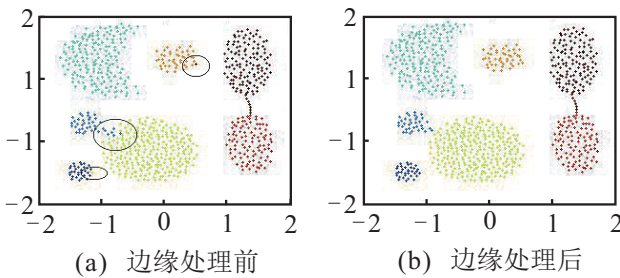


图4 边缘处理前后 DataSet1 的聚类结果

噪声点的处理不需要人为设定噪声点阈值截断,而是通过先寻找出类与类之间的边界网格对象,然后找出边界网格对象中密度最高的网格对象的密度作为阈值,记为  $\rho_b$ ,只需保留类中大于或等于  $\rho_b$  的网格对象,而低于  $\rho_b$  的网格对象即为噪声点.图5为样本数据集未去噪和去噪后的图示,其中黑点为噪声点.

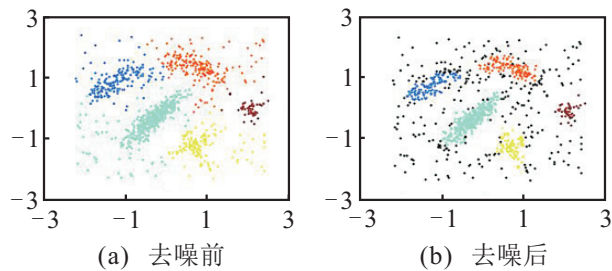


图5 噪声处理前后样本数据集的聚类结果

## 2 仿真实验与性能分析

实验所用操作系统为 Windows 7, 仿真软件 Matlab 7.12.0(R2011a)(64 位), 处理器为 Intel(R)Core (TM)i5, 安装内存 4.00 GB.

### 2.1 聚类性能对比

测试 DGCCD 算法聚类性能的实验数据集为 DataSet1、DataSet2、DataSet3、DataSet4、DataSet5、DataSet6 和 DataSet7(Iris 数据集), 数据集的具体信息如表 1 所示. 采用的聚类准确率  $r$  为正确分类数据量占总数据量的百分比. 实验中将每个算法在每一个数据集上进行 10 次聚类运算, 求得平均聚类准确度  $r_a$  和平均算法执行时间  $t_a$ , 提取各算法 10 次测试中最高聚类准确度  $r_{max}$ 、最低聚类准确度  $r_{min}$  和 10 次聚类准确度的方差  $\sigma^2$ , 用作分析聚类结果的稳定性.

表 1 7 个数据集的基本属性

聚类数据集	数据集属性		
	数据集维数	类属性数	数据量
DataSet 1	2	7	788
DataSet 2	2	31	3 100
DataSet 3	2	20	1 000
DataSet 4	2	15	600
DataSet 5	2	5	2 000
DataSet 6	2	2	1 829
DataSet 7	4	3	150

DataSet1 ~ DataSet6 数据集通过 DGCCD 算法聚类后的二维分布如图 6 所示. DataSet1~DataSet6 数据集网格化后网格对象集对应的  $\rho$  和  $\delta$  分布如图 7 所示, DataSet7 数据集网格化后网格对象对应的  $\rho$  和  $\delta$  分布如图 8 所示, 实验中选取密度值和距离值都较大的网格对象作为簇心网格对象. 7 个数据集网格化后网格对象集对应的  $N$  和  $f_G$ , 如表 2 所示.

表 2 7 个数据集网格化后的基本属性

聚类数据集	网格对象集属性	
	最终划分段数 $f_G$	网格对象数 $N$
DataSet 1	14	134
DataSet 2	31	540
DataSet 3	19	171
DataSet 4	22	106
DataSet 5	27	345
DataSet 6	23	324
DataSet 7	11	116

本文以  $K$ -means、Fuzzy  $K$ -means、DBSCA-N、Spectral Clustering (SC)、Rodriguez 等在 Science 提出的算法(下文用 Rodriguez-Clustering 表示)和 DGCCD

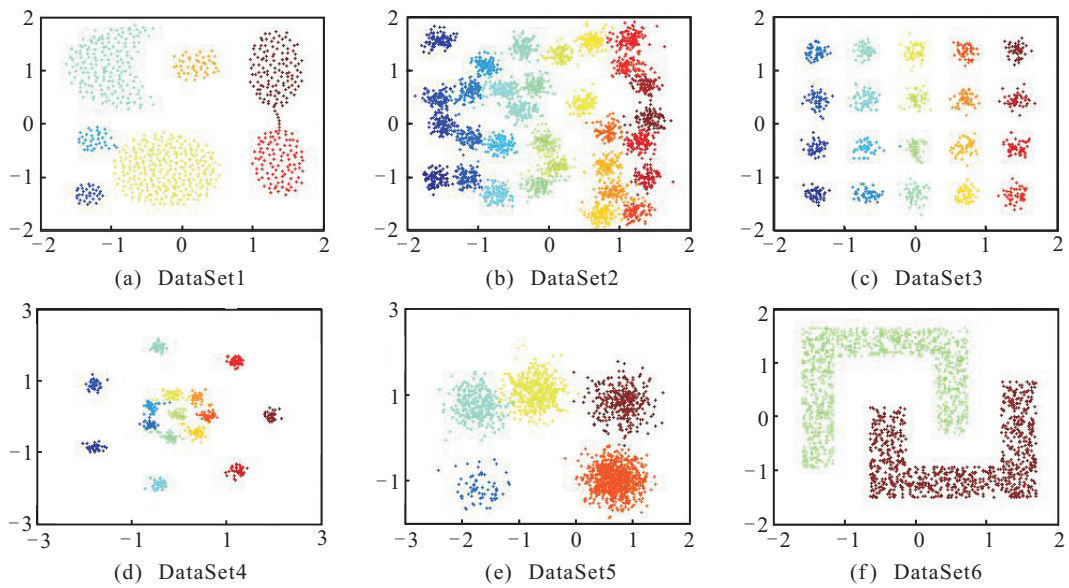


图6 6个数据集的二维分布以及DGCCD算法的聚类结果

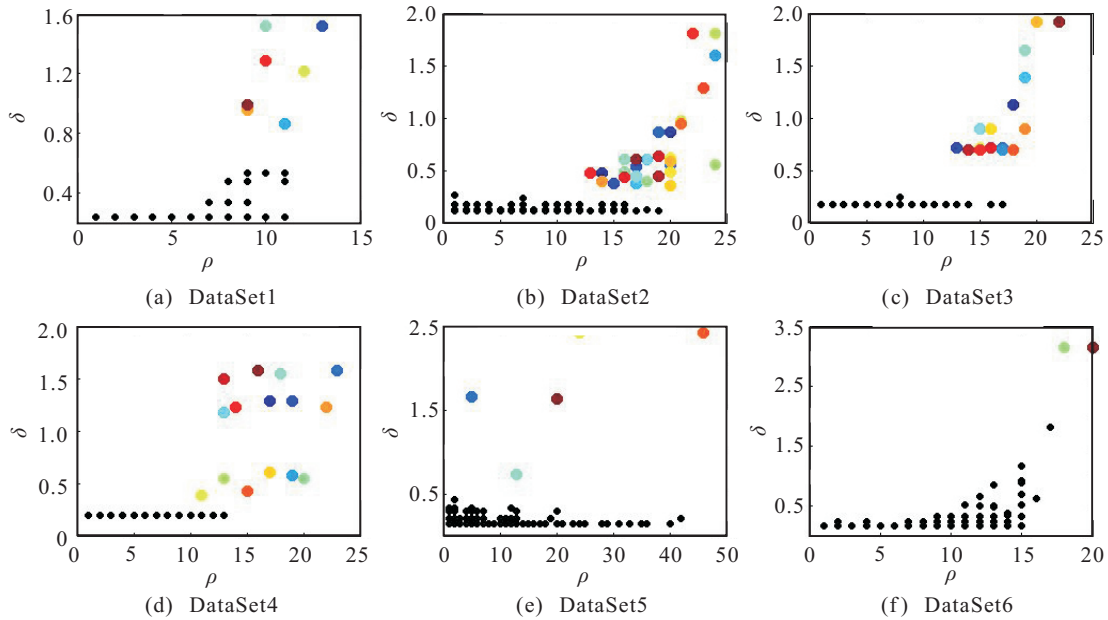


图7 6个数据集网络化后的网格对象对应的 $\rho$ 和 $\delta$ 分布

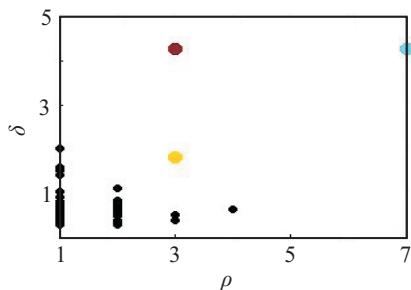


图8 DataSet的网格对象集的 $\rho$ 和 $\delta$ 分布

算法,在以上给出的7个数据集上的聚类性能进行对比实验,实验结果分别如表3~表9所示。

### 2.2 与Rodriguez-Clustering的算法执行时间对比

由表3~表9的实验结果可知,DGCCD算法与

Rodriguez-Clustering算法的聚类准确度较为接近,且相比其他算法有着较高的聚类准确性,并都具有绝对的聚类稳定性.表10为DGCCD算法与Rodriguez-Clustering算法在7个数据集上的平均执行时间的对比表.由此可见,DGCCD算法在7个数据集上的执行时间都小于Rodriguez-Clustering算法,其中DataSet2、DataSet5和DataSet6由于数据量相对较大,Rodriguez-Clustering算法的执行时间大幅度上升,而DGCCD算法的执行时间小于Rodriguez-Clustering算法.由此可得,DGCCD算法的执行速度快于Rodriguez-Clustering算法。

表3 6种算法在DataSet1数据集上的实验结果

聚类算法	DataSet1				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.7844	0.8680	0.7259	0.0033	0.0160
DBSCAN(Eps = 0.1348)	0.8249	0.8249	0.8249	0	0.0583
Fuzzy <i>K</i> -means	0.7236	0.7995	0.5888	0.0068	0.2145
SC( $\alpha = 1$ )	0.9175	0.9962	0.6802	0.0119	0.4000
Rodriguez-Clustering	0.9987	0.9987	0.9987	0	3.5079
DGCCD	0.9860	0.9860	0.9860	0	0.1694

表4 6种算法在DataSet2数据集上的实验结果

聚类算法	DataSet2				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.8097	0.8897	0.7181	0.0025	0.0866
DBSCAN(Eps = 0.0760)	0.4645	0.4645	0.4645	0	0.4614
Fuzzy <i>K</i> -means	0.8503	0.9087	0.7265	0.0029	1.4133
SC( $\alpha = 1$ )	0.7629	0.8194	0.6845	0.0020	9.7224
Rodriguez-Clustering	0.9674	0.9674	0.9674	0	179.5844
DGCCD	0.9413	0.9413	0.9413	0	2.2783

表5 6种算法在DataSet3数据集上的实验结果

聚类算法	DataSet3				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.8084	0.9130	0.6890	0.0060	0.0239
DBSCAN(Eps = 0.1191)	0.9900	0.9900	0.9900	0	0.1019
Fuzzy <i>K</i> -means	0.8904	0.9260	0.7870	0.0023	0.2274
SC( $\alpha = 1$ )	0.6875	0.8100	0.5280	0.0068	0.7167
Rodriguez-Clustering	1.0000	1.0000	1.0000	0	6.3399
DGCCD	0.9990	0.9990	0.9990	0	0.2373

表6 6种算法在DataSet4数据集上的实验结果

聚类算法	DataSet4				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.8153	0.9667	0.7000	0.0068	0.0155
DBSCAN(Eps = 0.0972)	0.9450	0.9450	0.9450	0	0.0530
Fuzzy <i>K</i> -means	0.9967	0.9967	0.9967	0	0.1404
SC( $\alpha = 1$ )	0.6713	0.8267	0.5217	0.0074	0.2720
Rodriguez-Clustering	0.9967	0.9967	0.9967	0	1.3270
DGCCD	0.9667	0.9667	0.9667	0	0.1760

表7 6种算法在DataSet5数据集上的实验结果

聚类算法	DataSet5				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.7528	0.9950	0.6305	0.0228	0.0210
DBSCAN(Eps = 0.1002)	0.7970	0.7970	0.7970	0	0.2279
Fuzzy <i>K</i> -means	0.7560	0.7560	0.7560	0	0.2787
SC( $\alpha = 1$ )	0.8902	0.9905	0.6905	0.0108	3.1721
Rodriguez-Clustering	0.9930	0.9930	0.9930	0	79.4415
DGCCD	0.9920	0.9920	0.9920	0	1.2259

表8 6种算法在DataSet6数据集上的实验结果

聚类算法	DataSet6				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.888 5	0.888 5	0.888 5	0	0.012 6
DBSCAN(Eps = 0.119 7)	1.000 0	1.000 0	1.000 0	0	0.184 3
Fuzzy <i>K</i> -means	0.890 7	0.890 7	0.890 7	0	0.089 7
SC( $\alpha = 1$ )	1.000 0	1.000 0	1.000 0	0	1.514 7
Rodriguez-Clustering	1.000 0	1.000 0	1.000 0	0	38.305 1
DGCCD	1.000 0	1.000 0	1.000 0	0	0.209 2

表9 6种算法在DataSet7数据集上的实验结果

聚类算法	DataSet7				
	平均聚类准确度 $r_a$	最高准确度 $r_{max}$	最低准确度 $r_{min}$	聚类准确度方差 $\sigma^2$	平均执行时间 $t_a/s$
<i>K</i> -means	0.781 1	0.825 1	0.564 3	0.002 4	0.018 9
DBSCAN(Eps = 1.081 4)	0.770 0	0.770 0	0.770 0	0	0.073 8
Fuzzy <i>K</i> -means	0.811 0	0.900 0	0.690 0	0.005 8	0.048 4
SC( $\alpha = 1$ )	0.843 3	0.990 0	0.550 0	0.021 4	0.430 4
Rodriguez-Clustering	0.910 0	0.910 0	0.910 0	0	0.115 7
DGCCD	0.853 0	0.852 0	0.853 0	0	0.097 9

表10 Rodriguez-Clustering算法与DGCCD算法在7个数据集上的执行时间对比表果

数据集	平均执行时间/s		
	数据量	Rodriguez-Clustering	DGCCD
DataSet1	788	3.507 9	0.169 4
DataSet2	3100	179.584 4	2.278 3
DataSet3	1000	6.339 9	0.237 3
DataSet4	600	1.327 0	0.176 0
DataSet5	2000	79.441 5	1.225 9
DataSet6	1829	38.305 1	0.209 2
DataSet7	150	0.115 71	0.097 9

2.3 算法复杂度分析

假设聚类数据集有  $n$  个  $m$  维数据集, 首先, 算法网格化并获取网格对象的密度, 该过程的计算代价为  $O(mnf_G^2)$ ; 然后计算每个网格对象的距离需要的计算代价为  $O((N^2 - N)/2)$ , 算法进行一次划分完成聚类, 其计算代价为  $O(N\log(N) + N/2)$ ; 最后边缘处理的过程的计算代价为  $O(bn)$ , 其中  $b$  表示边缘点的个数. 因此 DGCCD 算法的算法时间复杂度为  $O((mnf_G^2) + (N^2 - N)/2 + N\log(N) + N/2 + bn)$ .

表11中列出了本文算法和对比算法的算法复杂度. 由表11可以分析得到: 相比 Rodriguez-Clustering 和 SC, 本文算法的时间复杂度较低, 所以相比前两个算法, 本文算法有着较快的执行速度; 而相比 DBSCAN, *K*-means, Fuzzy *K*-means 算法, 本文算法

的时间复杂度较高, 其中主要消耗在获取每个网格对象距离值的过程中. 但本文算法的优点在于能确定簇心位置, 以及在确定簇心网格后基于密度的划分方式, 使得本文算法对于任意形状分布的数据集都有着较好的聚类结果, 因此有着较高的聚类准确性.

表11 算法时间复杂度对比表

聚类算法	时间复杂度
<i>K</i> -means	$O(tkn)$
DBSCAN	$O(n\log(n))$
Fuzzy <i>K</i> -means	$O(ndk^2t)$
SC	$O(n^3)$
Rodriguez-Clustering	$O(n^2 + (n^2 - n)/2 + n\log(n) + n/2)$
DGCCD	$O((mnf_G^2) + (N^2 - N)/2 + N\log(N) + N/2 + bn)$

2.4 实验结果分析

综上7个数据集的实验结果, 可得 Rodriguez-Clustering 算法聚类准确度最高, 但在处理稍大的数据集时(数据量在2000以上), 算法执行时间过长; *K*-means、Fuzzy *K*-means、DBSCAN、SC算法的聚类准确性都不如 DGCCD 算法.

实验结果表明: 由以上6种算法针对7个数据集的测试结果可知, DGCCD 算法与 *K*-means、Fuzzy *K*-means、DBSCAN 和 SC 算法相比, 具有着较高的聚类准确度; 与 *K*-means、Fuzzy *K*-means、SC 算法相比, 具有着较高的聚类稳定性; 与 Rodriguez-Clustering 算法相

比,具有较快的执行速度.

### 3 结 论

本文提出了一种基于密度和网格的簇心可确定聚类算法. 该算法对数据集进行网格化,减少了聚类过程中的计算量,加快了算法执行速度. 在算法网格化的过程中,通过对众多数据集的实验测试学习,设定了最终网格对象总数  $N_G \geq n/6$ . 在实验过程中,笔者发现对于处理一些数据类属性数较少且数据集形状复杂度不高的大规模数据集时,网格对象数量  $N_G$  在小于  $n/6$  的情况下也不会影响聚类准确性,能进一步减少算法的执行时间. 即在不影响聚类准确性的前提下,  $f_G$  值越小,算法执行速度越快. 因此,如何给定一个最为合适的  $f_G$  来优化数据集的网格化,提升算法的执行速度,将是笔者下一步的研究重点之一.

#### 参考文献(References)

- [1] Culler D, Estrin D, Srivastava M. Guest editors' introduction: Overview of sensor networks[J]. Computer, 2004, 37(8): 41-49.
- [2] Wood A D, Stankovic J A. Denial of service in sensor networks[J]. Computer, 2002, 35(10): 54-62.
- [3] Yang H. Data mining: Concepts and techniques[J]. San Francisco, 2001, 29(S1): 1-18.
- [4] 王骏, 王士同, 邓赵红. 聚类分析研究中的若干问题[J]. 控制与决策, 2012, 27(3): 321-0328.  
(Wang J, Wang S T, Deng Z H. Survey on challenges in clustering analysis research[J]. Control and Decision, 2012, 27(3): 321-328.)
- [5] Lloyd S. Least squares quantization in PCM[J]. IEEE Trans on Information, 1982, 28(2): 129-137.
- [6] 覃艳, 王洪, 周全华. 数据挖掘中聚类算法的研究[J]. 网络安全技术与应用, 2014(1): 65-66.  
(Qin Y, Wang H, Zhou Q H. The research of clustering algorithm in data mining[J]. Network Security Technology & Application, 2014(1): 65-66.)
- [7] 吴杨, 王韬, 李进东. 基于密度的划分式聚类过程参数选择算法[J]. 控制与决策, 2016, 31(1): 21-29.  
(Wu Y, Wang T, Li J D. Clustering parameters selection algorithm based on density for divisional clustering process[J]. Control and Decision, 2016, 31(1): 21-29.)
- [8] 夏宁霞, 苏一丹, 覃希. 一种高效的  $K$ -medoids 聚类算法[J]. 计算机应用研究, 2010, 27(12): 4517-4519.  
(Xia N X, Su Y D, Qin X. Efficient  $K$ -medoids clustering algorithm[J]. Application Research of Computers, 2010, 27(12): 4517-4519.)
- [9] Ng R T, Han J. CLARANS: A method for clustering objects for spatial data mining[J]. IEEE Trans on Knowledge & Data Engineering, 2002, 14(5): 1003-1016.
- [10] Ester B M, Kriegel H P, Sander J, et al. A density based algorithm for discovering clusters in large spatial databases[C]. Proc of Int Conf on Knowledge Discovery and Data Mining. Protland, 1996.
- [11] Sander J, Ester M, Kriegel H P, et al. Density-based clustering in spatial databases: The algorithm GDBSCAN and its applications[J]. Data Mining & Knowledge Discovery, 1998, 2(2): 169-194.
- [12] Ankerst M, Breunig M M, Kriegel H P, et al. OPTICS: Ordering points to identify the clustering structure[J]. Stanford Research Inst Memo Stanford University, 1999, 28(2): 49-60.
- [13] Guha S, Rastogi R, Shim K. CURE: An efficient clustering algorithm for large databases[C]. Proc of the ACM SIGMOD Int Conf on Management of Data. New York: ACM Press, 1998: 73-84.
- [14] Karypis G, Han E H, Kumar V. CHAMELEON: A hierarchical clustering algorithm using dynamic Modeling[J]. IEEE Computer, 1999, 32(8): 68-75.
- [15] Zhang T, Ramakrishnan R, Livny M. BIRCH: An efficient data clustering method for very large databases[J]. Acm Sigmod Record, 1996, 25(2): 103-114.
- [16] Wang W, Yang J, Muntz R R. STING: A statistical information grid approach to spatial data mining[C]. Proc of the 23rd Int Conf on Very Large Data Bases. Morgan Kaufmann Publishers Inc, 1997: 186-195.
- [17] Sheikholeslami G, Chatterjee S, Zhang A. WaveCluster: A multi-resolution clustering approach for very large spatial databases[C]. Int Conf on Very Large Data Bases. New York: Morgan Kaufmann Publishers Inc, 1998: 428-439.
- [18] Agrawal R, Gehrke J E, Gunopulos D, et al. Automatic subspace clustering of high dimensional data for data mining applications[C]. Acm Sigmod Int Conf on Management of Data. Seattle: ACM Press, 1998: 95-105.
- [19] Ng A Y, Jordan M I, Weiss Y. On spectral clustering: Analysis and an algorithm[J]. Proc of Advances in Neural Information Processing Systems, 2001, 14: 849-856.
- [20] 周林, 平西建, 徐森, 等. 基于谱聚类的聚类集成算法[J]. 自动化学报, 2012, 38(8): 1335-1342.  
(Zhou L, Ping X J, Xu S, et al. Cluster ensemble based on spectral clustering[J]. Acta Automatica Sinica, 2012, 38(8): 1335-1342.)
- [21] Gasch A P, Eisen M B. Exploring the conditional coregulation of yeast gene expression through fuzzy  $K$ -means clustering[J]. Genome Biology, 2002, 3(11): 129-137.
- [22] Rodriguez A, Laio A. Clustering by fast search and find of density peaks[J]. Science, 2014, 344(6191): 1492-1496.

(责任编辑: 孙艺红)