

# 基于密度峰值聚类的高斯混合模型算法

王卫东 徐金慧 张志峰 杨习贝

江苏科技大学计算机学院 江苏 镇江 212100

**摘要** 由于存在大量服从高斯分布的样本数据,采用高斯混合模型(Gaussian Mixture Models,GMM)对这些样本数据进行聚类分析,可以得到比较准确的聚类结果。通常采用EM算法(Expectation Maximization Algorithm)对GMM的参数进行迭代式估计。但传统EM算法存在两点不足:对初始聚类中心的取值比较敏感;迭代式参数估计的迭代终止条件是相邻两次估计参数的距离小于给定的阈值,这不能保证算法收敛于参数的最优值。为了弥补上述不足,提出采用密度峰值聚类(Density Peaks Clustering,DPC)来初始化EM算法,以提高算法的鲁棒性,采用相对熵作为EM算法的迭代终止条件,实现对GMM算法参数值的优化选取。在人工数据集及UCI数据集上的对比实验表明,所提算法不但提高了EM算法的鲁棒性,而且其聚类结果优于传统算法。尤其在服从高斯分布的数据集上的实验结果显示,所提算法大幅提高了聚类精度。

**关键词:** 密度峰值聚类;相对熵;高斯混合模型;EM算法;聚类算法

**中图分类号** TP391.4

## Gaussian Mixture Models Algorithm Based on Density Peaks Clustering

WANG Wei-dong, XU Jin-hui, ZHANG Zhi-feng and YANG Xi-bei

College of Computer Science, Jiangsu University of Science and Technology, Zhenjiang, Jiangsu 212100, China

**Abstract** Due to the existence of a large number of sample data which obey the Gaussian distribution, GMM (Gaussian mixture models) is used to cluster these sample data and get more accurate clustering results. In general, EM algorithm(expectation maximization algorithm) is used to estimate the parameters of GMM iteratively. However, the traditional EM algorithm has two shortcomings: it is sensitive to the initial clustering center; the iterative termination condition of iterative parameter estimation is to judge that the distance between two adjacent estimated parameters is less than a given threshold, which can't guarantee that the algorithm converges to the optimal value of the parameters. In order to overcome the above shortcomings, density peaks clustering (DPC) is proposed to initialize EM algorithm to improve the robustness of the algorithm. The relative entropy is used as the iteration termination condition of the EM algorithm to optimize the parameters of GMM algorithm. The comparative experiments on artificial datasets and UCI datasets show that the new algorithm not only improves the robustness of EM algorithm, but also outperforms the traditional clustering algorithm. On the datasets which obey Gaussian distribution, the new algorithm greatly improves the clustering accuracy.

**Keywords** Density peaks clustering, Relative entropy, Gaussian mixture models, Expectation maximization algorithm, Clustering algorithm

## 1 引言

聚类分析是进行数据挖掘<sup>[1-2]</sup>、实现商业智能的重要工具,被广泛应用于医疗诊断、图像处理<sup>[3]</sup>、信息检索<sup>[4]</sup>和生物信息学等领域。聚类分析是将一组未知分布的样本数据分类到不同的类或者簇中的一个过程。对于聚类结果的评价,要求同一个簇中的对象有很大的相似性,而不同簇间的对象有很大的相异性。聚类分析算法主要包括以下4类<sup>[1-9]</sup>:基于划分的聚类方法、基于层次的聚类方法、基于网格的聚类方法和基于集成技术的聚类方法。

由于来自工业企业及服务行业的数据急速增长,如何从

这些大数据中抽取出有价值的信息成为了目前的研究热点,而聚类分析是解决上述问题的有力工具。文献[1]分析了如何运用不同的聚类算法来满足稀疏工业数据集的各种聚类需求。文献[2]针对不同的聚类算法在处理不确定数据、多媒体数据、图形数据、生物数据、流数据、文本数据、时间序列数据等数据类型时的特点,系统地讨论了数据聚类算法的应用。

在各种聚类方法中,基于划分的聚类方法因其简单易用等特点而得到了广泛的应用,如著名的  $k$ -means<sup>[5]</sup>及其改进算法。文献[6]运用  $k$ -means 聚类算法对教学满意度的调查数据进行挖掘和分析,分析过程快速、高效、合理,从而为采取有效措施改进教学策略和教学设计提供了可靠、有价值的参

收稿日期:2020-08-28 返修日期:2020-11-30

基金项目:国家自然科学基金(61572242)

This work was supported by the National Natural Science Foundation of China(61572242).

通信作者:王卫东(78653221@qq.com)

考依据。为了解决  $k$ -means 算法随机选取初始类别中心的问题,文献[7]将谱聚类算法与  $k$ -means 算法结合,优化了  $k$ -means 的初始类别中心点坐标。文献[8]提出了一种动态  $k$ -means 聚类算法,该算法在聚类的过程中可以动态地修改  $k$  值,实现了对聚类类别数的优化。

但是  $k$ -means 算法的缺点也很明显,如对聚类中心均值的简单使用,且假设数据点是呈圆形分布的。而有限混合模型是分析复杂现象的一种灵活而强有力的建模工具,它提供了用简单结构模拟复杂密度的有效方法。其中,基于高斯分布的高斯混合模型 GMM<sup>[10-11]</sup> 用于解决同一集合中的数据包含多个不同的高斯分布的情况。GMM 采用有限个高斯分布的概率密度函数,按照一定的权重进行混合,来量化数据的分布,从而达到聚类的目的。

对于高斯混合模型 GMM 的参数,通常采用 EM<sup>[10-11]</sup> (期望最大化)算法进行估计。EM 算法的优点是简单稳定,可以在不知道待估计参数先验信息和观测数据不完备的情况下,通过迭代来计算参数的最大似然估计。但 EM 算法对输入的初始值比较敏感,采用传统的随机初始值方法时,聚类结果往往波动较大,即算法的鲁棒性很差。而聚类结果的大幅波动极大地限制了 GMM 算法的应用。EM 算法的另一个问题是,在采用迭代式参数估计时,其迭代终止条件并不能保证所估计的参数值是最优值。

对于 EM 算法的初始值敏感问题,主要有两类解决方案:1)对随机初始值进行优化;2)对初始值进行优化选取。如采用  $k$ -means 算法对随机选取的初始值进行迭代式优化,但这种方法的结果同样不稳定,也会有较大的偏差。在优化选取初始值方法中,有代表性的是文献[12]中引入的 binning 方法,即装箱法。该方法将数据空间在各维上划分成一个一个的箱子,再把数据点投射到对应的箱子中,属于密度估计方法。使用该方法的难点在于每一维上最优或者近似最优的 bin 宽度的选取。目前,尚没有针对迭代终止条件的优化算法。

2014年6月,Science上发表了一篇题为“采用密度峰值点快速搜索的聚类算法”的文章,其提出的算法简称为 DPC<sup>[9]</sup> 算法。DPC 定义了两个评价指标:局部密度  $\rho_i$  和距离  $\delta_i$ 。通过这两个指标,该算法可以较准确地刻画出聚类中心的特征。因此,对于长期困扰聚类分析的两个基本问题(如何预估类簇数和初始类簇中心),DPC 提出了行之有效的解决方案。但由于该算法将样本判为离其最近且密度比它大的样本所在的类,当类间样本相互重叠时,聚类的准确率就会大大降低,故该算法的聚类准确率并不高。

本文提出采用密度峰值聚类算法 DPC 来初始化 EM 算法,解决了 EM 算法的初始值敏感问题。另一方面,本文提出采用相对熵(也称为 KL 散度)作为 EM 算法的迭代终止条件,实现了对 GMM 算法参数估计值的优化选取。

## 2 高斯混合模型

高斯混合模型将聚类的数据看作是来自于  $K$  个高斯分布的混合概率分布,这些高斯分布分别代表不同的类。由中心极限定理可知,在样本数量足够多的情况下,每个  $K$  所代表的区域就可以用高斯分布来描述。高斯函数的另一个优点

是具有良好的计算性能,高斯混合模型如下:

$$P(y|\theta) = \sum_{k=1}^K \alpha_k f(y|\theta_k) \quad (1)$$

其中,  $\alpha_k$  是系数,  $\alpha_k \geq 0$ , 且  $\sum_{k=1}^K \alpha_k = 1$ ,  $f(y|\theta_k)$  是高斯密度函数, 参数  $\theta_k = (u, \sigma_k^2)$ 。

通常采用 EM 算法对 GMM 中的参数  $\alpha_k, \theta_k$  进行估计。EM 的输入数据为观测数据  $y_1, y_2, \dots, y_N$ , 算法的输出是高斯混合模型的参数值。EM 算法分为两步: E 步是求解目标函数期望, 通常是求解目标函数取对数之后的期望值; M 步是期望最大化, 采用极大似然估计、拉格朗日乘法对参数求偏导, 最终确定新的参数。EM 算法的步骤如下:

(1) 随机选取参数  $\alpha_k, \theta_k$  的初值, 开始迭代;

(2) E 步: 依据当前模型的参数, 计算分模型  $k$  对观测数据  $y_j$  的响应度;

$$\hat{r}_{jk} = \frac{\alpha_k f(y_j|\theta_k)}{\sum_{k=1}^K \alpha_k f(y_j|\theta_k)} \quad (2)$$

其中,  $j=1, 2, \dots, N; k=1, \dots, K$ 。

(3) M 步: 计算新一轮迭代的模型参数;

$$\hat{\mu}_k = \frac{\sum_{j=1}^N \hat{r}_{jk} y_j}{\sum_{j=1}^N \hat{r}_{jk}}, k=1, 2, \dots, K \quad (3)$$

$$\hat{\sigma}_k^2 = \frac{\sum_{j=1}^N \hat{r}_{jk} (y_j - \hat{\mu}_k)^2}{\sum_{j=1}^N \hat{r}_{jk}}, k=1, 2, \dots, K \quad (4)$$

$$\hat{\alpha}_k = \frac{\sum_{j=1}^N \hat{r}_{jk}}{N}, k=1, 2, \dots, K \quad (5)$$

(4) 当满足下述公式时, 迭代结束。

$$\|\hat{\alpha}_k - \hat{\alpha}_{k-1}\| \leq \epsilon \quad (6)$$

其中,  $\epsilon$  为给定的阈值, 当不满足式(6)时, 算法转步骤(2), 继续迭代。

EM 算法实现简单, 数值计算稳定, 因此得到了广泛的应用。但 EM 算法也存在明显的不足: 算法采用随机的方式选取初始值, 导致初始值敏感, 从而使 GMM 算法的聚类结果大幅波动; 迭代的终止条件是相邻两次迭代参数  $\alpha_k$  的距离小于给定的阈值, 但该条件并不能保证估计的参数是最优的。

## 3 DPC 算法

DPC 算法可以有效地预估数据集的类簇数和初始聚类中心。但使用该算法进行聚类时, 将样本判为离其最近且密度比它大的样本所在的类, 易发生所谓的“多米诺骨牌效应”<sup>[13]</sup>。一旦某一个样本分类错误, 就会带来一连串的分类错误, 因此 DPC 算法的聚类准确率并不高。本文仅采用 DPC 算法预估数据集的类簇数和初始聚类中心, 目的是解决 EM 算法的初始值敏感问题。

为了刻画聚类中心, DPC 定义了两个评价指标: 局部密度  $\rho_i$  和距离  $\delta_i$ 。其基本思想是: 聚类中心是其局部密度大于围绕它的邻居的局部密度; 同时, 不同类中心之间的距离相对较远。因此, 聚类中心就是  $\rho_i$  和  $\delta_i$  同时较大的数据点。局部密度  $\rho_i$  是通过截断函数 Cut-off kernel 来定义的, 如式(7)所示:

$$\rho_i = \sum_{j \in I_i} X(d_{ij} - d_i) \quad (7)$$

其中,  $d_{ij}$  表示样本  $i$  和  $j$  的距离,  $d_c$  为截断距离, 函数  $X(x)$  的定义如下:

$$X(x) = \begin{cases} 1, & x < 0 \\ 0, & x \geq 0 \end{cases} \quad (8)$$

从式(7)、式(8)可以看出, 局部密度  $\rho_i$  描述的是围绕样本点的样本密度, 相当于落入以  $d_c$  为半径的超球体内的样本数量。显然,  $d_c$  取不同的值,  $\rho_i$  的值也将发生相应的改变, 因此  $\rho_i$  的值依赖于  $d_c$  的取值。

另一个参数是距离  $\delta_i$ , 它描述的是数据点之间的距离, 定义如下:

$$\delta_i = \begin{cases} \min_{j \in I_i} \{d_{ij}\}, & I_i \neq \emptyset \\ \max_{j \in I_i} \{d_{ij}\}, & I_i = \emptyset \end{cases} \quad (9)$$

其中, 指标集为:

$$I_i = \{k \in I_s; \rho_k > \rho_i\} \quad (10)$$

从式(9)、式(10)可以看出, 当  $x_i$  具有最大局部密度时,  $\delta_i$  表示  $S$  中与  $x_i$  距离最大的数据点与  $x_i$  之间的距离; 否则,  $\delta_i$  表示在所有局部密度大于  $x_i$  的数据点中, 与  $x_i$  距离最小的那个数据点与  $x_i$  之间的距离。

DPC 算法将每个数据点的  $\rho_i$  值和  $\delta_i$  值表示在一个二维决策图上。用户根据决策图的分布情况, 对聚类中心点进行选择, 这是一个人工操作, 无法自动完成。

针对 DPC 算法的不足, 文献[14]提出了一种根据不同的数据集自动计算出  $d_c$  最优值的方法, 简称 DF-DPC 算法。但是, 该算法并没有很好地解决 DPC 算法聚类准确率低的问题。文献[15]提出了一种扩展的 E-CFSFDP 算法, 该算法先调用 DPC 算法, 然后执行一个子类的合并步骤, 其目的是解决类簇中存在多密度峰值(或称为无密度峰值)的情况, 但该算法的聚类准确率也没有明显提高。文献[16]针对 DPC 算法效率较低的问题, 提出了一种基于相对邻域和剪枝策略的密度峰值快速搜索聚类算法 RP-DPC, 该算法属于时间性能优化算法, 其聚类精度与 DPC 算法相同。

## 4 基于相对熵的迭代优化

EM 算法的迭代终止条件是: 当相邻两次迭代的参数值  $\alpha_k$  的距离小于给定的阈值时, 迭代终止。但是, 一方面阈值的取值并没有明确的标准; 另一方面也不能保证参数  $\alpha_k, \theta_k$  取得最优值。本文研究了高斯混合模型的理论基础, 基于最小错误的贝叶斯决策, 如图 1 所示。

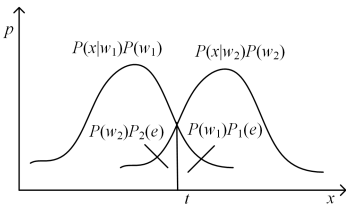


图 1 贝叶斯决策

Fig. 1 Bayesian decision making

图 1 中,  $P(\omega_1)P_1(e)$  为样本属于类别 1, 但被划分为类别 2 的情况。同样,  $P(\omega_2)P_2(e)$  为样本属于类别 2, 但被划分为类别 1 的情况。设总体的平均错误率为  $P(e)$ , 则  $P(e)$  为:

$$P(e) = P(\omega_2)P_2(e) + P(\omega_1)P_1(e)$$

$$P_1(e) = \int_t^{\infty} p(x|\omega_1)P(\omega_1)dx$$

$$P_2(e) = \int_{-\infty}^t p(x|\omega_2)P(\omega_2)dx$$

当采用 EM 算法估计的参数  $\alpha_k, \theta_k$  取得最优值时, 贝叶斯决策的平均错误率  $P(e)$  达到最小值。因此, EM 算法的迭代终止条件应该为: 当平均错误率  $P(e)$  最小时, 迭代结束。

如何判断平均错误率  $P(e)$  最小? 对于某个样本, 将其判别为两类的概率相差较小时, 其被错分的概率较大, 反之被错分的概率较小。本文引入信息论中的相对熵(relative entropy)<sup>[17]</sup>概念, 相对熵是对两个概率分布间差异的非对称性度量, 等价于两个概率分布的信息熵<sup>[17]</sup>(shannon entropy)的差值。以两类的分类问题为例, 当某个样本被判别为两类的概率差异越大, 则相对熵越大, 该样本被错分的概率就越小; 否则, 相对熵越小, 被判别为两类的概率差异越小, 该样本被错分的概率就越大。

当第  $i$  个样本的两类分类概率为  $P(i)$  和  $Q(i)$  时, 其相对熵的计算式如下:

$$D_i(p \| Q) = \sum_i P(i) \log_a \frac{P(i)}{Q(i)} \quad (11)$$

式(11)表示求  $P$  与  $Q$  之间的对数差在  $P$  上的期望值。当  $P(i) = Q(i)$  时, 相对熵为 0, 其他情况下相对熵大于 0。以相对熵为评价样本分类的度量指标, 其值越大, 聚类结果的可信度就越高, 否则可信度越低。

本文判断相对熵小于某个阈值的样本数, 当该样本数在某次迭代时取得极小值时认为 EM 算法在本次迭代中所估计的参数  $\alpha_k, \theta_k$  达到最优值。对于多类问题, 由于概率较大的两类发生错分的可能性较大, 因此取两个分类概率最大的概率值来计算其相对熵。

对于样本  $y_i$ , 概率最大的两个概率值分别是  $P(y_i)$  和  $Q(y_i)$ 。利用式(11)计算相对熵  $D_i(p \| Q)$ 。设相对熵的阈值为  $\delta$ , 采用下式判断相对熵小于阈值  $\delta$  的样本数。

$$f_j(y_i) = \begin{cases} 1, & D_i(p \| Q) < \delta \\ 0, & D_i(p \| Q) \geq \delta \end{cases} \quad (12)$$

其中,  $i = 1, \dots, N$  为样本数,  $j = 1, \dots, T$  为 EM 算法的迭代次数。采用式(13)计算某次迭代相对熵小于阈值的样本总数。

$$Sum_j = \sum_{i=1}^N f_j(y_i) \quad (13)$$

采用式(14)作为 EM 算法的迭代终止条件。

$$Sum_{t-1} \geq Sum_t \leq Sum_{t+1} \quad (14)$$

式(14)说明, 当某次迭代的相对熵小于阈值  $\delta$  的样本数与相邻两次迭代对比为最小值时, 迭代终止。最后, 取该次迭代的参数值作为 EM 算法的估计值输出。

## 5 实验结果及分析

### 5.1 实验数据集及评价指标

本文实验采用人工数据集<sup>[18-20]</sup>上的 5 个数据集(见表 1)和 UCI 机器学习数据库<sup>[21]</sup>中的 4 个数据集(见表 2)。这些数据集都是测试聚类算法的经典数据集。人工数据集是二维数据集, 可以直观地展示不同的聚类初始值选取方法对聚类结果的影响。所有实验的实验环境均为 Win10 64 bit 操作系统, Matlab R2018a 软件, 8 GB 内存, CPU Intel i7-9750H, GPU RTX2060。

表1 人工数据集

Table 1 Artificial datasets

数据集	样本数	特征数	类簇数
flame	240	2	2
Aggregation	788	2	7
Pathbased	300	2	3
Spiral	312	2	3
Jain	373	2	2

表2 UCI数据集

Table 2 UCI datasets

数据集	样本数	特征数	类簇数
iris	150	4	3
wine	178	13	3
WDBC	569	30	2
waveform	5000	21	3

UCI数据集中许多数据是真实的测量数据,iris数据集也称为鸢尾花卉数据集,wine数据集的13个属性是葡萄酒的13种化学成分,WDBC是乳腺癌数据集,waveform是噪声波形数据集,这些数据集在样本规模、特征数、类簇数等方面差别很大。

为了消除实验数据不同量纲对实验结果的影响,将其转化为无量纲的纯数值,以便于不同单位或量级的特征数据能够进行比较。本文采用min-max进行标准化处理。

$$x_{ij}^* = \frac{x_{ij} - \min_j}{\max_j - \min_j} \quad (15)$$

其中, $x_{ij}$ 表示第*i*个样本的第*j*个特征值, $\min_j$ 和 $\max_j$ 分别表示该特征的最小值和最大值。这样将特征数据归一化到 $[0,1]$ 区间。

采用两种度量指标对实验结果进行评价:聚类准确率(Accuracy, Acc)和调整兰德指数(Adjusted Rand Index, ARD<sup>[22]</sup>)。两种指标的取值上界为1,取值越大表示聚类结果越好。Acc是评价聚类结果的精度指标,需要计算正确聚类的样本数占样本总数的比例。ARI基于样本对计数,其中,RI计算样本预测值与真实值之间的相似度,RI的取值范围是 $[0,1]$ 。对于随机结果,ARI具有更高区分度,取值范围是 $[-1,1]$ ,值越大表示聚类结果和真实情况越吻合。

## 5.2 参数 $\alpha_k, \theta_k$ 的初始值

$\alpha_k$ 的初始值为各类相等的均值。对于参数 $\theta_k$ 的初始值,EM采用随机法选取初始类别中心。本文采用DPC算法得到初始类别中心。图2给出了随机法和DPC算法在flame数据集上选取的初始中心点。

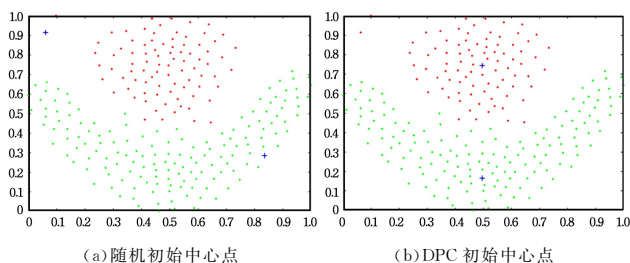


图2 在flame数据集上两种初始中心选取方法的对比

Fig. 2 Comparison of two initial center selection methods on flame dataset

从图2可以看出,随机选取的初始中心与真实的类别中

心偏差较大,而利用DPC算法得到的初始类别中心明显优于随机类别中心。图3给出了采用随机初始聚类中心方法、运行3次得到的3个随机初始聚类中心与DPC初始中心的聚类准确率对比。先采用*k*-means对3个随机初始中心进行优化,再与DPC进行对比。

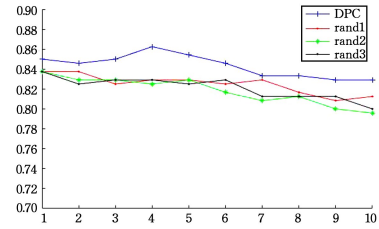


图3 迭代10次的聚类准确率

Fig. 3 Clustering accuracy of 10 iterations

图3给出了将EM算法迭代10次的各次聚类准确率。图3中,采用DPC选取初始中心后,GMM算法的准确率明显高于采用随机中心的传统算法。而3个随机初始中心的聚类准确率指标相差不大。图3还表明,将聚类准确率的均值及方差作为衡量算法的鲁棒性指标,DPC算法的优势明显。采用DPC算法后,EM估计的参数波动更小,提高了算法的鲁棒性。

随着迭代次数的增加,聚类准确率反而缓慢下降,原因是GMM算法的参数存在最优值,在达到最优值后继续迭代,反而会使错分的样本数增加,从而使得聚类的准确率下降。

## 5.3 采用相对熵的迭代优化

对于EM算法的迭代终止条件,本文不再采用相邻两次估计参数的距离,而是以相对熵小于某个阈值的样本数量作为衡量标准。当类别数大于两类时,则选取样本所有聚类概率中最大的两个概率,计算其相对熵。表3列出了在UCI的iris数据集上本文算法的聚类结果。

表3 iris数据集上本文算法的聚类结果

Table 3 Clustering results of our algorithm on iris dataset

阈值	1	2	3	4	5	6	7	8	9	10
$\delta=0.5$	14	9	9	8	4	1	4	3	1	1
$\delta=0.6$	17	10	10	8	5	2	4	3	3	2
$\delta=0.7$	18	10	10	8	5	2	4	3	2	2
$\delta=0.8$	20	13	11	9	5	7	4	3	3	3
Dist	0.009	0.007	0.016	0.006	0.008	0.012	0.013	0.012	0.010	0.010
Acc/%	89.33	93.33	96.67	97.33	99.33	98.67	98.00	98.00	96.67	97.33

实验将相对熵的阈值 $\delta$ 设置为0.5~0.8,考察选取不同的阈值对本文算法的影响,即本文算法对阈值参数 $\delta$ 的鲁棒性。用 $Dist = \|\hat{\alpha}_k - \hat{\alpha}_{k-1}\|$ 表示传统EM算法估计的相邻两个参数的距离,聚类结果采用聚类准确率指标Acc。实验中将EM算法迭代10次。

EM算法的初始参数是通过DPC算法得到的,本文算法的聚类准确率很高,最高达到了99.33%。这是目前在iris数据集上取得的最佳准确率指标,说明在服从高斯分布的数据集上,采用DPC+GMM可以获得最佳的聚类效果。

但并不是迭代的次数越多聚类的准确率就越高,最优值出现在第5次迭代中。实验计算了相邻两次迭代参数 $\alpha_k$ 的欧氏距离,可以看出该距离波动不大,最小值与其他距离

值相差很小。第4次迭代对应的聚类准确率并不是最优值,甚至不是次优值。因此,距离测度不能保证估计的参数值是最优的。

本文判断相对熵小于给定阈值  $\delta$  的样本数。当相对熵阈值  $\delta$  取值为 0.5~0.8 时,在第6次迭代时有3次取得了最小值,第5次迭代有1次得到了最小值。因此,采用本文的相对熵作为迭代终止条件存在明显的最小值,其所对应的聚类准确率指标3次是次优的98.67%,1次对应最优的99.33%。实验结果表明,本文的基于相对熵的迭代终止条件优于传统的基于距离的迭代终止条件,且算法对阈值  $\delta$  的取值并不敏感,算法的鲁棒性较好。

#### 5.4 在人工数据集上的对比实验

本文选取了5个二维人工数据集,这些数据集有的服从高斯分布,有的并不服从高斯分布。图4给出了采用DPC算法获得的初始聚类中心。

flame数据集的DPC选取结果见图2。图4给出了其余4个数据集的DPC选取结果。可以看出,在5个数据集上都得到了较准确的初始类别数与初始聚类中心点坐标。

从图2可以看出,flame数据集近似服从高斯分布。图4中,在Pathbased数据集的3个子类中,有2个子类服从高斯分布,Spiral数据集不服从高斯分布,Jain数据集近似服从高斯分布,而Aggregation数据集服从高斯分布。

可以预计,本文算法在Aggregation数据集上会取得较好的聚类结果,在其他3个数据集上的聚类结果为次优,而在Spiral数据集上的聚类准确率较低。

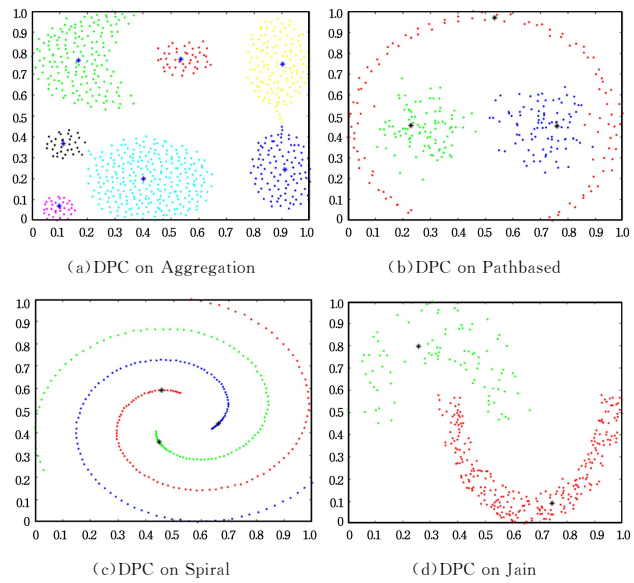


图4 DPC在4个人工数据集上得到的聚类中心  
Fig.4 Clustering centers obtained by DPC on four artificial datasets

从表4可以看出,当本文的相对熵阈值  $\delta$  取值为 0.5~0.8 时,本文算法在5个人工数据集上的ARI和Acc指标波动很小,说明本文算法对阈值参数  $\delta$  的取值并不敏感,算法的鲁棒性较好。本文算法在5个数据集上的ARI和Acc指标都好于传统的GMM算法。同时,本文算法在Aggregation数据集上取得了极高的聚类准确率,在Spiral数据集上的准确率最低,而在其他3个数据集上的结果为次优。

表4 人工数据集上的对比实验

Table 4 Comparative experiment on artificial datasets

算法阈值 数据集	$\delta=0.5$		$\delta=0.6$		$\delta=0.7$		$\delta=0.8$		GMM	
	ARI	Acc	ARI	Acc	ARI	Acc	ARI	Acc	ARI	Acc
flame	0.762	0.863	0.749	0.854	0.749	0.854	0.749	0.854	0.710	0.825
Aggregation	0.997	0.996	0.997	0.996	0.997	0.996	0.997	0.996	0.987	0.983
Pathbased	0.720	0.690	0.720	0.690	0.720	0.690	0.720	0.690	0.712	0.670
Spiral	0.496	0.349	0.496	0.349	0.496	0.349	0.496	0.349	0.331	0.323
Jain	0.771	0.868	0.771	0.868	0.759	0.860	0.759	0.860	0.751	0.845

#### 5.5 在UCI数据集上的对比实验

在UCI的4个数据集上,将本文算法与相关算法进行对

比,包括DPC算法和DF-DPC算法<sup>[14]</sup>、E-DPC算法<sup>[15]</sup>、传统的GMM算法及经典的基于划分的聚类算法k-means。实验采用Acc和ARI两个评价指标,实验结果如表5所列。

表5 UCI数据集上的对比实验

Table 5 Comparative experiment on UCI datasets

数据集	DPC		DF-DPC		E-DPC		GMM		k-means		本文算法	
	ARI	Acc	ARI	Acc	ARI	Acc	ARI	Acc	ARI	Acc	ARI	Acc
iris	0.720	0.887	0.561	0.727	0.862	0.933	0.856	0.866	0.660	0.825	0.982	0.986
wine	0.672	0.882	0.756	0.916	0.672	0.882	0.941	0.946	0.830	0.932	0.954	0.960
WDBC	0.471	0.845	0.451	0.838	0.471	0.845	0.875	0.933	0.826	0.928	0.884	0.938
Waveform	0.268	0.586	0.276	0.492	0.338	0.716	0.453	0.506	0.254	0.501	0.468	0.533

从表5可以看出,本文算法在前3个数据集上的聚类精度都是最高的,尤其是与DPC相关的3种算法相比,聚类准确率Acc及ARI指标的绝对值平均提高了近10%。同时,在这3个数据集上本文算法也显著优于传统GMM算法及k-means算法。

在Waveform数据集上,本文算法与对比算法的差别不

大,在Acc指标上甚至明显低于E-DCP算法。由于Waveform数据集的样本数据是噪声波形数据,这些数据显然并不服从高斯分布,因此本文算法不具备任何优势。

对于本文算法的算法复杂度,由于采用DPC进行初始中心的选取,而DPC算法首先需要计算数据集中任意两个样本间的欧氏距离,其时间复杂度为  $O(m * n^2)$ <sup>[16]</sup>,其中  $m$  为样

本特征个数,  $n$  为数据集样本个数。当处理海量高维数据时, 大量的高维欧氏距离计算会带来极大的时间开销, 可以采用文献[16]提出的相对邻域和剪枝策略(称为 RP-DPC 算法)。RP-DPC 算法具有与 DPC 相同的聚类效果, 并且时间性能显著优于 DPC 及其改进算法。

**结束语** 本文针对 EM 算法在估计高斯混合模型算法的参数时, 存在对参数的初始值敏感和迭代终止条件不能保证取得最优值的问题, 提出了采用密度峰值聚类算法 DPC 来初始化 EM 算法, 解决了 EM 算法的初始值敏感性问题。本文还提出采用相对熵小于某个阈值的样本数量作为 EM 算法的迭代终止条件, 实现了对参数估计值的优化选取。

在人工数据集和 UCI 数据集上的实验表明, 本文算法明显优于传统 GMM 算法。尤其是在数据集服从高斯分布的情况下, 本文算法取得了极佳的聚类效果。但当样本数据不服从高斯分布时, 本文算法的优势并不明显, 说明在具有不同分布特征的数据集上, 应采用不同的聚类策略。因此, 如何预估样本数据各个子类的分布特征成为进一步提高聚类准确率的关键。

### 参 考 文 献

[1] BENABDELLAH A C, BENGHABRIT A, BOUHADDOU I. A survey of clustering algorithms for an industrial context[J]. *Procedia Computer Science*, 2019, 148: 291-302.

[2] OYELADE J, ISEWON I, OLADIPUPO O, et al. Data Clustering: Algorithms and Its Applications[C]// 2019 19th International Conference on Computational Science and Its Applications (ICCSA). Saint Petersburg, Russia, 2019: 71-81.

[3] CHAI W Y, YANG F, YUAN S F, et al. Multi-class Gaussian Mixture Model and Neighborhood Information Based Gaussian Mixture Model for Image Segmentation[J]. *Computer Science*, 2018, 45(11): 272-277.

[4] ZOU C M, CHEN D. Unsupervised Anomaly Detection Method for High-dimensional Big Data Analysis[J]. *Computer Science*, 2021, 48(2): 121-127.

[5] MCQUEEN J. Some methods for classification and analysis of multivariate observations [C] // Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability. Los Angeles; University of California, 1967: 281-297.

[6] SHEN H, DUAN Z. Application Research of Clustering Algorithm Based on K-Means in Data Mining[C]// 2020 International Conference on Computer Information and Big Data Applications (CIBDA). Guiyang, China, 2020: 66-69.

[7] SAPKOTA N, ALSADOON A, PRASAD P W C, et al. Data Summarization Using Clustering and Classification: Spectral Clustering Combined with k-Means Using NFPH[C]// 2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon). Faridabad, India, 2019: 146-151.

[8] ZAKIR H, NASIM A, AHMADR B, et al. A dynamic K-means clustering for data mining[J]. *Indonesian Journal of Electrical Engineering and Computer Science*, 2019, 13(2): 521-526.

[9] ALEX R, ALESSANDRO L. Clustering by fast search and find of density peaks[J]. *Science*, 2014, 344: 1492-1496.

[10] DEMPSTERA P, LAIRDN M, RUBIND B. Maximum Likelihood from Incomplete Data via the EM Algorithm[J]. *Journal of the Royal Statistical Society*, 1997, 39(1): 1-38.

[11] YANG M S, LAI C Y, LIN C Y. A robust EM clustering algorithm for Gaussian mixture models[J]. *Pattern Recognition*, 2012, 45(11): 3950-3961.

[12] YUE J, WANG S T. Algorithm EM and Its Initialization in Gaussian Mixture Model Based Clustering[J]. *Microcomputer Information*, 2006, 11(22): 244-247.

[13] XIE J Y, GAO H C, XIE W X. K-nearest neighbors optimized clustering algorithm by fast search and finding the density peaks of a data set[J]. *Scientia Sinica Informationis*, 2016, 46(2): 258-280.

[14] WANG S L, WANG D K, LI C Y, et al. Clustering by fast search and find of density peaks with data field[J]. *Chinese Journal of Electronics*, 2016, 3(25): 397-402.

[15] ZHANG W K, LI J. Extended fast search clustering algorithm: widely density clusters, no density peaks [J]. arXiv: 1505.05160, 2015.

[16] JI X, YAO S, ZHAO P. Relative Neighborhood and Pruning Strategy Optimized Density Peaks Clustering Algorithm [J]. *ACTA Automatica Sinica*, 2020, 46(3): 562-575.

[17] YANG W, CAI L, WU F. Image segmentation based on gray level and local relative entropy two dimensional histogram[J]. *PLOS ONE*, 2020, 15(3): 1-9.

[18] GIONIS A, MANNILA H, TSAPARAS P. clustering aggregation[C]// Proceedings of ACM Transactions on Knowledge Discovery from Data. 2007, 1(1): 1-30.

[19] LIMIN F, ENZO M. Flame, a novel fuzzy clustering method for the analysis of DNA microarray data [J]. *BMC Bioinformatics*, 2007, 8(1): 3-17.

[20] CHANG H, YEUNG D Y. Robust path-based spectral clustering[J]. *Pattern Recognition*, 2008, 41(1): 191-203.

[21] LI C M. UCI machine learning repository [EB/OL]. <http://archive.ics.uci.edu/ml>.

[22] VINH N X, EPPS J, BAILEY J. Information theoretic measures for clustering comparison: is a correction for chance necessary? [C]// Proceedings of ICML'09. Montreal, 2009: 1073-1080.



**WANG Wei-dong**, Ph.D, associate professor. His main research interests include pattern recognition and intelligent information processing.