



## 全局数据驱动的模糊聚类有效性评价指标

唐益明, 刘子龙, 高健玮

引用本文:

唐益明, 刘子龙, 高健玮. 全局数据驱动的模糊聚类有效性评价指标[J]. *智能系统学报*, 2026, 21(3): 598-616.

TANG Yiming, LIU Zilong, GAO Jianwei. Global data-driven fuzzy cluster validity index[J]. *CAAI Transactions on Intelligent Systems*, 2026, 21(3): 598-616.

在线阅读 View online: <https://dx.doi.org/10.11992/tis.202507010>

## 您可能感兴趣的其他文章

### 基于结构相似性与模板校正的织物瑕疵检测方法

Fabric defect detection based on structural similarity and template correction

智能系统学报. 2020, 15(3): 475-483 <https://dx.doi.org/10.11992/tis.201810011>

### 基于可拓距的改进k-means聚类算法

Improved k-means algorithm based on extension distance

智能系统学报. 2020, 15(2): 344-351 <https://dx.doi.org/10.11992/tis.201811020>

### 基于模糊不一致对的多标记属性约简

Multi-label attribute reduction based on fuzzy inconsistency pairs

智能系统学报. 2020, 15(2): 374-385 <https://dx.doi.org/10.11992/tis.201905046>

### 面向一致性样本的属性约简

Attribute reduction over consistent samples

智能系统学报. 2019, 14(6): 1170-1178 <https://dx.doi.org/10.11992/tis.201905051>

### 重要度集成的属性约简方法研究

Research on ensemble significance based attribute reduction approach

智能系统学报. 2018, 13(3): 414-421 <https://dx.doi.org/10.11992/tis.201706080>

### 增量极坐标编码的贝赛尔曲线智能优化算法

Intelligent optimized Bezier curves based on incremental polar coordinate coding

智能系统学报. 2017, 12(6): 841-847 <https://dx.doi.org/10.11992/tis.201706076>

DOI: 10.11992/tis.202507010

网络出版地址: <https://link.cnki.net/urlid/23.1538.tp.20251222.1754.005>

# 全局数据驱动的模糊聚类有效性评价指标

唐益明, 刘子龙, 高健玮

(合肥工业大学计算机与信息学院, 安徽合肥 230601)

**摘要:** 现有模糊聚类有效性指标 (cluster validity index, CVI) 在处理具有噪声干扰以及簇规模差异较大的数据集时, 往往难以保持准确的评估性能。为此, 本文提出了一个全局数据驱动指标 (global data driven index, GDD), 该指标针对聚类结果中各簇规模的差异, 设计了一套针对簇内紧致度与簇间分离度的基于簇规模的加权机制, 以此来应对该差异对指标结果的影响。GDD 指标设计时采用簇内紧致度与簇间分离度比值的形式, 对于可能会出现 CVI 结果根据聚类结果数目单调递增的情况进行遏制; 在考虑模糊隶属度的簇内点的平均距离的基础上, 基于不同规模的簇占数据集中的比重不同, 对于较大的簇给予其更高的权重, 在计算单个簇内紧致度的基础上引入了该簇样本数占总样本的比值, 该加权机制将直接影响该簇紧致度结果对于加和后的总紧致度的贡献, 从而构建了更客观的簇内紧致度的表达; 综合考虑了所有类中心之间的均值与不同规模的簇占数据集的比重, 对于较大的簇给予更高的权重, 在计算每个聚类中心到聚类中心均值的距离的基础上加入了该簇样本数占总样本的比值, 从而锚定了上述距离结果占所有簇加和的距离结果的比重, 这种加权机制构建了更合理的簇间分离度的表达。实验结果显示, GDD 能够很好地适应各种模糊聚类算法, 而且在面对复杂结构和噪声时表现出较强的鲁棒性。本文提出的 GDD 指标可以在复杂结构与噪声环境下较好地完成对各类模糊聚类算法的评价。

**关键词:** 聚类算法; 模糊聚类; 聚类有效性指标; 模糊 C 均值算法; 加权紧致度; 加权分离度; 鲁棒性; 抗噪声

**中图分类号:** TP181; TN99 **文献标志码:** A **文章编号:** 1673-4785(2026)03-0598-19

中文引用格式: 唐益明, 刘子龙, 高健玮. 全局数据驱动的模糊聚类有效性评价指标 [J]. 智能系统学报, 2026, 21(3): 598-616.

英文引用格式: TANG Yiming, LIU Zilong, GAO Jianwei. Global data-driven fuzzy cluster validity index [J]. CAAI transactions on intelligent systems, 2026, 21(3): 598-616.

## Global data-driven fuzzy cluster validity index

TANG Yiming, LIU Zilong, GAO Jianwei

(School of Computer Science and Information, Hefei University of Technology, Hefei 230601, China)

**Abstract:** Existing cluster validity index (CVI) for fuzzy clustering often struggle to maintain accurate evaluation performance when handling datasets containing noise interference or exhibiting significant differences in cluster sizes. To address these limitations, this study proposes a global data-driven (GDD) index. The GDD index incorporates a scale-aware weighting mechanism for intra-cluster compactness and inter-cluster separation to mitigate the adverse impact of imbalanced cluster sizes on validity assessment. First, the GDD index adopts a ratio-based formulation of intra-cluster compactness to inter-cluster separation. This design prevents the undesirable monotonic increase of the index value as the number of clusters grows. Second, to obtain a more objective measure of intra-cluster compactness, the index computes the average distance among data points within each cluster, incorporating fuzzy membership degrees. Crucially, recognizing that clusters of different scales contribute unequally to the overall dataset structure, larger clusters are assigned higher weights. Specifically, the ratio of the number of samples in each cluster to the total number of samples is introduced into the compactness calculation. This weighting scheme directly influences each cluster's contribution to the overall compactness, thereby enhancing representational fairness and accuracy. Third, for inter-cluster separation, the index comprehensively considers both the centroid distribution and the relative sizes of different clusters. Rather than treating all centroids equally, the index assigns higher weights to centroids of larger clusters. When computing the distance from each cluster centroid to the mean of all centroids, the sample-size ratio of the corresponding cluster is incorporated. This adjustment anchors the contribution of each centroid's distance to the total separation measure, resulting in a more reasonable and balanced expression of inter-cluster separation. To evaluate the effectiveness and robustness of the GDD index, extensive experiments were conducted using three representative fuzzy clustering algorithms. Experimental results demonstrate that the GDD index consistently identifies the optimal number of clusters with high accuracy, adapts well across various fuzzy clustering frameworks, and demonstrates strong robustness in challenging scenarios with noise and highly imbalanced cluster sizes. The proposed index provides a more comprehensive and reliable evaluation of fuzzy clustering algorithms in complex, noisy environments.

**Keywords:** clustering algorithms; fuzzy clustering; clustering validity index; fuzzy C-means algorithms; weighted compactness; weighted separation; robustness; noise robustness

收稿日期: 2025-07-07. 网络出版日期: 2025-12-23.

基金项目: 国家自然科学基金项目 (62576130, 62176083).

通信作者: 高健玮. E-mail: [jwgao810@163.com](mailto:jwgao810@163.com).

聚类作为一种无监督的机器学习方法<sup>[1-2]</sup>, 其在数据分析和数据清洗方向的能力适用于处理不

同领域和类型的复杂数据<sup>[3]</sup>,尤其在数据挖掘、模式识别和图像处理等领域得到了广泛的研究和应用。聚类是将一个未标识的数据集划分成不同类簇的过程,尽可能使相似的数据被划分到同一个簇中,不相似的数据划分到不同的簇中<sup>[4]</sup>。聚类在研究和应用过程中有两个重要问题需要考虑:一是聚类算法,通过聚类算法来划分数据集,得到聚类结果;二是聚类验证<sup>[5]</sup>,通过评估聚类结果进而评估聚类算法的效果。

随着计算机科学的发展,尤其是大数据和人工智能的兴起,聚类算法得到了长足的发展和广泛的应用。经典的K均值聚类算法(K-Means)<sup>[6]</sup>和层次聚类方法<sup>[7]</sup>最早被提出,随后催生了基于密度<sup>[8-9]</sup>、网格和模型等的聚类方法<sup>[10-12]</sup>。以上这些通常都是硬聚类算法,在划分过程中要求每个数据点只能划分到一个唯一的簇中。这种非此即彼的划分方式难以适用某些不确定的数据集,在这种特殊的情况下,数据点可以同时归属于多个簇,具有一定的模糊性。这种需求催生了模糊聚类算法。Dunn<sup>[13]</sup>于1973年提出了模糊C均值聚类算法(fuzzy C-means, FCM),后续Bezdek等<sup>[14]</sup>在此基础上发展并推广了FCM算法。FCM是模糊聚类的代表算法,通过优化隶属度函数来计算每个数据点对于簇的隶属度。此算法成为模糊聚类的重要基础,在处理不确定性数据方面有了显著进展。但FCM也存在缺点,如对噪声数据敏感。针对这一问题,Krishnapuram等<sup>[15]</sup>提出了可能性C均值聚类算法(possibilistic C-means, PCM),PCM引入了可能性隶属度,而非单一的隶属度,放宽了对隶属度的限制,赋予噪声点更低的隶属度,从而增强了聚类的鲁棒性。后续随着模糊聚类应用需求的增加,研究人员提出了更多种类的模糊聚类算法。例如,Pal等<sup>[16]</sup>提出的模糊可能性C均值聚类算法(fuzzy possibilistic C-means, FPCM)、Antoine等<sup>[17]</sup>提出的可能性模糊C均值聚类算法(possibilistic fuzzy C-means, PFCM)、Zhang等<sup>[18]</sup>提出的基于核的模糊聚类算法(kernel-based fuzzy C-means, KFCM),KFCM将核空间的概念引入到模糊聚类算法中,增强了算法的鲁棒性。

聚类验证是对聚类的结果进行评估,以验证聚类算法的划分是否正确。聚类验证常用的方法是通过聚类有效性指标(cluster validity index, CVI)<sup>[19]</sup>来对聚类算法划分的结果进行计算,对比不同划分状态下的计算结果,选取最优的划分结果。CVI大体上可以划分为3种<sup>[20]</sup>,即外部有效性<sup>[21]</sup>、内部有效性<sup>[22]</sup>和相对有效性<sup>[23]</sup>,文中涉及

的指标主要是内部有效性指标。内部有效性指标通过聚类结果的内部结构来评估聚类质量,不依赖于外部标签。常见的内部有效性指标有Calinski等<sup>[24]</sup>提出的CH(Calinski-Harabasz)指标、Davies等<sup>[25]</sup>提出的DB(Davies-Bouldin)指标,还有Dunn<sup>[13]</sup>提出的Dunn指标,但Dunn指标对于环状或者线性数据集的适应效果较差。此外还有,Xie等<sup>[26]</sup>的XBI(Xie-Beni)指标、Fukuyama等<sup>[27]</sup>提出的FSI(Fukuyama-Sugeno)指标、Wu等<sup>[28]</sup>等提出的WLI(Wu-and-Li)指标、Liu等<sup>[29]</sup>提出的不平衡指数(imbalanced index, IMI)、Mittal等<sup>[30]</sup>提出的SMI(Saraswat-and-Mittal)指标、Zhu等<sup>[31]</sup>提出的基于方差的聚类有效性(the variance based clustering validity, VCVI)指标、Maulikh等<sup>[32]</sup>提出的MB(Maulik-Bandyopadhyay)指标和Tang等<sup>[33]</sup>提出的三重中心关系(triple center relation, TCR)指标等,这些指标都在CVI发展进程中起到一定作用。

上文提到的这些CVI在聚类质量的评估上起到重要作用,是聚类分析过程中不可或缺的一部分,但是它们也存在一些缺点和局限性:1)部分CVI对于簇的形状较为敏感,且对于现实数据中常出现的多噪声情况难以有合适的应对方法;2)基于FCM算法提出的CVI会因为FCM的错误分类导致出现错误的结果。如FCM算法的均匀效应会导致FCM算法在处理不平衡数据集时会出现错误结果<sup>[34-35]</sup>。

考虑以上问题,本文提出了一种新的CVI,称之为全局数据驱动指标(global data driven index, GDD)。该CVI共有2个部分组成。GDD全面考虑了簇内外的数据特征,同时2个部分的组合,将簇内紧致性和簇间分离性结合起来分析,从而更好地度量聚类结果的合理性。

针对提出的CVI,本文采用UCI数据集<sup>[36]</sup>、人工数据集和Olivetti Face数据集共19个数据集以及11个CVI进行了对比实验。为了增加实验的可信度,分别选取了一些含有噪声的数据集和一些数据分布不均匀的数据集进行对比实验。与此同时,为了验证GDD能否同时适应多种模糊聚类算法,本文采用了FCM、PFCM和KFCM这3种不同的模糊聚类算法。通过不同的模糊聚类算法进行实验验证,证明GDD有更强的适应力和更广泛的实用价值。此外,本文论证了GDD的收敛性,以证明新CVI在理论上的正确性。

## 1 部分内部有效性评价指标分析

内部有效性评价指标通过分析聚类结果中簇

的结构来衡量聚类划分的质量, 无需使用外部的分类标签来对比判断。具体指标主要关注簇内紧致性(foc)和簇间分离性(fos)两个方面。这里介绍较为经典的 5 种有效性指标。

1) Dunn 指标

Dunn 指标是一种经典且有广泛应用的内部评价指标, 其对非球形数据分布有较好的适应性。但是其对噪声敏感且计算较复杂。计算公式为

$$Dunn^{(+)} = \frac{\min_{1 \leq i < j \leq K} \{dis(C_i, C_j)\}}{\max_{1 \leq k \leq K} \{diam(C_k)\}}$$

式中:  $dis(C_i, C_j) = \min_{x_i \in C_i, x_j \in C_j} \{\|x_i - x_j\|\}$ ,  $diam(C_k) = \max_{x_i, x_j \in C_k} \{\|x_i - x_j\|\}$ 。  $C_i$  为第  $i$  个簇。

2) XBI 指标

XBI 指标通过计算簇内部数据点到聚类中心的距离来度量簇内紧致性, 计算所有聚类中心之间的距离最小值来度量簇间分离性, 考虑因素比较全面。但是其单独选用聚类中心之间的距离最小值来度量簇间分离性, 会导致对不规则数据集处理的效果不佳。同时, 指标容易随着聚类簇数增加而单调变化。其计算公式为

$$XBI^{(-)} = \frac{\sum_{k=1}^K \sum_{i=1}^N \mu_{ik}^m \|x_i - v_k\|}{N \times \min_{i \neq j} \{\|v_i - v_j\|\}}$$

式中:  $N$  为样本总数,  $K$  为聚类数。这里的距离通常采用欧氏距离。

3) WLI 指标

WLI 指标在度量簇内紧致性的公式中加入了

$$TCR^{(-)} = \frac{\sum_{k=1}^K \frac{\sum_{i=1}^N u_{ik}^2 \|x_i - v_k\|^2}{\sum_{i=1}^N \max_{1 \leq k \leq K} u_{ik}}}{\frac{N}{K-1} \left( \min_{i \neq j} \{\|v_i - v_j\|^2\} \times \text{mean}_{i \neq j} \{\|v_i - v_j\|^2\} \times \sum_{k=1}^K \|v_k - \bar{v}\|^2 \right)}$$

式中:  $N$  为样本数,  $K$  为聚类数,  $\min$  和  $\text{mean}$  分别为簇中心之间距离的最小值与均值。通常采用欧氏距离进行计算。

## 2 全局数据驱动指标

GDD 指标主要由簇内紧致性和簇间分离性两个部分组成。在参考 XBI 指标的整体结构基础上, GDD 指标对簇内紧致性和簇间分离性的表达进行一些改进。

### 2.1 簇内紧致性表达

GDD 指标的簇内紧致性的表达分为两个部

模糊隶属度的和, 更好地度量簇内数据的紧致性, 同时采用聚类中心之间距离的最小值和中值来度量簇间分离性, 对分离性的度量更加全面。其计算公式为

$$WLI^{(-)} = \frac{\sum_{k=1}^K \left( \frac{\sum_{i=1}^N \mu_{ik}^2 \|x_i - v_k\|^2}{\sum_{i=1}^N \mu_{ik}} \right)}{\min_{i \neq j} \{\|v_i - v_j\|^2\} + \text{median}_{i \neq j} \{\|v_i - v_j\|^2\}}$$

式中  $\min$  与  $\text{median}$  为簇中心之间的距离的最小值与中位数值。通常采用欧氏距离进行计算。

4) IMI 指标

IMI 指标在簇间分离性的度量上加入了不同簇之间的数据点数量比值, 考虑了不同簇之间的不平衡比, 对分离性的刻画更加准确。其公式为

$$IMI^{(-)} = \frac{\sum_{k=1}^K \frac{\sum_{i=1}^N u_{ik}^q \|x_i - v_k\|}{\sum_{i=1}^N u_{ik}}}{\min_{i \neq j} \{\delta_{ij} \|v_i - v_j\|^2\} + \text{median}_{i \neq j} \{\delta_{ij} \|v_i - v_j\|^2\}}$$

式中:  $\delta_{ij} = f_i / f_j, f_i > f_j, f_i = \sum_{n=1}^N u_{ni}$ , 可以看出  $\delta_{ij}$  是簇  $C_i$  和簇  $C_j$  的不平衡比。

5) TCR 指标

TCR 指标通过计算聚类中心之间距离最小值和均值, 以及聚类中心方差来度量簇间分离性, 三重因素相互作用, 互相平衡, 更好地刻画了分离性。其计算公式为

分。第一部分 foc 为簇内所有数据点之间的平均距离, 在此基础上, 考虑到模糊聚类中的模糊隶属度, 在计算数据点之间的欧氏距离时加上两点当前簇的模糊隶属度, 定义的公式为

$$f_{oc1} = \frac{\sum_{1 \leq i < j \leq n_k} u_{ik} u_{jk} \|x_i - x_j\|}{n_k(n_k - 1)}$$

式中:  $n_k$  为第  $k$  个簇的数据点数目,  $u_{ik}$  为在第  $k$  个簇中第  $i$  个样本的模糊隶属度,  $x_i, x_j$  分别为该簇中两个互不相同的样本。距离计算采用欧氏距离。

第二部分  $f_{oc}$  为单个簇内部数据点的标准差,  $\bar{x}$  表示簇内样本的样本均值, 表达式为

$$f_{oc2} = \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} \|x_i - \bar{x}\|^2}$$

将两个部分相加得到对单个簇内部紧致性的表达式, 同时考虑到不同的簇所含数据点数量的不同, 其在整个数据集中所占的比重不同。所以每个簇对整体数据紧凑程度的贡献不应相同, 所占比重高的簇应当在整体数据集紧凑程度表达上有更高的权重。因此整体紧致性表达式由每个簇内部紧致性表达式加权相加得来, 权重为当前簇的数据点数目占总体数据点数量的比值。最终表达式为

$$f_{oc} = \sum_{k=1}^K \frac{n_k}{N} \left( \frac{\sum_{1 \leq i < j \leq n_k} u_{ik} u_{jk} \|x_i - x_j\|}{\frac{n_k(n_k - 1)}{2}} + \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} \|x_i - \bar{x}\|^2} \right)$$

使用加权距离和来衡量数据紧凑程度能更好地反映簇内的数据几何结构, 也能更好应对形状各异的数据集。此外标准差的加入能增强 GDD 指标对簇内数据点紧致性和离散程度的评估能力。两个部分的结合使得 GDD 指标对簇内紧致性的评估更加全面。

### 2.2 簇间分离性表达

GDD 指标的簇间分离性的表达分为两个部分。第一部分  $f_{os}$  以所有聚类中心到聚类中心均值距离之和为基底, 同时考虑到不同聚类中心所在簇所含数据点数量的不同, 其在整个数据集中所占的比重不同。所以每个聚类中心所在簇对整体数据的贡献不应相同, 所占比重高的簇应当在整体数据集上有更高的权重。因此, 在每个聚类中心到聚类中心均值的距离上加上对应的权重。权重为当前聚类中心所在簇的数据点数目占总体数据点数量的比值。最终表达式为

$$f_{os1} = \sum_{k=1}^K \left( \frac{n_k}{N} \|v_k - \bar{v}\| \right)$$

式中:  $\bar{v}$  为所有簇中心的均值;  $N$  为样本总数。同样采用欧氏距离进行计算。

第二部分  $f_{os}$  为所有聚类中心之间距离的均值, 表达式为

$$f_{os2} = \text{mean}_{i \neq j} \{\|v_i - v_j\|\}$$

将两个部分相加即得到 GDD 指标对于簇间分离性的刻画, 表达式为

$$f_{os} = \sum_{k=1}^K \left( \frac{n_k}{N} \|v_k - \bar{v}\| \right) + \text{mean}_{i \neq j} \{\|v_i - v_j\|\}$$

关于簇间分离性刻画的第一部分, 加权距离和在整体上对所有聚类中心与中心点的离散程度和分离性做评估, 权重的引入进一步考虑了不同大小的簇对整体数据分离程度的影响。关于第二部分, 聚类中心之间距离均值从另一个方面对所有聚类中心之间的分离程度做全面的评估。两者的互相结合使得 GDD 指标对簇间分离性的刻画更加全面深刻。

综上所述, GDD 指标由两个部分组成: 第一部分  $f_{oc}$  用来刻画簇内数据点之间的紧致性; 第二部分  $f_{os}$  用于刻画聚类中心之间分离性, 也就是簇间分离性。

$$\text{GDD}(K)^{(-)} = \frac{f_{oc}}{f_{os}} = \frac{\sum_{k=1}^K \frac{n_k}{N} \left( \frac{\sum_{1 \leq i < j \leq n_k} u_{ik} u_{jk} \|x_i - x_j\|}{\frac{n_k(n_k - 1)}{2}} + \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} \|x_i - \bar{x}\|^2} \right)}{\sum_{k=1}^K \left( \frac{n_k}{N} \|v_k - \bar{v}\| \right) + \text{mean}_{i \neq j} \{\|v_i - v_j\|\}} \quad (1)$$

### 2.3 CVI 有效性证明

下面将对 GDD 指标的有效性做出证明, 其中用到了 Dunn 指标作为参照。其公式为

$$\text{Dunn}^{(+)} = \frac{\min_{1 \leq i < j \leq K} \{\text{dis}(C_i, C_j)\}}{\max_{1 \leq k \leq K} \{\text{diam}(C_k)\}}$$

其中  $\text{dis}(C_i, C_j)$  和  $\text{diam}(C_k)$  定义分别为

$$\text{dis}(C_i, C_j) = \min_{x_i \in C_i, x_j \in C_j} \{\|x_i - x_j\|\}$$

$$\text{diam}(C_k) = \max_{x_i, x_j \in C_k} \{\|x_i - x_j\|\}$$

Dunn 指标的论文中证明当  $\text{Dunn} > 1$  时, Dunn 指标会指出正确的聚类数目并且 Dunn 指标的数值越大标识当前的聚类效果越好。分析 GDD 指标的构成, 猜测当在正确的聚类划分结构下, Dunn 指标变得足够大, GDD 指标会变得足够小。接下来将对这一猜测进行证明。

假设在聚类数目为  $K$  时, 对数据集进行划分得到划分结果。分别使用 GDD 指标和 Dunn 指标对结果进行评估, 得到评价数值  $\text{GDD}(K)$  和  $\text{Dunn}(K)$ 。

**定理 1** 设  $k \in \{2, 3, \dots, N-1\}$ , 且  $\max_{1 \leq k \leq K} \{\text{dia}^2(C_k)\} \geq 1$  成立。同时,  $\mu_{ik} (1 \leq i \leq N, 1 \leq k \leq K)$  是模糊隶属度, 而  $\omega_{ik} (1 \leq i \leq N, 1 \leq k \leq K)$  是相应硬划分的隶属度,

其公式为

$$\omega_{ik} = \begin{cases} 1, & k = \operatorname{argmax}_{1 \leq k' \leq K} \{\mu_{ik'}\} \\ 0, & \text{其他} \end{cases}$$

则可得

$$\text{GDD}(K) \leq \frac{2}{\text{Dunn}(K)} \quad (2)$$

**证明** 设  $K$  类划分为数据集  $X$  的优化划分, 其中  $X = \{x_i | 1 \leq i \leq N\}$ , 聚类中心为  $v_k (1 \leq k \leq K)$ , 隶属度为  $\mu_{ik} (1 \leq i \leq N, 1 \leq k \leq K)$ 。  $n_k$  为第  $k$  个簇内部数据点的数目。推导  $f_{oc}$  和  $\text{diam}(C_k)$  的大小关系可得公式为

$$\text{diam}(C_k) = \max_{x_i, x_j \in C_k} \{\|x_i - x_j\|\} \geq \frac{\sum_{1 \leq i < j \leq n_k} \|x_i - x_j\|}{n_k(n_k - 1)} = f_{oc1} \quad (3)$$

$$\text{diam}(C_k) = \max_{x_i, x_j \in C_k} \{\|x_i - x_j\|\} \geq \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} \|x_i - \bar{x}\|^2} = f_{oc2} \quad (4)$$

结合式 (3) 和式 (4), 且聚类结果中所有的簇都满足式 (3) 和式 (4), 所以对式 (3) 和式 (4) 的加权和也满足。推出下式:

$$f_{oc} = \sum_{k=1}^K \frac{n_k}{N} \left( \frac{\sum_{1 \leq i < j \leq n_k} u_{ik} u_{jk} \|x_i - x_j\|}{n_k(n_k - 1)} + \sqrt{\frac{1}{n_k} \sum_{i=1}^{n_k} \|x_i - \bar{x}\|^2} \right) \leq 2 \max_{1 \leq k \leq K} \{\text{diam}(C_k)\} \quad (5)$$

同时, 由于

$$f_{os} \geq \operatorname{mean} \left\{ \left\| \|v_i - v_j\| \right\}^2 \right\} \geq \operatorname{dis}(C_i, C_j) \geq \min_{1 \leq i < j \leq K} \operatorname{dis}(C_i, C_j) \quad (6)$$

结合式 (5)、(6) 以及  $\max_{1 \leq k \leq K} \{\text{dia}^2(C_k)\} \geq 1$ , 可推出:

$$\text{GDD}(K) = \frac{f_{oc}}{f_{os}} \leq \frac{2 \max_{1 \leq k \leq K} \{\text{diam}(C_k)\}}{\min_{1 \leq i < j \leq K} \{\operatorname{dis}(C_i, C_j)\}} \leq \frac{2}{\text{Dunn}(K)}$$

至此证明完毕。

从定理 1 的式 (2) 可知, 在正确的聚类划分结构下, Dunn 指标变得足够大, GDD 指标会变得足够小, 说明对应的聚类划分效果是非常理想的。

### 2.4 CVI 时间复杂度分析

时间复杂度是衡量一个 CVI 面对由聚类算法划分后的数据集的处理效率的重要参照。一个优秀的 CVI 往往在拥有优秀识别正确率的基础上尽可能的达到更小的时间复杂度。时间复杂度因为有对 CVI 的实用性评估的功效, 需要对其进行分析。

分析式 (1), 将 GDD 分成以下部分进行时间复杂度的讨论。

1) 计算全局质心。对于全局质心, 有

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

对此, 需要遍历数据集中所有共  $n$  个点, 每个点  $d$  维, 所以该部分的时间复杂度为  $O(nd)$ 。

2) 计算模糊加权平均距离。对于每个簇  $j$  分别进行处理。每个簇的簇大小  $|C_j| = n_j \leq \tau$ , 其中  $\tau$  为所有簇中最大的样本数量, 所需计算的所有数据点对共  $O(\tau^2)$  对, 每对计算距离  $\|x_i - x_j\|$  的复杂度为  $O(d)$ 。每对需乘其隶属度  $O(1)$ 。因此该部分的总计算量为  $\sum_{j=1}^k O(\tau^2 d)$ , 考虑其最坏时间复杂度为  $O(K\tau^2 d)$ , 其中  $K$  为簇的数量。

3) 计算簇内标准差之和。计算每个簇的质心  $\bar{x}_j$  所需的时间复杂度为  $O(\tau d)$ , 接着计算  $\sum_{x_i \in C_j} \|x_i - \bar{x}_j\|^2$ , 其时间复杂度为  $O(\tau d)$ , 这只是单个簇的复杂度, 该部分总的时间复杂度为  $\sum_j O(\tau d) = O(nd)$ 。

4) 计算加权类间距离。由于在以上部分已经计算了全局质心  $\bar{x}$  和每个簇的质心  $\bar{x}_j$ , 因此该部分计算  $\|\bar{x}_j - \bar{x}\|^2$  的时间复杂度为  $O(d)$ , 这同样是单个簇的复杂度, 一共有  $K$  个簇, 从而类间距离部分总共  $O(Kd)$ , 又因为需要遍历数据集内所有共  $n$  个点, 这部分需要  $O(nd)$  因此该部分的总复杂度为  $O(Kd + nd)$ 。

5) 计算簇中心间的平均距离。单个簇中心间的距离  $\|\bar{x}_i - \bar{x}_j\|$  的时间复杂度为  $O(d)$ , 而该簇中心对一共有  $K^2$  对, 因此该部分的总复杂度为  $O(K^2 d)$ 。

至此, 可以得到 GDD 指标的时间复杂度为  $O(K\tau^2 d + nd + Kd + nd + K^2 d)$ , 而当聚类正确的情况下聚类数  $K$  应当远小于样本数  $n$ , 因此 GDD 的时间复杂度最终记为  $O(K\tau^2 d + nd)$ 。

表 1 中给出了文中参与实验的 12 个指标的时间复杂度。

表 1 12 个指标的时间复杂度  
Table 1 Time complexity of 12 indicators

序号	CVI	时间复杂度
1	CH	$O(K\tau d)$
2	Dunn	$O(K^2 \tau^2 d)$
3	DB	$O(K\tau^2 d)$
4	MB	$O(Knd)$
5	IMI	$O(Kn(d + K))$
6	XBI	$O(Knd)$
7	VCVI	$O(K\tau d)$

续表 1

序号	CVI	时间复杂度
8	FSI	$O(Knd)$
9	WLI	$O(Knd)$
10	SMI	$O(K^2\tau^2d + Knd)$
11	TCR	$O(Kn(d + K))$
12	GDD	$O(K\tau^2d + nd)$

从表 1 中可以看出, GDD 的时间复杂度为  $O(K\tau^2d + nd)$ , 属于高开销指标, 同为高开销指标的还有 Dunn 指标与 SMI 指标, 而其余指标在时间复杂度上均优于上述三者, 这意味着 GDD 指标在实际使用时可能需要更多的时间来处理数据并得出结果。

### 3 GDD 指标的仿真与对比实验

本文选取 3 个模糊聚类算法、11 个对比 CVI 以及 19 个数据集来做对比实验。3 个模糊聚类算法分别是 FCM 算法、PFCM 算法和 KFCM 算法; 11 个对比 CVI 分别是 CH、Dunn、DB、MB、IMI、XBI、VCVI、FSI、WLI、SMI 和 TCR 指标; 19 个数据集分别是 9 个 UCI 数据集、9 个人造数据集和 Olivetti Face 数据集。实验的软件环境如下: 操作系统是 Windows 11OS, 编程软件是专业数学分析软件。实验的硬件环境如下: CPU 是 Intel(R) Core(TM)i9-12900 KF, 显卡是 NVIDIA GeForce RTX3090, 内存是 64 GB。

#### 3.1 数据集与对比指标

在本实验中, 主要采用 3 类数据集进行研究, 分别为 UCI 数据集、人工合成数据集以及 Olivetti Face 数据集。本实验选取 UCI 数据集中的 9 个数据集, 分别是 Zoo、Hayes-Roth、Iris、Glass、Dermatology、Breast Cancer、Balance Scale、Libras、Letter 数据集。其中 Zoo 数据集来自动物数据, 数据维度为 16 维, 共有 101 个样本点, 正确的聚类数目为 7 类; Hayes-Roth 数据集来自社科研究, 数据维度为 4 维, 共有 132 个样本点, 正确的聚类数目为 3 类; Iris 数据集来自不同品种的鸢尾花数据, 数据维度为 4 维, 共有 150 个样本点, 正确的聚类数目为 3 类; Glass 数据集来自具有不同元素含量的玻璃的数据, 数据维度为 9 维, 共有 214 个样本点, 正确的聚类数目为 6 类; Dermatology 数据集来自一种疾病的医疗数据, 数据维度为 34 维, 共有 366 个样本点, 正确的聚类数目为 6 类; Breast Cancer 数据集来自肿瘤数据, 用于预测肿瘤的性质, 数据维度为 30 维, 共有 569 个样本点, 正确的

聚类数目为 2 类; Balance Scale 数据集来自模拟心理学实验结果, 数据维度为 4 维, 共有 625 个样本点, 正确的聚类数目为 3 类; Libras 数据集来自人手部动作数据, 数据维度为 90 维, 共有 360 个样本点, 正确的聚类数目为 15 类; Letter 数据集来自字符图像数据, 数据维度为 16 维, 共有 20 000 个样本点, 正确的聚类数目为 26 类。为了避免实验的数据集的类型单一, 本实验同样选取部分人造数据集进行实验来增强实验的可信性。人工数据集作为评估聚类算法在不同结构、密度、维度和噪声条件下性能的工具, 其允许研究者精确控制簇的结构、形状、密度、重叠度及噪声水平的特性, 能够针对不同的被测对象设计对应的测试内容, 从而系统地评估被测对象的鲁棒性、可拓展性和对特定模式的识别能力。人工数据集常常通过高斯分布采样或者数据点, 再通过高斯噪声引入坐标扰动, 通过不同的设计与配比得到满足要求的数据集合。其中包含有 9 个人造数据集, 分别为 Data\_60、Data\_150、Circle、Jain、X8D5K、Data\_77、E6、Dim\_128 和 Dim\_256 数据集。Data\_60 数据集的数据维度为 2 维, 共有 60 个样本点, 正确的聚类数目为 3 类; Data\_150 数据集的数据维度为 2 维, 共有 150 个样本点, 正确的聚类数目为 3 类; Circle 数据集的数据维度为 2 维, 共有 150 个样本点, 正确的聚类数目为 2 类; Jain 数据集的数据维度为 2 维, 共有 373 个样本点, 正确的聚类数目为 2 类; X8D5K 数据集的数据维度为 8 维, 共有 1 000 个样本点, 正确的聚类数目为 5 类; Data\_77 数据集的数据维度为 2 维, 共有 1000 个样本点, 正确的聚类数目为 7 类; E6 数据集的数据维度为 2 维, 共有 8 537 个样本点, 正确的聚类数目为 4 类; Dim\_128 数据集的数据维度为 128 维, 共有 1 024 个样本点, 正确的聚类数目为 16 类; Dim\_256 数据集的数据维度为 256 维, 共有 1024 个样本点, 正确的聚类数目为 16 类。

同时, 本实验选取 Olivetti Face 数据集, 来验证实验在图像数据上的准确性。实验前对 Olivetti Face 数据集进行切割处理, 处理后的数据集包含 40 类人脸图像, 每类 10 张图像, 10 张图像是同一个人不同状态下的面部图像, 每张图像大小为 4096 个像素点。

图 1 给出了部分上述数据集在正确聚类数目下的分布情况。其中维度大于 2 的数据集采用 TSNE 方法降维到 3 维后给出。Olivetti Face 数据集的缩略图如图 2 所示。

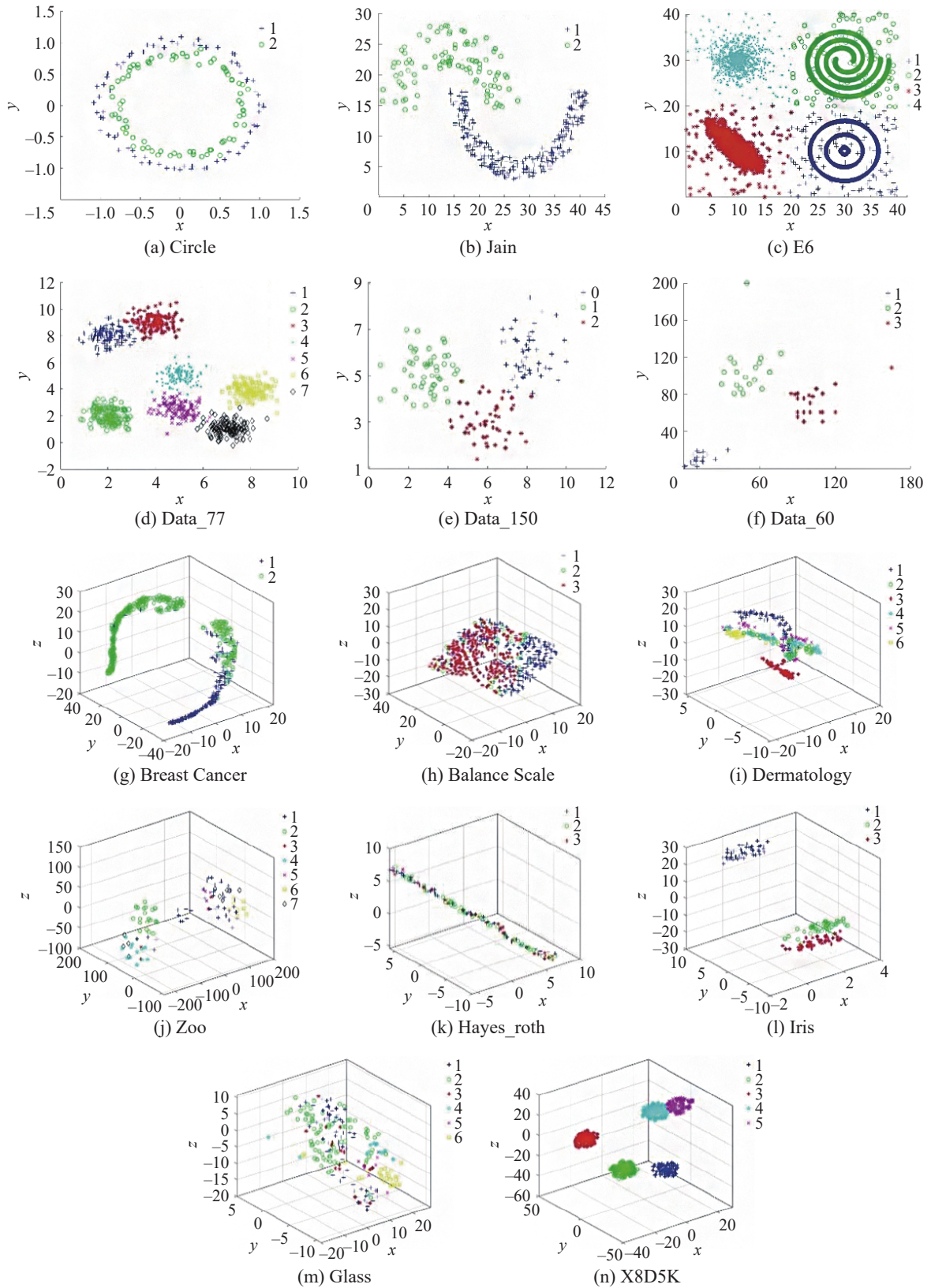


图 1 数据集的分布

Fig. 1 Distribution of datasets

本实验选取另外 11 个 CVI 作为对比指标, 分别是 CH、Dunn、DB、MB、IMI、XBI、VCVI、FSI、WLI、SMI 和 TCR 指标。通过与各种不同 CVI 进行对比实验, 来验证提出的 GDD 指标具有更强

的适应性和更高的准确性。

### 3.2 对比实验

本实验采用 FCM 算法、PFCM 算法和 KFCM 算法对 CVI 的准确性进行验证。在正确聚类数目为

2~10 的数据集上, 每个数据集进行 10 轮实验, 每轮实验的  $K$  值从 2 变到 10; 在正确聚类数目为 2~30 的数据集上, 每个数据集进行 30 轮实验, 每轮实验的  $K$  值从 2 变到 30。实验结束后, 收集所有实验数据, 并根据不同 CVI 最优值的分布统计各轮实验中 CVI 对应的最优聚类数目。表 2 给出了 UCI 数据集在上述 3 种算法下的实验结果, 其中的数据表示各轮实验中最优聚类数目及其出现的频次。例如, 数据“2<sup>7</sup>3<sup>3</sup>”表示在 10 轮实验中, 聚类数目为 2 的最优值出现了 7 次, 而聚类数目为 3 的最优值出现了 3 次。

表 2 第一部分是 FCM 算法下 CVI 的实验结果。以其中的 Dermatology 数据集为例, Dermatology 数据集的正确聚类数目为 6 类, MB 指标 10 轮的结果中 7 轮是 4 类为最优值, 2 轮是 5 类为最优值, 1 轮是 6 类为最优值, 其他 CVI 不再赘述。统计出现次数最多的最优值聚类数目作为该指标的评价结果, MB 指标结果为 4 类, WLI 指标结果为 5 类, GDD 指标结果为 6 类。可以看出, MB 指标和 WLI 指标的结果不正确, GDD 指标正确得出了 Dermatology 数据集的聚类数目。

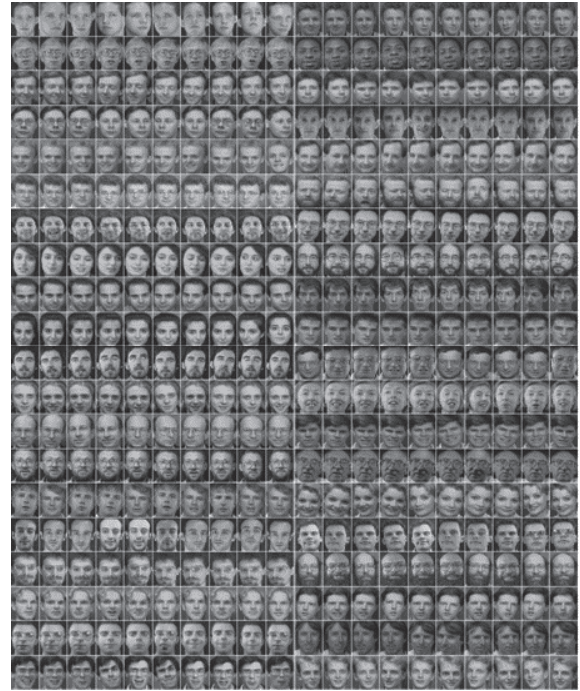


图 2 Olivetti Face 数据集  
Fig. 2 Olivetti Face dataset

表 2 UCI 数据集的实验结果  
Table 2 Experimental results of UCI dataset

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	Zoo	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	7 <sup>10</sup>	5 <sup>2</sup> 7 <sup>8</sup>	2 <sup>10</sup>	2 <sup>8</sup> 3 <sup>1</sup> 4 <sup>1</sup>	3 <sup>10</sup>	2 <sup>3</sup> 7 <sup>7</sup>	7 <sup>10</sup>	7 <sup>10</sup>	7 <sup>10</sup>
	Hayes-Roth	2 <sup>10</sup>	9 <sup>3</sup> 10 <sup>7</sup>	3 <sup>10</sup>	3 <sup>10</sup>	2 <sup>3</sup> 3 <sup>7</sup>	2 <sup>1</sup> 3 <sup>7</sup> 4 <sup>2</sup>	2 <sup>1</sup> 3 <sup>5</sup> 4 <sup>4</sup>	3 <sup>2</sup> 4 <sup>8</sup>	3 <sup>10</sup>	3 <sup>10</sup>	2 <sup>3</sup> 3 <sup>6</sup> 4 <sup>1</sup>	3 <sup>10</sup>
	Iris	2 <sup>10</sup>	3 <sup>9</sup> 5 <sup>1</sup>	3 <sup>10</sup>	3 <sup>10</sup>	3 <sup>8</sup> 5 <sup>2</sup>	3 <sup>10</sup>	3 <sup>6</sup> 4 <sup>1</sup> 5 <sup>3</sup>	10 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>	2 <sup>1</sup> 3 <sup>9</sup>	3 <sup>10</sup>
	Glass	2 <sup>10</sup>	2 <sup>6</sup> 3 <sup>4</sup>	3 <sup>10</sup>	2 <sup>10</sup>	2 <sup>2</sup> 5 <sup>2</sup> 6 <sup>6</sup>	5 <sup>3</sup> 6 <sup>7</sup>	2 <sup>10</sup>	3 <sup>8</sup> 5 <sup>1</sup> 6 <sup>1</sup>	4 <sup>7</sup> 5 <sup>1</sup> 6 <sup>2</sup>	5 <sup>9</sup> 6 <sup>1</sup>	5 <sup>2</sup> 6 <sup>8</sup>	6 <sup>10</sup>
	Dermatology	2 <sup>8</sup> 3 <sup>2</sup>	2 <sup>2</sup> 3 <sup>2</sup> 5 <sup>6</sup>	2 <sup>10</sup>	4 <sup>7</sup> 5 <sup>2</sup> 6 <sup>1</sup>	6 <sup>7</sup> 7 <sup>3</sup>	6 <sup>7</sup> 7 <sup>3</sup>	2 <sup>10</sup>	2 <sup>2</sup> 3 <sup>8</sup>	4 <sup>1</sup> 5 <sup>5</sup> 6 <sup>4</sup>	2 <sup>7</sup> 6 <sup>3</sup>	4 <sup>3</sup> 6 <sup>7</sup>	4 <sup>2</sup> 6 <sup>8</sup>
	Breast Cancer	2 <sup>10</sup>	2 <sup>8</sup> 6 <sup>2</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>8</sup> 3 <sup>2</sup>	2 <sup>8</sup> 3 <sup>2</sup>	2 <sup>10</sup>	2 <sup>8</sup> 3 <sup>1</sup> 5 <sup>1</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>
	Balance Scale	2 <sup>8</sup> 3 <sup>1</sup> 5 <sup>1</sup>	3 <sup>2</sup> 4 <sup>8</sup>	2 <sup>7</sup> 3 <sup>3</sup>	2 <sup>8</sup> 5 <sup>2</sup>	2 <sup>3</sup> 3 <sup>2</sup> 4 <sup>5</sup>	2 <sup>1</sup> 3 <sup>3</sup> 4 <sup>6</sup>	3 <sup>8</sup> 6 <sup>2</sup>	2 <sup>5</sup> 3 <sup>2</sup> 5 <sup>3</sup>	3 <sup>1</sup> 5 <sup>7</sup> 6 <sup>2</sup>	2 <sup>7</sup> 3 <sup>2</sup> 5 <sup>1</sup>	2 <sup>2</sup> 4 <sup>7</sup> 6 <sup>1</sup>	4 <sup>1</sup> 5 <sup>9</sup>
	Libras	2 <sup>22</sup> 3 <sup>3</sup> 5 <sup>4</sup> 7 <sup>1</sup>	2 <sup>15</sup> 3 <sup>12</sup> 4 <sup>3</sup>	2 <sup>15</sup> 3 <sup>5</sup> 4 <sup>10</sup>	18 <sup>22</sup> 20 <sup>8</sup>	14 <sup>4</sup> 15 <sup>20</sup> 16 <sup>6</sup>	13 <sup>9</sup> 14 <sup>15</sup> 15 <sup>6</sup>	2 <sup>25</sup> 3 <sup>3</sup> 6 <sup>2</sup>	30 <sup>30</sup>	14 <sup>7</sup> 15 <sup>17</sup> 17 <sup>6</sup>	2 <sup>16</sup> 10 <sup>6</sup> 11 <sup>4</sup> 15 <sup>4</sup>	10 <sup>3</sup> 13 <sup>2</sup> 14 <sup>20</sup> 15 <sup>5</sup>	11 <sup>4</sup> 14 <sup>2</sup> 15 <sup>22</sup> 16 <sup>2</sup>
	Letter	2 <sup>20</sup> 3 <sup>6</sup> 4 <sup>2</sup> 5 <sup>2</sup>	2 <sup>6</sup> 3 <sup>17</sup> 4 <sup>7</sup>	2 <sup>30</sup>	2 <sup>10</sup> 3 <sup>18</sup> 5 <sup>2</sup>	26 <sup>10</sup> 27 <sup>18</sup> 30 <sup>2</sup>	20 <sup>5</sup> 22 <sup>15</sup> 23 <sup>4</sup> 25 <sup>6</sup>	2 <sup>19</sup> 5 <sup>3</sup> 10 <sup>3</sup> 14 <sup>2</sup> 15 <sup>3</sup>	28 <sup>25</sup> 29 <sup>1</sup> 30 <sup>4</sup>	24 <sup>2</sup> 25 <sup>18</sup> 26 <sup>4</sup> 27 <sup>6</sup>	2 <sup>30</sup>	22 <sup>1</sup> 26 <sup>25</sup> 27 <sup>2</sup> 30 <sup>2</sup>	22 <sup>2</sup> 26 <sup>23</sup> 27 <sup>5</sup>
PFCM	Zoo	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	7 <sup>10</sup>	5 <sup>1</sup> 6 <sup>1</sup> 7 <sup>8</sup>	2 <sup>10</sup>	2 <sup>8</sup> 4 <sup>2</sup>	3 <sup>10</sup>	2 <sup>3</sup> 7 <sup>7</sup>	4 <sup>2</sup> 7 <sup>8</sup>	7 <sup>10</sup>	7 <sup>10</sup>
	Hayes-Roth	3 <sup>10</sup>	2 <sup>1</sup> 9 <sup>2</sup> 10 <sup>7</sup>	3 <sup>10</sup>	3 <sup>10</sup>	2 <sup>6</sup> 3 <sup>3</sup> 5 <sup>1</sup>	3 <sup>8</sup> 4 <sup>2</sup>	3 <sup>7</sup> 4 <sup>3</sup>	3 <sup>5</sup> 4 <sup>3</sup> 7 <sup>2</sup>	3 <sup>4</sup> 4 <sup>1</sup>	3 <sup>10</sup>	2 <sup>3</sup> 3 <sup>6</sup> 4 <sup>2</sup>	3 <sup>9</sup> 4 <sup>1</sup>
	Iris	2 <sup>10</sup>	3 <sup>9</sup> 5 <sup>1</sup>	3 <sup>10</sup>	2 <sup>2</sup> 3 <sup>8</sup>	3 <sup>8</sup> 4 <sup>1</sup> 5 <sup>1</sup>	3 <sup>10</sup>	3 <sup>6</sup> 4 <sup>2</sup> 5 <sup>2</sup>	10 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>
	Glass	2 <sup>10</sup>	2 <sup>7</sup> 3 <sup>3</sup>	2 <sup>1</sup> 3 <sup>8</sup> 5 <sup>1</sup>	2 <sup>10</sup>	5 <sup>3</sup> 6 <sup>7</sup>	4 <sup>1</sup> 5 <sup>1</sup> 6 <sup>8</sup>	2 <sup>10</sup>	3 <sup>10</sup>	4 <sup>3</sup> 6 <sup>7</sup>	4 <sup>1</sup> 5 <sup>8</sup> 6 <sup>1</sup>	5 <sup>2</sup> 6 <sup>8</sup>	6 <sup>10</sup>
	Dermatology	2 <sup>8</sup> 3 <sup>2</sup>	2 <sup>2</sup> 3 <sup>1</sup> 5 <sup>7</sup>	2 <sup>8</sup> 5 <sup>2</sup>	4 <sup>4</sup> 5 <sup>1</sup> 6 <sup>5</sup>	4 <sup>1</sup> 6 <sup>7</sup> 7 <sup>2</sup>	6 <sup>7</sup> 7 <sup>3</sup>	2 <sup>10</sup>	2 <sup>2</sup> 3 <sup>8</sup>	5 <sup>7</sup> 6 <sup>3</sup>	2 <sup>7</sup> 6 <sup>3</sup>	4 <sup>2</sup> 5 <sup>1</sup> 6 <sup>7</sup>	4 <sup>1</sup> 6 <sup>9</sup>
	Breast Cancer	2 <sup>10</sup>	2 <sup>8</sup> 6 <sup>2</sup>	2 <sup>10</sup>	2 <sup>9</sup> 4 <sup>1</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>8</sup> 3 <sup>2</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>
	Balance Scale	2 <sup>8</sup> 3 <sup>2</sup>	2 <sup>1</sup> 3 <sup>2</sup> 4 <sup>7</sup>	2 <sup>7</sup> 3 <sup>3</sup>	2 <sup>8</sup> 4 <sup>1</sup> 5 <sup>1</sup>	2 <sup>1</sup> 3 <sup>3</sup> 4 <sup>6</sup>	3 <sup>3</sup> 4 <sup>7</sup>	3 <sup>8</sup> 6 <sup>2</sup>	2 <sup>7</sup> 5 <sup>3</sup>	5 <sup>8</sup> 6 <sup>2</sup>	2 <sup>7</sup> 3 <sup>2</sup> 5 <sup>1</sup>	2 <sup>2</sup> 4 <sup>8</sup>	5 <sup>10</sup>
	Libras	2 <sup>22</sup> 3 <sup>3</sup> 5 <sup>4</sup>	2 <sup>15</sup> 3 <sup>10</sup> 4 <sup>3</sup> 6 <sup>2</sup>	2 <sup>15</sup> 3 <sup>5</sup> 4 <sup>9</sup> 7 <sup>1</sup>	18 <sup>21</sup> 20 <sup>6</sup> 24 <sup>3</sup>	14 <sup>2</sup> 15 <sup>21</sup> 16 <sup>6</sup> 20 <sup>1</sup>	13 <sup>6</sup> 14 <sup>15</sup> 15 <sup>6</sup> 17 <sup>3</sup>	2 <sup>24</sup> 3 <sup>3</sup> 6 <sup>2</sup> 10 <sup>1</sup>	30 <sup>30</sup>	14 <sup>4</sup> 15 <sup>17</sup> 17 <sup>2</sup> 20 <sup>3</sup>	2 <sup>16</sup> 10 <sup>10</sup> 15 <sup>4</sup>	10 <sup>3</sup> 14 <sup>23</sup> 15 <sup>4</sup>	11 <sup>4</sup> 15 <sup>25</sup> 16 <sup>1</sup>
	Letter	2 <sup>22</sup> 3 <sup>6</sup> 4 <sup>2</sup>	2 <sup>6</sup> 3 <sup>17</sup> 4 <sup>8</sup> 8 <sup>2</sup>	2 <sup>30</sup>	2 <sup>10</sup> 3 <sup>18</sup> 5 <sup>1</sup> 6 <sup>1</sup>	26 <sup>10</sup> 27 <sup>20</sup>	22 <sup>15</sup> 23 <sup>9</sup> 25 <sup>6</sup>	2 <sup>23</sup> 10 <sup>3</sup> 14 <sup>2</sup> 15 <sup>2</sup>	28 <sup>25</sup> 29 <sup>2</sup> 30 <sup>3</sup>	25 <sup>18</sup> 26 <sup>6</sup> 27 <sup>6</sup>	2 <sup>30</sup>	26 <sup>25</sup> 27 <sup>3</sup> 30 <sup>2</sup>	22 <sup>2</sup> 26 <sup>20</sup> 27 <sup>5</sup> 28 <sup>3</sup>

续表 2

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
KFCM	Zoo	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>179</sup>	4 <sup>15178</sup>	2 <sup>10</sup>	2 <sup>941</sup>	3 <sup>753</sup>	2 <sup>24177</sup>	5 <sup>179</sup>	7 <sup>10</sup>	7 <sup>10</sup>
	Hayes-Roth	3 <sup>10</sup>	2 <sup>1109</sup>	3 <sup>10</sup>	2 <sup>238</sup>	2 <sup>10</sup>	3 <sup>10</sup>	2 <sup>13841</sup>	2 <sup>13743</sup>	3 <sup>10</sup>	3 <sup>763</sup>	3 <sup>64272</sup>	3 <sup>10</sup>
	Iris	2 <sup>10</sup>	3 <sup>941</sup>	3 <sup>10</sup>	2 <sup>238</sup>	2 <sup>139</sup>	3 <sup>852</sup>	3 <sup>75281</sup>	10 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>	3 <sup>10</sup>
	Glass	2 <sup>10</sup>	2 <sup>73261</sup>	3 <sup>10</sup>	2 <sup>10</sup>	5 <sup>268</sup>	2 <sup>15168</sup>	2 <sup>10</sup>	3 <sup>10</sup>	2 <sup>14663</sup>	2 <sup>25761</sup>	3 <sup>15267</sup>	4 <sup>268</sup>
	Dermatology	2 <sup>832</sup>	4 <sup>258</sup>	2 <sup>10</sup>	4 <sup>466</sup>	2 <sup>16772</sup>	6 <sup>872</sup>	2 <sup>10</sup>	3 <sup>852</sup>	3 <sup>25761</sup>	2 <sup>75261</sup>	4 <sup>268</sup>	6 <sup>971</sup>
	Breast Cancer	2 <sup>10</sup>	2 <sup>83161</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>10</sup>	2 <sup>832</sup>	2 <sup>10</sup>	2 <sup>10</sup>
	Balance Scale	2 <sup>832</sup>	2 <sup>149</sup>	2 <sup>53342</sup>	2 <sup>53342</sup>	2 <sup>13445</sup>	2 <sup>23345</sup>	3 <sup>76271</sup>	2 <sup>10</sup>	5 <sup>862</sup>	2 <sup>75162</sup>	2 <sup>14653</sup>	5 <sup>961</sup>
	Libras	2 <sup>2232</sup>	2 <sup>1539</sup>	2 <sup>1533</sup>	15 <sup>31821</sup>	14 <sup>41522</sup>	13 <sup>71417</sup>	2 <sup>2134</sup>	30 <sup>30</sup>	15 <sup>17176</sup>	2 <sup>1631</sup>	10 <sup>3132</sup>	14 <sup>21522</sup>
	Letter		5 <sup>4102</sup>	4 <sup>472</sup>	4 <sup>1052</sup>	20 <sup>241</sup>	16 <sup>4</sup>	15 <sup>6</sup>	6 <sup>5</sup>	20 <sup>7</sup>	10 <sup>9154</sup>	14 <sup>21154</sup>	16 <sup>6</sup>
			2 <sup>1934</sup>	2 <sup>6319</sup>	2 <sup>30</sup>	2 <sup>10318</sup>	24 <sup>526</sup>	20 <sup>42215</sup>	2 <sup>2142103</sup>	28 <sup>24294</sup>	25 <sup>18266</sup>	2 <sup>30</sup>	20 <sup>4222</sup>
		4 <sup>25273</sup>	4 <sup>5</sup>	5 <sup>2</sup>	10 <sup>2715</sup>	23 <sup>5256</sup>	14 <sup>2152</sup>	30 <sup>2</sup>	27 <sup>3303</sup>	26 <sup>21273</sup>	27 <sup>3291</sup>		

统计表 2 中的数据, 选取实验结果中出现次数最多的聚类数目作为最终的指标评价结果, 形成表 3。例如 FCM 算法下, Dermatology 数据

集 Dunn 指标的实验结果是“4<sup>75261</sup>”, 比较不同聚类数目出现的次数, 得出最终的聚类数目为 4 类。

表 3 UCI 数据集的 Eff 统计结果  
Table 3 Eff statistical results of UCI dataset

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	Zoo	2	2	2	<u>7</u>	<u>7</u>	2	2	3	7	7	7	7
	Hayes-Roth	2	10	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	4	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Iris	2	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	10	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Glass	2	2	3	2	<b>6</b>	<b>6</b>	2	3	4	5	<b>6</b>	<b>6</b>
	Dermatology	2	5	2	4	<b>6</b>	<b>6</b>	2	3	5	2	<b>6</b>	<b>6</b>
	Breast Cancer	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	Balance Scale	2	4	2	2	4	4	<b>3</b>	2	5	2	4	5
	Libras	2	2	2	18	<b>15</b>	14	2	30	<b>15</b>	2	14	<b>15</b>
	Letter	2	3	2	3	27	22	2	28	25	2	<b>26</b>	<b>26</b>
	Eff	0.11	0.22	0.33	0.44	0.78	0.56	0.44	0.11	0.56	0.44	0.78	0.89
PFCM	Zoo	2	2	2	<b>7</b>	<b>7</b>	2	2	3	7	7	7	7
	Hayes-Roth	<b>3</b>	10	<b>3</b>	<b>3</b>	2	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Iris	2	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	10	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Glass	2	2	3	2	<b>6</b>	<b>6</b>	2	3	<b>6</b>	5	<b>6</b>	<b>6</b>
	Dermatology	2	5	2	<b>6</b>	<b>6</b>	<b>6</b>	2	3	5	2	<b>6</b>	<b>6</b>
	Breast Cancer	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	Balance Scale	2	4	2	2	4	4	<b>3</b>	2	5	2	4	5
	Libras	2	2	2	18	<b>15</b>	14	2	30	<b>15</b>	2	14	<b>15</b>
	Letter	2	3	2	3	27	22	2	28	25	2	<b>26</b>	<b>26</b>
	Eff	0.22	0.22	0.33	0.44	0.67	0.56	0.44	0.22	0.67	0.44	0.78	0.89
KFCM	Zoo	2	2	2	<b>7</b>	<b>7</b>	2	2	3	7	7	7	7
	Hayes-Roth	<b>3</b>	10	<b>3</b>	<b>3</b>	2	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Iris	2	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	10	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Glass	2	2	3	2	<b>6</b>	<b>6</b>	2	3	4	5	<b>6</b>	<b>6</b>
	Dermatology	2	5	2	<b>6</b>	<b>6</b>	<b>6</b>	2	3	5	2	<b>6</b>	<b>6</b>
	Breast Cancer	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	Balance Scale	2	4	2	2	4	4	<b>3</b>	2	5	2	4	5
	Libras	2	2	2	18	<b>15</b>	14	2	30	<b>15</b>	2	14	<b>15</b>
	Letter	2	3	2	3	27	22	2	28	25	2	<b>26</b>	<b>26</b>
	Eff	0.22	0.22	0.33	0.56	0.67	0.56	0.44	0.22	0.56	0.44	0.78	0.89

注: 加粗代表该数值和正确的聚类数目一致。

此时需要借助评价因子来对 CVI 的性能进行评价。选取的评价因子为有效性比率 (effectiveness ratio, Eff)<sup>[37]</sup>, 其值为 CVI 得到正确聚类数目结果的数据集个数与总的数据集个数的比值。Eff 的计算公式为

$$L_{\text{Eff}} = \frac{1}{D} \sum_{d=1}^D \theta, \theta = \begin{cases} 1, & k \text{ 是正确聚类数} \\ 0, & \text{其他} \end{cases}$$

分析表 3 的统计数据可以看出, 本文提出的 GDD 指标在 UCI 数据集上的实验结果中相较于其他 CVI 准确率更高。GDD 指标在 FCM 算法、PFCM 算法和 KFCM 算法上均是有 1 个数据集的结果错误, 准确率达到 89%。其他准确率较高的指标, 例如 IMI 指标和 TCR 指标, 其在 FCM 算法上的准确率大致在 78% 左右, 在 PFCM 算法和 KFCM 算法上的准确率稍低一点在 70% 左右。剩余的指标准准确率大多都在 60% 以下。从实验

结果可以分析出, 本文提出的 GDD 指标在 3 个算法上的准确率均高于其他 CVI, 说明 GDD 指标对于不同算法的适应力均优于其他 CVI。同时对比不同 CVI 在不同算法下的准确率, 例如 WLI 指标在 FCM 算法上的准确率为 56%, 而在 PFCM 算法下的准确率为 67%, 这在一定程度上说明不同的 CVI 有着更契合 CVI 本身的算法。

为了更精细地确定一个 CVI 的评价效果和稳定程度, 在此引入另外一个评价因子来对实验结果进行评价。引入的评价因子为有效性偏差 (Bias), 其值为实验所得的聚类数目和正确聚类数目之间差值的绝对值。Bias 的计算公式为

$$L_{\text{Bias}} = \sum_{d=1}^D |k - k_{\text{true}}|$$

应用 Bias 评价因子后得到表 4。

表 4 UCI 数据集的 Bias 统计结果  
Table 4 Bias statistical results of UCI dataset

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	Zoo	5	5	5	0	0	5	5	4	0	0	0	0
	Hayes-Roth	1	7	0	0	0	0	0	1	0	0	0	0
	Iris	1	0	0	0	0	0	0	7	0	0	0	0
	Glass	4	4	3	4	0	0	4	3	2	1	0	0
	Dermatology	4	1	4	2	0	0	4	3	1	4	0	0
	Breast Cancer	0	0	0	0	0	0	0	0	0	0	0	0
	Balance Scale	1	1	1	1	1	1	0	1	2	1	1	2
	Libras	13	13	13	3	0	1	13	15	0	13	1	0
	Letter	24	23	24	23	1	4	24	2	1	24	0	0
	Bias	53	54	50	33	2	11	50	36	6	43	2	2
PFCM	Zoo	5	5	5	0	0	5	5	4	0	0	0	0
	Hayes-Roth	0	7	0	0	1	0	0	0	0	0	0	0
	Iris	1	0	0	0	0	0	0	7	0	0	0	0
	Glass	4	4	3	4	0	0	4	3	0	1	0	0
	Dermatology	4	1	4	0	0	0	4	3	1	4	0	0
	Breast Cancer	0	0	0	0	0	0	0	0	0	0	0	0
	Balance Scale	1	1	1	1	1	1	0	1	2	1	1	2
	Libras	13	13	13	3	0	1	13	15	0	13	1	0
	Letter	24	23	24	23	1	4	24	2	1	24	0	0
	Bias	52	54	50	31	3	11	50	35	4	43	2	2
KFCM	Zoo	5	5	5	0	0	5	5	4	0	0	0	0
	Hayes-Roth	0	7	0	0	1	0	0	0	0	0	0	0
	Iris	1	0	0	0	0	0	0	7	0	0	0	0
	Glass	4	4	3	4	0	0	4	3	2	1	0	0
	Dermatology	4	1	4	0	0	0	4	3	1	4	0	0
	Breast Cancer	0	0	0	0	0	0	0	0	0	0	0	0
	Balance Scale	1	1	1	1	1	1	0	1	2	1	1	2
	Libras	13	13	13	3	0	1	13	15	0	13	1	0
	Letter	24	23	24	23	1	4	24	2	1	24	0	0
	Bias	52	54	50	31	3	11	50	35	6	43	2	2

分析表 4 中的数据,以 KFCM 算法下的运行结果为例, Bias 值比较低的有 IMI 指标、WLI 指标、TCR 指标和 GDD 指标,其中 GDD 指标的 Bias 值最低,足以说明本文提出的 GDD 指标在 UCI 数据集上有更好的评价效果和更稳定的评价能力,即使在 GDD 指标给出错误评价结果的情况下,得到的错误聚类数目和正确的聚类数目也不会相差太多。相反,其他 CVI,如 CH 指标、Dunn 指标和 DB 指标,这些 CVI 的 Bias 值都在 50 上下,相较于其他 CVI 高出很多。这说明 CH 指标、Dunn 指标和 DB 指标这些 CVI 在错误评价后得到的聚类数目会和正确的聚类数目相差较大,其稳定性较差。在其他几个算法下, GDD 的 Bias 值

也是最低的。

同时,为了更直观地体现不同 CVI 在不同数据集下的评价结果和正确值之间的差距波动,选取 FCM 算法下实验结果绘制图 3,其中红色点画线是当前数据集的正确聚类数。分析图 3 中的折线图可以发现,本文提出的 GDD 指标是 12 个 CVI 中比较稳定的指标,在绝大多数数据集上评价结果都是正确的,在评价错误的数据集 Balance Scale 上, GDD 指标的结果也和正确的聚类数目相差不多。相较于其他 CVI,例如 CH 指标、Dunn 指标和 FSI 指标,这些 CVI 在大部分数据集上评价结果都和正确的聚类数目相差较大。

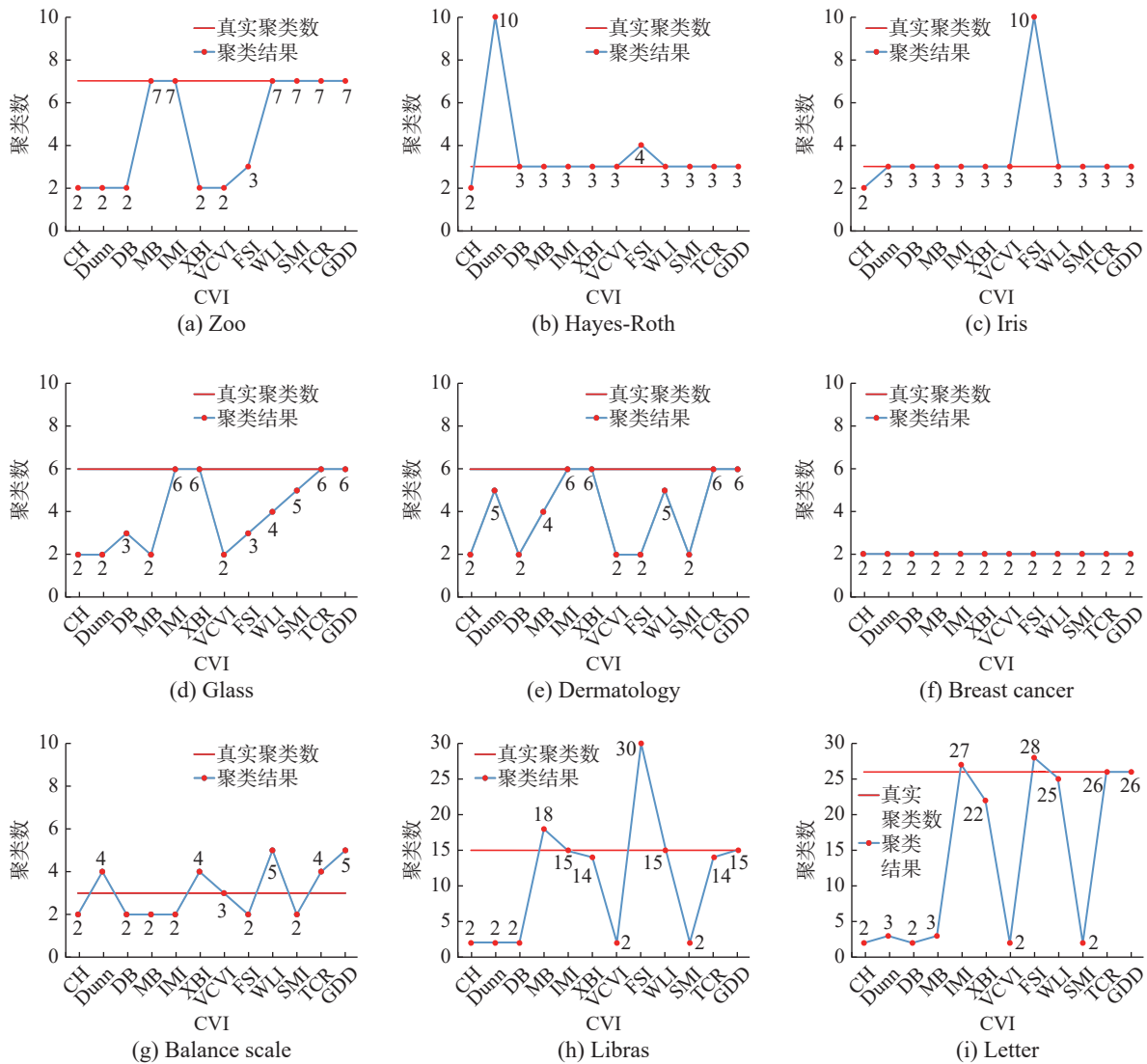


图 3 FCM 算法下 UCI 数据集的统计结果

Fig. 3 Statistical results of UCI dataset under the FCM algorithm

图 4 给出了 12 个指标在 FCM 算法下的 Dermatology 数据集的计算结果数值变化,其中红色“·”形状标记表示当前聚类数目下的 CVI 取最优值。这里, FSI 本应为值越小越好,在图 4 中将其

纵坐标翻转后标记为最大值为最优聚类数。通过分析图 3、图 4,可以发现 IMI 指标、XBI 指标、TCR 指标和 GDD 指标在当前数据集上的聚类结果是正确的,其中正确的聚类数目为 6 类。然而分析





续表 6

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	Circle	3	3	<b>2</b>	3	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	Jain	3	3	3	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	4	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	X8D5K	2	<b>5</b>	2	3	<b>5</b>	3	4	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	Data_77	2	2	3	6	5	5	5	5	5	5	5	7
	E6	2	2	5	2	5	4	4	2	4	4	4	4
	Dim_128	2	2	2	2	<b>16</b>	14	2	30	14	2	<b>16</b>	<b>16</b>
	Dim_256	2	2	2	3	<b>16</b>	15	3	30	15	2	15	<b>16</b>
	Eff	0.22	0.33	0.33	0.33	0.78	0.56	0.56	0.33	0.67	0.67	0.78	1
PFCM	Data_60	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Data_150	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Circle	3	<b>2</b>	<b>2</b>	3	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	Jain	<b>2</b>	3	3	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	4	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	X8D5K	2	<b>5</b>	2	3	<b>5</b>	3	4	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	Data_77	2	2	3	6	5	5	5	5	5	5	5	7
	E6	2	2	5	2	5	4	4	2	4	4	4	4
	Dim_128	2	2	<sup>2</sup>	2	<b>16</b>	<b>16</b>	2	30	<b>16</b>	2	<b>16</b>	<b>16</b>
Dim_256	2	2	2	2	<b>16</b>	15	3	30	14	2	15	<b>16</b>	
Eff	0.33	0.44	0.33	0.33	0.78	0.67	0.56	0.33	0.78	0.67	0.78	1	
KFCM	Data_60	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Data_150	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>
	Circle	3	3	<b>2</b>	3	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
	Jain	3	3	3	<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>	4	<b>2</b>	<b>2</b>	<b>2</b>	3
	X8D5K	2	<b>5</b>	2	3	4	3	4	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	Data_77	2	2	3	6	5	5	5	5	5	5	5	7
	E6	2	2	5	2	5	4	4	2	4	4	4	4
	Dim_128	2	2	2	2	<b>16</b>	14	3	30	14	2	<b>16</b>	<b>16</b>
Dim_256	2	2	2	3	<b>16</b>	15	3	30	14	2	15	<b>16</b>	
Eff	0.22	0.33	0.33	0.33	0.67	0.56	0.56	0.33	0.67	0.67	0.78	0.89	

注: 加粗代表该数值和正确的聚类数目一致。

表 7 人造数据集的 Bias 统计结果

Table 7 Bias statistical results of artificial dataset

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	Data_60	0	0	0	0	0	0	0	0	0	0	0	0
	Data_150	0	0	0	0	0	0	0	0	0	0	0	0
	Circle	1	1	0	1	0	0	0	0	0	0	0	0
	Jain	1	1	1	0	0	0	0	2	0	0	0	0
	X8D5K	3	0	3	2	0	2	1	1	0	0	0	0
	Data_77	5	5	4	1	2	2	2	2	2	2	2	0
	E6	2	2	1	2	1	0	0	2	0	0	0	0
	Dim_128	14	14	14	14	0	2	14	14	2	14	0	0
	Dim_256	14	14	14	13	0	1	13	14	1	14	1	0
	Bias	40	37	37	33	3	7	30	35	5	30	3	0

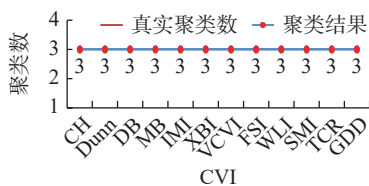
续表 7

算法	数据集	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
PFCM	Data_60	0	0	0	0	0	0	0	0	0	0	0	0
	Data_150	0	0	0	0	0	0	0	0	0	0	0	0
	Circle	1	0	0	1	0	0	0	0	0	0	0	0
	Jain	0	1	1	0	0	0	0	2	0	0	0	0
	X8D5K	3	0	3	2	0	2	1	1	0	0	0	0
	Data_77	5	5	4	1	2	2	2	2	2	2	2	0
	E6	2	2	1	2	1	0	0	2	0	0	0	0
	Dim_128	14	14	14	14	0	0	14	14	0	14	0	0
	Dim_256	14	14	14	14	0	1	13	14	2	14	1	0
	Bias	39	36	37	34	3	5	30	35	4	30	3	0
KFCM	Data_60	0	0	0	0	0	0	0	0	0	0	0	0
	Data_150	0	0	0	0	0	0	0	0	0	0	0	0
	Circle	1	1	0	1	0	0	0	0	0	0	0	0
	Jain	1	1	1	0	0	0	0	2	0	0	0	1
	X8D5K	3	0	3	2	1	2	1	1	0	0	0	0
	Data_77	5	5	4	1	2	2	2	2	2	2	2	0
	E6	2	2	1	2	1	0	0	2	0	0	0	0
	Dim_128	14	14	14	14	0	2	13	14	2	14	0	0
	Dim_256	14	14	14	13	0	1	13	14	2	14	1	0
	Bias	40	37	37	33	4	7	29	35	6	30	3	1

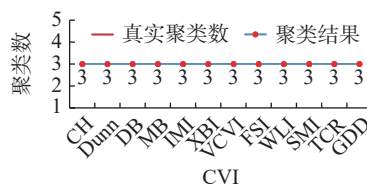
分析表 7 中的数据, 以 FCM 算法下的运行结果为例, Bias 值比较低的有 IMI 指标、WLI 指标、TCR 指标和 GDD 指标, 其中 GDD 指标的 Bias 值最低, 足以说明本文提出的 GDD 指标在人造数据集上有更好的评价效果和更稳定的评价能力, 即使在 GDD 指标给出错误的评价结果的情况下, 得到的错误聚类数目和正确的聚类数目也不会相差太多。相反, 其他 CVI 例如 CH 指标、Dunn 指标和 DB 指标, 这些 CVI 的 Bias 值都在 40 上下, 相较于其他 CVI 高出很多。这说明 CH 指标、Dunn 指标和 DB 指标这些 CVI 在错误评价后得到的聚类数目会和正确的聚类数目相差较大, 其

稳定性较差。

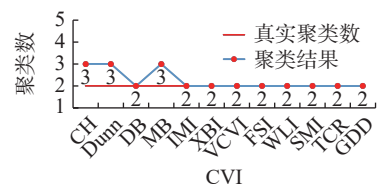
同时, 为了更直观地体现不同 CVI 在不同数据集下的评价结果和正确值之间的差距波动, 选取 FCM 算法下实验结果绘制图 5, 其中红色点画线是当前数据集的正确聚类数。分析图 5 中的折线图, 可以发现, 本文提出的 GDD 指标是 12 个 CVI 中比较稳定的指标, 在绝大多数数据集上评价结果都是正确的, 在 KFCM 算法下评价错误的数据集 Jain 上, GDD 指标的结果也和正确的聚类数目相差不多。如 CH 指标、Dunn 指标和 FSI 指标, 这些 CVI 在大部分数据集上评价结果都和正确的聚类数目相差较大。



(a) Data\_60



(b) Data\_150



(c) Circle

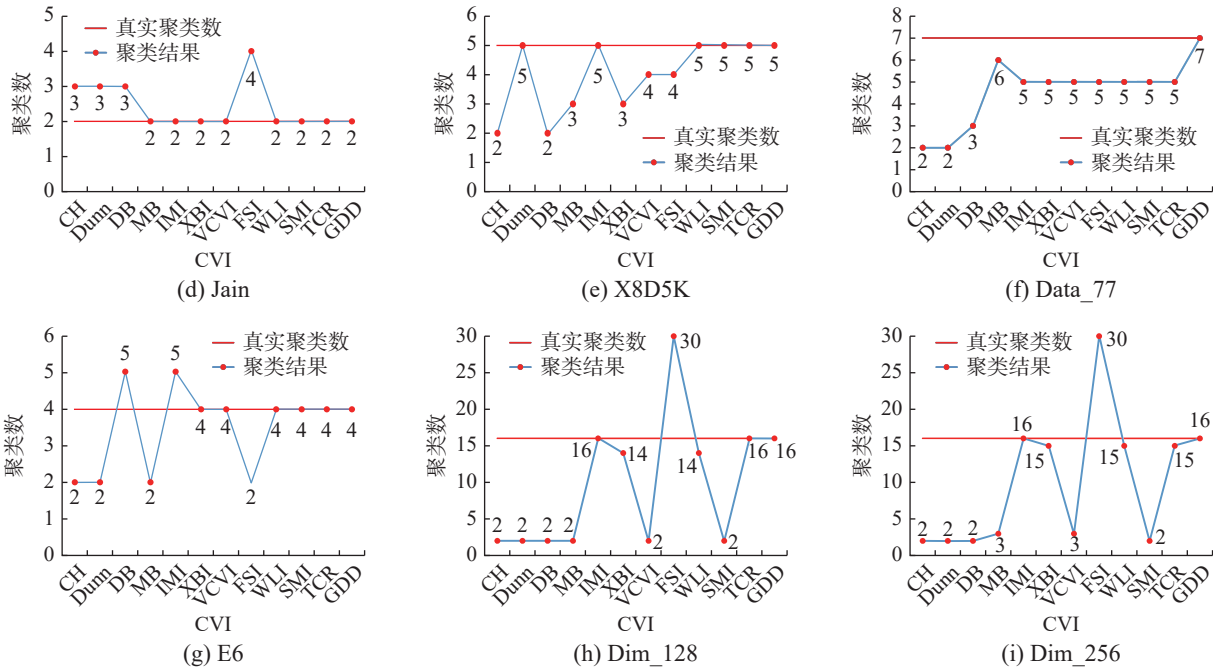


图 5 FCM 算法下人造数据集的统计结果

Fig. 5 Statistical results of artificial dataset under FCM algorithm

接下来使用 Olivetti Face 数据集进行实验来检验本文提出的 GDD 指标面对图像类数据集的效果, 验证其能否适应不同类型的数据。表 8 是 Olivetti Face 数据集在 FCM、PFCM 和 KFCM 这

3 种算法下的运行结果。统计表 8 中的数据, 选取实验结果中出现次数最多的聚类数目作为最终的指标评价结果, 同时应用 Eff 和 Bias 评价因子以形成表 9。

表 8 Olivetti Face 数据集的实验结果

Table 8 Experimental results of Olivetti Face dataset

算法	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	30 <sup>30</sup>	2 <sup>17</sup> 3 <sup>5</sup>	2 <sup>30</sup>	2 <sup>23</sup> 3 <sup>2</sup>	9 <sup>5</sup> 10 <sup>23</sup>	25 <sup>6</sup> 29 <sup>5</sup>	2 <sup>1</sup> 3 <sup>4</sup> 4 <sup>25</sup>	30 <sup>30</sup>	5 <sup>2</sup> 8 <sup>3</sup>	2 <sup>30</sup>	10 <sup>20</sup> 11 <sup>10</sup>	8 <sup>5</sup> 9 <sup>3</sup>
		4 <sup>5</sup> 6 <sup>3</sup>		4 <sup>1</sup> 6 <sup>2</sup> 9 <sup>2</sup>	12 <sup>2</sup>	30 <sup>19</sup>			9 <sup>5</sup> 10 <sup>20</sup>			10 <sup>22</sup>
PFCM	30 <sup>30</sup>	2 <sup>21</sup> 3 <sup>1</sup>	2 <sup>30</sup>	2 <sup>21</sup> 5 <sup>3</sup>	2 <sup>1</sup> 8 <sup>2</sup>	22 <sup>3</sup> 25 <sup>4</sup>	3 <sup>5</sup> 4 <sup>25</sup>	30 <sup>30</sup>	7 <sup>3</sup> 9 <sup>1</sup>	2 <sup>30</sup>	8 <sup>2</sup> 10 <sup>19</sup>	5 <sup>1</sup> 9 <sup>4</sup>
		4 <sup>5</sup> 5 <sup>3</sup>		9 <sup>4</sup> 10 <sup>2</sup>	9 <sup>4</sup> 10 <sup>23</sup>	28 <sup>4</sup> 30 <sup>19</sup>			10 <sup>23</sup> 11 <sup>3</sup>		11 <sup>5</sup> 14 <sup>4</sup>	10 <sup>21</sup> 12 <sup>4</sup>
KFCM	25 <sup>1</sup> 29 <sup>1</sup>	2 <sup>10</sup> 3 <sup>16</sup>	2 <sup>30</sup>	2 <sup>25</sup> 4 <sup>5</sup>	8 <sup>2</sup> 9 <sup>4</sup>	22 <sup>4</sup> 29 <sup>5</sup>	3 <sup>2</sup> 4 <sup>21</sup>	30 <sup>30</sup>	9 <sup>5</sup> 10 <sup>22</sup> 11 <sup>3</sup>	2 <sup>30</sup>	9 <sup>3</sup> 10 <sup>22</sup>	9 <sup>5</sup> 10 <sup>19</sup>
	30 <sup>28</sup>	5 <sup>3</sup> 7 <sup>1</sup>			10 <sup>22</sup> 20 <sup>2</sup>	30 <sup>21</sup>	5 <sup>5</sup> 10 <sup>2</sup>				11 <sup>3</sup> 14 <sup>2</sup>	11 <sup>2</sup> 13 <sup>4</sup>

表 9 Olivetti Face 数据集的统计结果

Table 9 Statistical results of Olivetti Face dataset

算法	评价因子	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	Value	30	2	2	2	<b>10</b>	30	4	30	<b>10</b>	2	<b>10</b>	<b>10</b>
	Eff	0	0	0	0	1	0	0	0	1	0	1	1
	Bias	20	8	8	8	0	20	6	20	0	8	0	0
PFCM	Value	30	2	2	2	<b>10</b>	30	4	30	<b>10</b>	2	<b>10</b>	<b>10</b>
	Eff	0	0	0	0	1	0	0	0	1	0	1	1
	Bias	20	8	8	8	0	20	6	20	0	8	0	0
KFCM	Value	30	3	2	2	<b>10</b>	30	4	30	<b>10</b>	2	<b>10</b>	<b>10</b>
	Eff	0	0	0	0	1	0	0	0	1	0	1	1
	Bias	20	7	8	8	0	20	6	20	0	8	0	0

注: 加粗代表该数值和正确的聚类数目一致。

评估一个 CVI 的优劣不仅需要各种类型、各种结构的数据集进行实验,同样也需要检测 CVI 在含有噪声的数据集上的效果。一个稳定的 CVI 应该能适应添加有不同程度噪声数据的数据集,并且能得出正确的结果。接下来选取 X8D5K 数据集, X8D5K 数据集的数据维度为

8 维,共有 1 000 个样本点,正确的聚类数目为 5 类。向其中加入不同程度的噪声数据,分别为 2%、4%、6%、8% 和 10%。在这 5 种噪声下的数据集上进行实验,分别使用 FCM 算法、PFCM 算法和 KFCM 算法。统计实验后的结果如表 10 所示。

表 10 噪声数据集的统计结果  
Table 10 Statistical results of noisy dataset

算法	噪声比例/%	CH	Dunn	DB	MB	IMI	XBI	VCVI	FSI	WLI	SMI	TCR	GDD
FCM	0	2	<b>5</b>	2	3	<b>5</b>	3	4	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	2	3	<b>5</b>	2	3	<b>5</b>	4	<b>5</b>	4	<b>5</b>	4	<b>5</b>	<b>5</b>
	4	2	<b>5</b>	2	3	4	4	4	4	<b>5</b>	4	<b>5</b>	<b>5</b>
	6	2	4	3	3	3	3	3	6	4	4	<b>5</b>	<b>5</b>
	8	2	4	3	3	3	3	4	3	6	3	6	<b>5</b>
	10	2	3	3	3	3	2	3	3	<b>5</b>	2	3	4
PFCM	0	2	<b>5</b>	2	3	<b>5</b>	3	4	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	2	2	<b>5</b>	2	3	<b>5</b>	3	4	3	<b>5</b>	4	<b>5</b>	<b>5</b>
	4	3	<b>5</b>	3	3	<b>5</b>	4	3	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	6	3	4	3	3	4	<b>5</b>	3	3	4	7	4	<b>5</b>
	8	2	3	4	3	3	<b>5</b>	3	2	3	4	3	<b>5</b>
	10	2	3	<b>5</b>	3	3	4	2	2	3	4	3	<b>5</b>
KFCM	0	2	<b>5</b>	2	3	4	3	4	4	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	2	2	<b>5</b>	<b>5</b>	3	<b>5</b>	3	4	3	<b>5</b>	<b>5</b>	<b>5</b>	<b>5</b>
	4	2	4	<b>5</b>	3	4	4	4	4	5	4	<b>5</b>	<b>5</b>
	6	2	3	4	3	4	3	3	3	4	<b>5</b>	4	<b>5</b>
	8	2	4	3	3	3	3	3	3	3	4	4	<b>5</b>
	10	2	2	3	4	3	2	4	3	3	3	3	6

注:加粗代表该数值和正确的聚类数目一致。

#### 4 结束语

1) 本文提出了一个新的模糊聚类有效性评价指标 GDD。GDD 指标采用簇内紧致度与簇间分离度比值的形式抑制了聚类结果数目单调递增的情况,同时在簇内紧致度中在计算单个簇紧致度的基础上引入了该簇样本占总样本的比值,构建了更客观的紧致度表达;在簇间分离度中在计算聚类中心到聚类中心均值的基础上增加了该簇样本占总样本的比值,构建了更合理的簇间分离度表达。

2) 在复杂度分析中,虽然 GDD 指标拥有与 Dunn 和 SMI 指标相同的时间复杂度,但是其拥有更加优异的性能,并且性能同样优于其他时间复杂度更低的指标,说明了 GDD 在能接受的时间

复杂度下达到了比较指标中最好的性能。

3) 在面对 3 个模糊聚类算法、11 个对比 CVI 的实验中,GDD 指标在 UCI 数据集中均优于其他指标,在人造数据集的实验中 GDD 仅在 KFCM 算法下的 Jain 数据集出现了错误,在人脸数据集 Olivetti Face 中 GDD 同样得到了较为优秀的效果。

4) 在噪声实验中,FCM 算法下噪声为 8% 时只有 GDD 指标得到了正确结果,而噪声为 10% 时只有 FSI 得到了正确结果。PFCM 算法下当噪声为 10% 时只有 DB 和 GDD 能够得到正确结果。KFCM 算法下当噪声为 8% 时只有 GDD 得到了正确结果,而噪声为 10% 时则没有指标得到正确聚类数。

综上,GDD 指标具有全新的指标设计思路、可接受的时间复杂度与较好的实验表现,在噪声

实验中体现了其较强的鲁棒性。在未来研究中, 将试图结合模糊逻辑的思想来进行聚类有效性评价指标的设计与应用。

## 参考文献:

- [1] TANG Yiming, PAN Zhifu, HU Xianghui, et al. Knowledge-induced multiple kernel fuzzy clustering[J]. *IEEE transactions on pattern analysis and machine intelligence*, 2023, 45(12): 14838–14855.
- [2] TANG Yiming, PAN Zhifu, PEDRYCZ W, et al. Viewpoint-based kernel fuzzy clustering with weight information granules[J]. *IEEE transactions on emerging topics in computational intelligence*, 2023, 7(2): 342–356.
- [3] TANG Yiming, REN Fuji, PEDRYCZ W. Fuzzy C-Means clustering through SSIM and patch for image segmentation[J]. *Applied soft computing*, 2020, 87: 1–16.
- [4] TANG Yiming, WU Wenbin, PEDRYCZ W, et al. Clustering interval and triangular granular data: modeling, execution, and assessment[J]. *IEEE transactions on neural networks and learning systems*, 2025, 36(6): 10000–10014.
- [5] TANG Yiming, LI Bing, PEDRYCZ W, et al. A clustering validity index with multi-granularity fusion for multiple fuzzy clustering algorithms[J]. *IEEE transactions on pattern analysis and machine intelligence*, 2025, 47(10): 8379–8396.
- [6] MACQUEEN J. Some methods for classification and analysis of multivariate observations[C]//Proceedings of Fifth Berkeley Symposium on Mathematical Statistics and Probability. Berkeley: University of California, 1965.
- [7] ROSS H H, SOKAL R R, SNEATH P H A, et al. Principles of numerical taxonomy[J]. *Systematic zoology*, 1964, 13(2): 106.
- [8] CAMPELLO R J G B, MOULAVI D, SANDER J. Density-based clustering based on hierarchical density estimates[C]//The 17th Pacific-Asia Conference on Knowledge Discovery and Data Mining. Gold Coast: PAKDD, 2013.
- [9] 吕莉, 陈威, 肖人彬, 等. 面向密度分布不均数据的加权逆近邻密度峰值聚类算法[J]. *智能系统学报*, 2024, 19(1): 165–175.  
LYU Li, CHEN Wei, XIAO Renbin, et al. Density peak clustering algorithm based on weighted reverse nearest-neighbor for uneven density datasets[J]. *CAAI transactions on intelligent systems*, 2024, 19(1): 165–175.
- [10] ZHAO Yanchang, SONG Junde. GDILC: a grid-based density-isoline clustering algorithm[C]//Proceedings of the International Conference on Information Technology and Intelligent Network. Beijing: IEEE, 2001.
- [11] VIJAY R K, NANDA S J, SHARMA A. A spatio-temporal binary grid-based clustering model for seismicity analysis[J]. *Pattern analysis and applications*, 2024, 27(1): 14.
- [12] 孙林, 梁娜, 徐久成. 基于邻域互信息与 K-means 特征聚类的特征选择[J]. *智能系统学报*, 2024, 19(4): 983–996.  
SUN Lin, LIANG Na, XU Jiucheng. Feature selection using neighborhood mutual information and feature clustering with K-means[J]. *CAAI transactions on intelligent systems*, 2024, 19(4): 983–996.
- [13] DUNN J C. A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters[J]. *Journal of cybernetics*, 1973, 3(3): 32–57.
- [14] BEZDEK J C, EHRLICH R, FULL W. FCM. The fuzzy c-means clustering algorithm[J]. *Computers & geosciences*, 1984, 10(2/3): 191–203.
- [15] KRISHNAPURAM R, KELLER J M. A possibilistic approach to clustering[J]. *IEEE transactions on fuzzy systems*, 1993, 1(2): 98–110.
- [16] PAL N R, PAL K, BEZDEK J C. A mixed C-means clustering model[C]//Proceedings of 6th International Fuzzy Systems Conference. Barcelona: IEEE, 1997.
- [17] ANTOINE V, GUERRERO J A, ROMERO G. Possibilistic fuzzy c-means with partial supervision[J]. *Fuzzy sets and systems*, 2022, 449: 162–186.
- [18] ZHANG Daoqiang, CHEN Songcan, PAN Zhisong, et al. Kernel-based fuzzy clustering incorporating spatial constraints for image segmentation[C]//Proceedings of the 2003 International Conference on Machine Learning and Cybernetics. Washington: IEEE, 2003.
- [19] PUNIT R, ZAHRA G, BEZDEK J C, et al. Approximating Dunn's cluster validity indices for partitions of big data[J]. *IEEE transactions on cybernetics*, 2018, 49(5): 1629–1641.
- [20] HASSAN B A, TAYFOR N B, HASSAN A A, et al. From a-to-z review of clustering validation indices[J]. *Neurocomputing*, 2024, 601: 128–198.
- [21] FUKUI K, NUMAO M. Neighborhood-based smoothing of external cluster validity measures[C]//Proceedings of the 16th Pacific-Asia Conference on Knowledge Discovery and Data Mining. Berlin: Springer, 2012: 354–365.
- [22] 唐益明, 陈仁好, 李冰. 面向模糊 C 均值算法的 MAME 聚类有效性指标[J]. *智能系统学报*, 2023, 18(5): 945–956.  
TANG Yiming, CHEN Renhao, LI Bing. A clustering validity index called MAME for the fuzzy c-means algorithm[J]. *CAAI transactions on intelligent systems*, 2023, 18(5): 945–956.
- [23] VENDRAMIN L, CAMPELLO R J G B, HRUSCHKA E

- R. Relative clustering validity criteria: a comparative overview[J]. *Statistical analysis & data mining*, 2010, 3(4): 209–235.
- [24] CALINSKI R B, HARABASZ J. A dendrite method for cluster analysis[J]. *Communications in statistics*, 1974, 3(1): 1–27.
- [25] DAVIES D L, BOULDIN D W. A cluster separation measure[J]. *IEEE transactions on pattern analysis and machine intelligence*, 1979, 2(1): 224–227.
- [26] XIE X L, BENI G. A validity measure for fuzzy clustering[J]. *IEEE transactions on pattern analysis and machine intelligence*, 1991, 13(8): 841–847.
- [27] FUKUYAMA Y, SUGENO M. A new method of choosing the number of clusters for the fuzzy c-means method[C]// *Proceedings of the 5th Fuzzy Systems Symposium*. Tokyo: Japan Society for Fuzzy Theory and Systems, 1989.
- [28] WU Chihhung, OUYANG Chensen, CHEN Liwen, et al. A new fuzzy clustering validity index with a median factor for centroid-based clustering [J] *IEEE transactions on fuzzy systems*, 2015, 23(3): 701–718.
- [29] LIU Yun, JIANG Yanfang, HOU Tao, et al. A new robust fuzzy clustering validity index for imbalanced data sets[J]. *Information sciences*, 2021, 547: 579–591.
- [30] MITTAL H, SARASWAT M. A new fuzzy cluster validity index for hyperellipsoid or hyperspherical shape close clusters with distant centroids[J]. *IEEE transactions on fuzzy systems*, 2020, 29(11): 3249–3258.
- [31] ZHU Erzhuo, MA Zhujuan, LI Xuejun, et al. An effective partitional clustering algorithm based on new clustering validity index[J]. *Applied soft computing*, 2018, 71: 608–621.
- [32] MAULIKH U, BANDYOPADHYAY S. Performance evaluation of some clustering algorithms and validity indices[J]. *IEEE transactions on pattern analysis and machine intelligence*, 2002, 24(12): 1650–1654.
- [33] TANG Yiming, HUANG Jiajia, PEDRYCZ W, et al. A fuzzy clustering validity index induced by triple center relation[J]. *IEEE transactions on cybernetics*, 2023, 53(8): 5024–5036.
- [34] LIU Yun, HOU Tao, LIU Fu, et al. Improving fuzzy c-means method for unbalanced dataset[J]. *Electronics letters*, 2015, 51(23): 1880–1881.
- [35] ZHOU Kaile, YANG Shanlin. Exploring the uniform effect of FCM clustering: a data distribution perspective[J]. *Knowledge-based systems*, 2016, 96: 76–83.
- [36] DUA D, GRAFF C. UCI machine learning repository [EB/OL]. (2017–01–01)[2025–03–20]. <http://archive.ics.uci.edu/ml>.
- [37] SALEM S A, NANDI A K. Development of assessment criteria for clustering algorithms[J]. *Pattern analysis and applications*, 2009, 12(1): 79–98.

### 作者简介:



唐益明, 教授, 博士, 主要研究方向为聚类、模糊逻辑与推理、情感计算和图像处理。主持国家自然科学基金项目 4 项。发表学术论文 100 余篇, 获国家发明专利授权 8 项。E-mail: [tym608@163.com](mailto:tym608@163.com)。



刘子龙, 硕士研究生, 主要研究方向为聚类和聚类有效性指标。E-mail: [2024170934@mail.hfut.edu.cn](mailto:2024170934@mail.hfut.edu.cn)。



高健玮, 博士研究生, 主要研究方向为聚类、粒计算和模糊推理。E-mail: [jwgao810@163.com](mailto:jwgao810@163.com)。

[ 责任编辑: 刘冰洁 ]