

# 控制与决策

Control and Decision

## 增量学习的模糊风格K平面聚类

顾苏杭, 王士同

引用本文:

顾苏杭, 王士同. 增量学习的模糊风格K平面聚类[J]. 控制与决策, 2020, 35(9): 2081–2093.

在线阅读 View online: <https://doi.org/10.13195/j.kzyjc.2019.0023>

---

## 您可能感兴趣的其他文章

Articles you may be interested in

### 基于密度的模糊代表点聚类算法

A density-based fuzzy exemplar clustering algorithm

控制与决策. 2020, 35(5): 1123–1133 <https://doi.org/10.13195/j.kzyjc.2018.1179>

### 基于目标特征选择和去除的改进K-means聚类算法

Improved K-means clustering algorithm based on feature selection and removal on target point

控制与决策. 2019, 34(6): 1219–1226 <https://doi.org/10.13195/j.kzyjc.2017.1548>

### 一种基于相对密度和决策图的聚类算法

A novel clustering algorithm based on relative density and decision graph

控制与决策. 2018, 33(11): 1921–1930 <https://doi.org/10.13195/j.kzyjc.2017.0822>

### 基于划分自适应融合的多视角模糊聚类算法

Multi-view fuzzy clustering algorithm based on partition adaptive-fusion

控制与决策. 2016, 31(4): 593–600 <https://doi.org/10.13195/j.kzyjc.2015.0057>

### 基于多代表点的大规模数据模糊聚类算法

Fuzzy clustering algorithm based on multiple medoids for large-scale data

控制与决策. 2016, 31(12): 2122–2130 <https://doi.org/10.13195/j.kzyjc.2015.1488>

# 增量学习的模糊风格 $K$ 平面聚类

顾苏杭<sup>1,2†</sup>, 王士同<sup>1</sup>

(1. 江南大学 数字媒体学院, 江苏 无锡 214122;  
2. 常州轻工职业技术学院 信息工程与技术学院, 江苏 常州 213164)

**摘要:** 提出利用特征增量学习和数据风格信息双知识表达约束的模糊  $K$  平面聚类 (ISF-KPC) 算法. 为了获得更好的泛化性, 聚类前利用高斯核函数对原输入特征进行增长式的特征扩维. 考虑数据集中来源于同一聚类的样本具有相同的风格, 以矩阵的形式表达数据风格信息, 并采用迭代的方式确定每个聚类的风格矩阵. 大量实验结果表明, 双知识表达约束的 ISF-KPC 与对比算法相比能够取得竞争性的聚类性能, 尤其在具有典型风格数据集上能够取得优异的聚类性能.

**关键词:**  $K$  平面聚类; 风格信息; 特征扩维; 模糊聚类

中图分类号: TP391 文献标志码: A

DOI: 10.13195/j.kzyjc.2019.0023

引用格式: 顾苏杭, 王士同. 增量学习的模糊风格  $K$  平面聚类[J]. 控制与决策, 2020, 35(9): 2081-2093.

## Incremental learning based fuzzy style $K$ -plane clustering

GU Su-hang<sup>1,2†</sup>, WANG Shi-tong<sup>1</sup>

(1. School of Digital Media, Jiangnan University, Wuxi 214122, China; 2. School of Information Engineering and Technology, Changzhou Institute of Industry Technology, Changzhou 213164, China)

**Abstract:** A fuzzy  $K$ -plane clustering algorithm based on double knowledge representations about incremental feature learning and homogeneous style of data (ISF-KPC) is proposed. Before partitioning data samples into different groups, i.e., clusters, the feature augmentation is firstly conducted in an incremental manner based on the original inputs. Since data samples originating from a group share a same homogeneous style, the style information of each group, denoted by a style matrix, will be iteratively determined. By extensive experiments, the proposed ISF-KPC based on the double knowledge representations can obtain comparative clustering performance when compared to the adopted comparative clustering methods, especially it has its best clustering performance on the datasets with clearly homogeneous styles.

**Keywords:**  $K$ -plane clustering; style information; feature augmentation; fuzzy clustering

## 0 引言

在许多实际应用中, 来源于同一组或同一源的数据往往潜在或明显地表现出独特的数据风格<sup>[1-4]</sup>. 典型的应用包括癫痫脑电信号识别<sup>[5-6]</sup> (采集于健康人群的脑电信号所表现出来的波形特征明显不同于患有癫痫的人群)、手写体识别<sup>[1,4,7-8]</sup> (不同的作者写出来的字体风格完全不一样)、元音识别<sup>[1,8]</sup> (英文中的元音发音互不相同). 这些数据所表现的风格信息完全有别于数据的物理特征, 如距离、颜色和相似性. 另外, 按照人的识别思维, 人们习惯将具有相同风格的事物进行归类, 因此考虑数据风格信息的识别模型符合数据的实际情况.

聚类是一种无监督数据分析, 在未给出样本标签信息的情况下可将样本划分到不同的类别中. 聚类分析已广泛应用于识别、图像处理、计算机视觉以及文本分析<sup>[9-10]</sup>. 基于  $K$  平面的聚类算法以其新颖的聚类方式已获得越来越多科研人员关注<sup>[11-14]</sup>, 其主要思想是通过生成  $K$  个平面以代替基于中心点聚类算法中<sup>[15-17]</sup> 的  $K$  个聚类中心点. 基于  $K$  平面的聚类算法与基于中心点的聚类算法相比, 中心点浓缩了数据风格信息, 平面更利于挖掘并表达数据风格信息. 另外, 将数据风格信息用于数据分类并改进分类行为已取得相关研究进展<sup>[1-4]</sup>. Huang 等<sup>[1]</sup> 认为同一类数据拥有独特的数据风格, 不同类数据之间的风格

收稿日期: 2019-01-04; 修回日期: 2019-03-21.

基金项目: 国家自然科学基金项目 (61572236, 61300151); 常州工业职业技术学院博士基金项目 (BSJJ13101010); 常州工业职业技术学院新一代信息技术团队项目 (YB201813101005).

责任编委: 刘宝碇.

†通讯作者. E-mail: gusuhang09@163.com.

相互区别,并以矩阵的形式计算每一类数据的风格信息用于改善分类模型的性能;Jiang等<sup>[2]</sup>提出的风格集中自动解码器能够提取鲁棒的图像风格特征代表,从而有效提高时尚、建筑以及漫画等图像的分类性能;Veeramachaneni等<sup>[3-4]</sup>提出一种风格约束的文本分类方法,通过可供选择的风格假设赋予文本中每个类(即每个风格)不同高斯密度,用以区别文本中不同类的风格特征.该方法在文本数据集上的分类效果明显优于基于语义、词典(单一的数据物理特征)的文本分类方法.然而,现有针对数据风格的研究成果主要集中于有监督学习算法,对将数据风格信息用于聚类分析并提高聚类性能的关注很少.因此,针对以上所述内容,所提出算法主要集中于聚类分析且以 $K$ 平面聚类算法为基础.

在聚类分析中,由于样本标签信息未知,与有监督数据分析分类算法相比,聚类算法获得相对较少的数据特征信息.为了弥补可利用数据信息的不平衡,在聚类算法中对数据进行双知识表达,即对原输入特征增量式扩维进行增量学习并挖掘扩维后的数据风格信息.双知识表达约束下的ISF-KPC分割数据样本时使得具有同一风格的样本成为一类.本文主要贡献在于:1)双知识表达中的增量学习体现了高斯型支持向量机(support vector machine, SVM)中特征向高维映射的思想;2)双知识表达中的数据风格信息有别于数据物理特征;3)大量实验结果表明双知识表达能够很好地约束聚类行为,与所选对比算法相比,ISF-KPC至少能够取得竞争性的聚类性能,尤其在具有典型风格数据集上能够取得最好的聚类结果.

## 1 ISF-KPC

### 1.1 $K$ 平面聚类算法

由于ISF-KPC算法以 $K$ 平面聚类算法为基础,本文首先简单介绍 $K$ 平面聚类算法<sup>[1]</sup>.考虑一个数据集 $\mathbf{X}$ 包含 $N$ 个样本,即 $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ 且有 $\mathbf{x}_i \in \mathbf{R}^d$ .根据 $\mathbf{X}$ ,平面聚类算法( $K$ -plane clustering, KPC)首先生成 $K$ 个平面,然后将每个样本划分到距离最近的平面中,即

$$\begin{aligned} \min_{\omega_j, b_j} & \frac{1}{2} \|\mathbf{X}_j \omega_j + b_j \mathbf{e}\|; \\ \text{s.t.} & \|\omega_j\| = 1. \end{aligned} \quad (1)$$

其中: $\mathbf{X}_j$ 为与第 $j$ 个聚类相对应的数据子集; $\mathbf{e}$ 为具有合适维度的列向量,其元素均为1; $\omega_j$ 、 $b_j$ 为聚类平面参数,可通过迭代收敛的方式求解所有聚类平面参数.

典型的 $K$ 平面聚类为硬聚类算法,即一个样本

只能属于一个 $K$ 平面,而在ISF-KPC中引入模糊隶属度约束,即一个样本可被分割到多个聚类平面,因此式(1)可变为

$$\begin{aligned} \min_{\mu_{ij}, \omega_j, b_j} & = \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m (\mathbf{X}_i \omega_j + b_j)^2; \\ \text{s.t.} & \sum_{j=1}^k \mu_{ij} = 1, \|\omega_j\|^2 = 1. \end{aligned} \quad (2)$$

### 1.2 动机-双知识表达

ISF-KPC的聚类行为由双知识表达所约束,即增量学习和挖掘数据风格信息.增量学习保证了ISF-KPC的泛化性,数据风格信息能够提高ISF-KPC的实际聚类精确度.以下分别介绍这两个有关数据的知识表达.

与高斯型SVM相似,ISF-KPC通过增量学习将样本特征进行增长式扩维,在更高维空间中生成 $K$ 个聚类平面并完成聚类分析.首先,将原输入样本 $\mathbf{x}_i$ 一一进行投射生成增强特征 $z_i^1$ ,投射函数选择高斯函数,进而将所有增强特征组成增强节点 $\mathbf{E}^1$ ,如图1所示.将生成的增强节点嵌入到原特征空间中,此时原输入样本 $\mathbf{X}$ 变为新输入样本 $\mathbf{X}^1 = [\mathbf{X}, \mathbf{E}^1]$ .将新生成的输入样本同样利用高斯函数进行一一投射,并将生成的增强节点再次嵌入到原输入特征空间中,经过一系列特征扩维操作直至ISF-KPC取得最好聚类性能,此时原输入样本 $\mathbf{X}$ 变为新输入样本 $\mathbf{X}^s = [\mathbf{X}, \mathbf{E}^s]$ .图1形象地解释了增量学习的过程,其中增强节点 $\mathbf{E}^s$ 与原输入特征共同影响ISF-KPC的聚类性能.大量实验结果表明,在特征增量学习过程中, $s$ 的取值范围为1~3,即特征扩维操作只需执行1~3次的情况下ISF-KPC能够取得最佳聚类性能.与高斯型SVM将特征通过核化映射到无穷维相比,本文增量学习过程简单且易于执行,无需求解无穷维特征空间中的超平面参数.

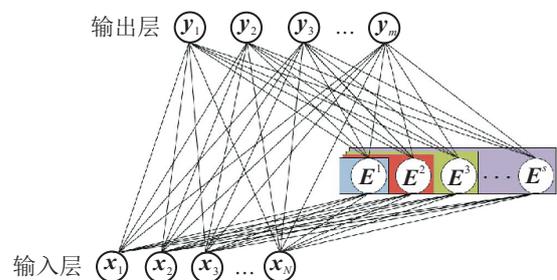


图1 增量学习

数据集中的每一类数据除了物理特征(如距离、颜色或相似性)外,还潜在或明显地拥有相应的风格,挖掘并利用数据风格信息可有效提高决策模型的性能<sup>[1-4]</sup>.如图2为基于数据物理特征的决策模型(简称

为模型1)与基于数据风格信息的决策模型(简称为模型2)之间的区别. 其中,图2左边部分展现了包含若干字母的数据集,这些字母可组成不同的字体风格. 依据数据的相似性物理特征建立的模型1将具有相似外观的字母归为一类,如图2中间部分所示. 如果在建立决策模型的过程中考虑数据风格信息,模型2则将具有相同风格的字母归为一类,如图2右边部分所示. 模型2与模型1得到截然不同的归类结果,且模型2的归类结果一方面符合数据集中每一类数据具有各自的风格,另一方面也符合人们的认知习惯. ISF-KPC致力于挖掘并利用数据集中每一类数据风格信息改善其分类行为.

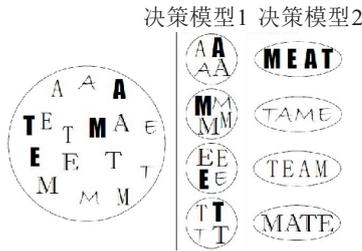


图2 两种决策模型之间的区别

### 1.3 ISF-KPC目标函数及其优化

通过第2.2节动机描述可得双知识表达约束下的ISF-KPC目标函数为

$$J_{\text{ISF-KPC}} = \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \boldsymbol{\omega}_j + b_j)^2 + \lambda \sum_{j=1}^k \|\mathbf{A}_j - \mathbf{I}\|;$$

$$\text{s.t. } \sum_{j=1}^k \mu_{ij} = 1, \|\boldsymbol{\omega}_j\|^2 = 1. \quad (3)$$

其中:  $\mathbf{x}_i^s$  为在原输入基础上连续执行增量学习后生成的新的输入;  $\mu_{ij}$  为第  $i$  个样本归属第  $j$  个聚类的模糊隶属度,且有  $\sum_{j=1}^k \mu_{ij} = 1, \mu_{ij} > 0, \forall i; \mathbf{A}_j \in R^{d \times d}$  为与数据集中第  $j$  个聚类相对应的风格矩阵,通过风格矩阵  $\mathbf{A}_j$  迭代地挖掘每个聚类的数据风格信息;  $\lambda$  为风格调节参数,用于惩罚过度的数据风格信息,可人为确定或通过网格搜索结合交叉验证的方法确定其值,第2.1节会给出根据具体数据集确定其值的指导方法,另外,当  $\lambda \rightarrow \infty$  时,数据风格信息不再约束 ISF-KPC 聚类行为,当  $\lambda$  值设置过小时,ISF-KPC 会过度地使用数据风格信息;  $\mathbf{I}$  为具有合适维度的单位矩阵. 大量的实验结果表明,当数据集中每一类数据具有典型数据风格时,风格矩阵  $\mathbf{A}_j$  中的元素值趋向于某一相对较大值,反之,风格矩阵  $\mathbf{A}_j$  中的元素值(除对角线元素)趋向于0.

与KPC算法相比,ISF-KPC具有以下显著区别: 1) ISF-KPC的目标函数中嵌入双知识表达,如式(3)所示,ISF-KPC的聚类行为受双知识表达约束,尤其数据风格信息  $\mathbf{A}_j$  能够有效地保证 ISF-KPC 的聚类精确度; 2) ISF-KPC 目标函数带有正则项  $\|\mathbf{A}_j - \mathbf{I}\|_F^2$  ( $\|\cdot\|_F^2$  为 Frobenius 范数),其作用于约束数据风格信息,由于该正则项基于数据风格信息,显著区别于 SVM 目标函数中正则项(能够很好地优化 SVM 模型性能).

### 1.4 ISF-KPC参数优化

通过 ISF-KPC 目标函数(式(3))可知,所提出聚类算法 ISF-KPC 共涉及 5 个参数,即样本模糊隶属度  $\mu_{ij}$ 、聚类平面参数  $\boldsymbol{\omega}_j$  和  $b_j$ 、风格调节参数  $\lambda$ 、风格矩阵  $\mathbf{A}_j$ . 其中:  $\lambda$  在第2.1节给出确定其值的指导方法,  $\mu_{ij}$ 、 $\boldsymbol{\omega}_j$ 、 $b_j$ 、 $\mathbf{A}_j$  四个参数通过以下两个独立程序迭代地确定其值.

#### 1) 模糊K平面参数学习.

固定所有风格矩阵  $\{\mathbf{A}_j\} (j = 1, 2, \dots, k)$  以及模糊隶属度矩阵,目标函数变为

$$J'_{\text{ISF-KPC}} = \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \boldsymbol{\omega}_j + b_j)^2,$$

$$\text{s.t. } \|\boldsymbol{\omega}_j\|^2 = 1. \quad (4)$$

初始化  $\{\mathbf{A}_j\}$  使得任意  $\mathbf{A}_j$  为单位矩阵,  $\mathbf{U}$  可通过随机赋值或 fuzzy C means (FCM) 算法<sup>[16]</sup>进行初始化,且有  $\sum_{j=1}^k \mu_{ij} = 1, \mu_{ij} > 0, \forall i$ . 与式(1)相比,此时目标函数优化问题已转换成求解标准KPC中  $K$  个聚类平面参数问题<sup>[11]</sup>. 因此可通过求解特征值问题实现  $K$  平面参数学习. 考虑拉格朗日函数将式(4)转换为

$$L'(\boldsymbol{\omega}_j, b_j, \xi_j) = \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \boldsymbol{\omega}_j + b_j)^2 - \sum_{j=1}^k \xi_j (\boldsymbol{\omega}_j^T \boldsymbol{\omega}_j - 1). \quad (5)$$

将式(5)分别对  $\boldsymbol{\omega}_j$ 、 $b_j$  的偏导求解,有

$$\frac{\partial L'}{\partial \boldsymbol{\omega}_j} = 0, \frac{\partial L'}{\partial b_j} = 0. \quad (6)$$

进而得到

$$\mathbf{D}_j \boldsymbol{\omega}_j = \xi_j \boldsymbol{\omega}_j, \quad (7)$$

$$b_j \sum_{i=1}^N \mu_{ij}^m = \mathbf{U}(:, j)^T (\mathbf{X}^s)^T \mathbf{A}_j \boldsymbol{\omega}_j. \quad (8)$$

其中

$$\mathbf{D}_j = \frac{\left( \sum_{i=1}^N \mu_{ij}^m \mathbf{A}_j^T \mathbf{x}_i^s \right) \left( \sum_{i=1}^N \mu_{ij}^m (\mathbf{x}_i^s)^T \mathbf{A}_j \right)}{\sum_{i=1}^N \mu_{ij}^m}$$

$$\sum_{i=1}^N \mu_{ij}^m \mathbf{A}_j^T \mathbf{x}_i^s (\mathbf{x}_i^s)^T \mathbf{A}_j, \quad (9)$$

$$\mathbf{D}_j = \frac{\mathbf{A}_j^T (\mathbf{X}^s)^T \mathbf{U}(:, j) \mathbf{U}(:, j)^T \mathbf{X}^s \mathbf{A}_j - \mathbf{e}^T \mathbf{U}(:, j)}{\mathbf{A}_j^T (\mathbf{X}^s)^T \Lambda \mathbf{X}^s \mathbf{A}_j}. \quad (10)$$

$\mathbf{e}$ 为具有合适维度的列向量,其元素均为1; $\Lambda$ 为对角矩阵,其对角元素分别为 $\mu_{1j}^m, \mu_{2j}^m, \dots, \mu_{Nj}^m$ . 根据文献[11]求解KPC中 $K$ 个聚类平面参数 $\{\omega_j\}$ 与 $\{b_j\}$ 的方法,式(7)中 $\xi_j$ 为 $\mathbf{D}_j$ 的最小特征值,且 $\omega_j$ 为与最小特征值 $\xi_j$ 相对应的特征向量,ISF-KPC中 $K$ 个聚类平面参数 $\{\omega_j\}$ 即可求出,此时 $b_j$ 为

$$b_j = \frac{\mathbf{U}(:, j)^T (\mathbf{X}^s)^T \mathbf{A}_j \omega_j}{\sum_{i=1}^N \mu_{ij}^m}. \quad (11)$$

至此,固定 $\{\mathbf{A}_j\}$ 和 $\mathbf{U}$ ,可求出 $K$ 个聚类平面参数 $\{\omega_j\}$ 、 $\{b_j\}$ . 此时,在 $\{\mathbf{A}_j\}$ 、 $\{\omega_j\}$ 、 $\{b_j\}$ 都已知的情况下关于目标函数(如式(3)所示)的优化问题可转变为

$$\begin{aligned} J''_{\text{ISF-KPC}} &= \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j + b_j)^2, \\ \text{s.t. } \sum_{j=1}^k \mu_{ij} &= 1. \end{aligned} \quad (12)$$

同样考虑拉格朗日函数,式(12)变为

$$\begin{aligned} L''(\mu_{ij}, \varepsilon_i) &= \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j + b_j)^2 - \\ &\sum_{i=1}^N \varepsilon_i \left( \sum_{j=1}^k \mu_{ij} - 1 \right). \end{aligned} \quad (13)$$

求解样本模糊隶属度可参照FCM算法<sup>[16]</sup>,计算公式如下:

$$\mu_{ij} = \frac{1}{\sum_{j'=1}^k \left( \frac{(\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j - b_j}{(\mathbf{x}_i^s)^T \mathbf{A}_{j'} \omega_{j'} - b_{j'}} \right)^{\frac{2}{m-1}}}. \quad (14)$$

2) 数据风格信息学习. 当单独的模糊 $K$ 平面参数学习程序结束后,即参数 $\{\omega_j\}$ 、 $\{b_j\}$ 、 $\mathbf{U}$ 已知,此时目标函数的优化问题转变为

$$\begin{aligned} J'''_{\text{ISF-KPC}} &= \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j + b_j)^2 + \\ &\lambda \sum_{j=1}^k \|\mathbf{A}_j - \mathbf{I}\|_F^2. \end{aligned} \quad (15)$$

由式(15)可知, $J'''_{\text{ISF-KPC}}$ 与 $k$ 个独立的风格矩阵相关. 考虑拉格朗日函数,式(15)变为

$$L'''(\mathbf{A}_j) = \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m ((\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j + b_j)^2 +$$

$$\lambda \sum_{j=1}^k \|\mathbf{A}_j - \mathbf{I}\|_F^2. \quad (16)$$

很明显,式(16)是一个凸函数问题,通过求偏导的方法,有

$$\frac{\partial L'''}{\partial \mathbf{A}_j} = 0, \quad (17)$$

很难求出所有风格的矩阵 $\{\mathbf{A}_j\}$ ,参照文献[18],可通过迭代的方式求解 $\{\mathbf{A}_j\}$ . 首先定义中间变量 $v_i$ 并令其为

$$v_i = (\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j + b_j, \quad (18)$$

此时每个风格矩阵 $\mathbf{A}_j$ 固定为单位矩阵. 式(16)变为

$$\begin{aligned} L'''(\mathbf{A}_j) &= \sum_{i=1}^N \sum_{j=1}^k \mu_{ij}^m v_i ((\mathbf{x}_i^s)^T \mathbf{A}_j \omega_j + b_j) + \\ &\lambda \sum_{j=1}^k \|\mathbf{A}_j - \mathbf{I}\|_F^2. \end{aligned} \quad (19)$$

由式(17)可得

$$\mathbf{A}_j = \mathbf{I} - \frac{1}{\lambda} \sum_{i=1}^N \mu_{ij}^m v_i \mathbf{x}_i^s \omega_j. \quad (20)$$

在式(18)~(20)之间迭代计算 $\mathbf{A}_j$ 与 $v_i$ ,直到满足 $\sum_{i=1}^N \|v_i^{p+1} - v_i^p\|^2 < \varphi$ 或达到最大迭代次数 $P$ <sup>[18]</sup>. 其中 $\varphi$ 为一阈值,可根据实际实验结果人为设定.

根据模糊 $K$ 平面参数学习和数据风格信息学习两个独立的程序可迭代计算出 $\mathbf{U}$ 、 $\{\omega_j\}$ 、 $\{b_j\}$ 、 $\mathbf{A}_j$ 四个参数直到满足 $\sum_{j=1}^k \|\mathbf{A}_j^{h+1} - \mathbf{A}_j^h\|^2 < \theta$ 或达到最大迭代次数 $H$ ,其中 $\mathbf{A}_j^h$ 为第 $h$ 次迭代过程中第 $j$ 个聚类的风格矩阵.

## 1.5 ISF-KPC算法及其复杂度分析

通过第1.4节对ISF-KPC目标函数优化问题的描述,可得到以下算法过程描述.

输入: 给定数据集 $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ ,每个样本 $\mathbf{x}_i \in R^d$ ,聚类数 $k$ ,双知识表达中的增量学习执行次数 $S$ ,所选高斯函数核宽度 $\sigma$ ,风格调节参数 $\lambda$ ,最大迭代次数 $P$ 和 $H$ ,迭代终止条件 $\varphi$ 和 $\theta$ ,模糊隶属度指数 $m$ ;

输出: 各聚类平面参数 $\{\omega_j\}$ 与 $\{b_j\}$ ,风格矩阵 $\{\mathbf{A}_j\}$ ,模糊隶属度矩阵 $\mathbf{U}$ ,预测标签集 $\mathbf{Y}$ .

step 1: 设定 $s = 1, \mathbf{X}^s = \mathbf{X}^0$ .

step 2: 执行双知识表达中的增量学习,生成增强节点 $\mathbf{Z}^s = [z_1^s, z_2^s, \dots, z_N^s]^T$ ,继而生成新的输入 $\mathbf{X}^s = [\mathbf{X}^{s-1}, \mathbf{Z}^s]^T$ .

step 3: 设定 $h = 0$ .

step 4: 循环直至程序收敛或达到最大迭代次数  $H$ .

step 4.1: 初始化  $\{A_j\}$ 、 $U$ .

step 4.2:  $h = h + 1$ .

step 4.3: 利用式(7)和(11)求解  $\{\omega_j\}$ 、 $\{b_j\}$ .

step 4.4: 利用式(14)确定  $U$ .

step 4.5: 迭代确定  $\{A_j\}$ .

step 4.5.1: 设定  $p = 1$ ;

step 4.5.2: 循环直至程序收敛或达到最大迭代次数  $P$ ;

step 4.5.3: 利用式(18)和(20)确定  $\{A_j\}$ .

根据 ISF-KPC 算法流程描述可得出以下算法复杂度分析(每个步骤均考虑最高阶情况下的复杂度): step 2 中需要计算每一对样本间的距离, 因此 step 2 占用的复杂度为  $O(N^2)$ . 对于 step 4.1 中初始化操作, 如果考虑模糊隶属度矩阵  $U$  中每个元素随机赋值且  $\sum_{j=1}^k \mu_{ij} = 1, \mu_{ij} > 0, \forall i$ , 则 step 4.1 占用的复杂度为  $O(2N)$ . 对于 step 4.3, 在确定各个聚类平面参数  $\omega_j$ 、 $b_j$  前需要由式(10)计算  $D_j$ , 该操作相应的复杂度为  $O(d^2N + 2dN^2 + d^3)$ .  $\omega_j$  为与  $D_j$  最小特征值相对应的特征向量, 因此确定  $\omega_j$  占用的复杂度为  $O(d^3)$ . 根据式(11)可确定  $b_j$ , 占用的复杂度为  $O(d^3 + dN^2)$ , 因此当考虑最高阶时 step 4.3 占用的复杂度大小为  $O[k(3d^3 + 3dN^2 + d^2N)]$ . 根据式(14)可确定  $U$  中每个元素, step 4.4 占用的复杂度为  $O(d^4N)$ . step 4.5 每一次迭代计算  $A_j$  的过程中需要首先由式(18)确定  $v_i$ , 该操作占用的复杂度为  $O(d^3)$ . 另外, 根据式(20), 需要  $O(N^2d + 2d^2)$  复杂度计算  $A_j$ , 因此 step 4.5 计算  $\{A_j\}$  的复杂度为  $O[kP(dN^2 + d^3)]$ . 综上所述, ISF-KPC 整个算法复杂度为  $O[SH(d^4N + 3dN^2 + kPdN^2)]$ . 在迭代次数有限的情况下, ISF-KPC 算法复杂度主要与样本维度和样本总数相关. 因此 ISF-KPC 适合于样本维度和样本总数适当情况下的聚类分析.

## 2 实验与分析

本节将在人造数据集以及具有典型数据风格的真实数据集上验证所提 ISF-KPC 的聚类性能, 并通过与其他算法的比较表明 ISF-KPC 中双知识表达的有效性.

### 2.1 对比算法及参数设置

由于 ISF-KPC 以 KPC 算法<sup>[11]</sup>为基础, 将 KPC 及其模糊化版本 F-KPC 作为对比算法. 为了验证双知

识表达的有效性, 将 ISF-KPC 的简易版本 ISF-KPC\_0, 即仅考虑单一知识表达(数据风格信息学习)作为对比算法之一. ISF-KPC 中引入模糊隶属度约束, 并将聚类中心点替换成聚类平面, 选择典型的基于中心点的聚类算法  $K$ -medoids<sup>[15]</sup> 和 FCM<sup>[16]</sup> 作为对比算法. 由于 affinity propagation (AP)<sup>[19-20]</sup> 和 density-based spatial clustering of applications with noise (DBSCAN)<sup>[21-22]</sup> 两种算法能够识别具有任意形状的数据, 将这两种算法也作为对比算法.

由式(3)和增量学习过程可知, ISF-KPC 共涉及 6 个参数, 分别为聚类平面参数  $\{\omega_j\}$ 、 $\{b_j\}$ 、模糊隶属度矩阵  $U$ 、风格矩阵  $\{A_j\}$ 、风格调节参数  $\lambda$  和增量学习高斯函数中的核宽度  $\sigma$ .  $\{\omega_j\}$ 、 $\{b_j\}$ 、 $U$  及  $\{A_j\}$  可由模糊  $K$  平面参数学习与数据风格信息学习两个独立的程序根据真实数据集迭代确定, 因此 ISF-KPC 主要确定  $\lambda$  和  $\sigma$  两个参数. 根据大量的实验结果可给出以下关于这两个参数的推荐设置: 参数  $\sigma$  的搜索范围为  $\{10^{-2}, 10^{-1}, \dots, 10^2, 10^3\}$ , 对于  $\sigma$  的每一个值, 增量学习后所有新生成的数据间的平均欧氏距离为  $\bar{D} = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N d_{ij}$ . 其中:  $N$  为样本总数,  $d_{ij}$  为第  $i$  个样本与第  $j$  个样本之间的欧氏距离. 由此可设置参数  $\lambda$  的值等于  $\bar{D}$  的量级, 或在  $\bar{D}$  的量级附近进行搜索. 对于  $\lambda$  和  $\sigma$ , 均采用网格搜索相结合交叉验证的方法进行确定<sup>[5,23]</sup>.

关于其他对比算法, ISF-KPC\_0 并没有考虑增量学习, 其涉及的参数  $\{\omega_j\}$ 、 $\{b_j\}$ 、 $U$ 、 $\{A_j\}$  以及  $\lambda$  均可参考 ISF-KPC 进行设置. KPC 参数  $k$  设置为数据集中聚类数. F-KPC 参数  $k$  设置为数据集中聚类数, 模糊隶属度矩阵  $U$  及模糊隶属度指数可参照 FCM 算法.  $K$ -medoids 参数  $k$  设置为数据集中聚类数, FCM 算法采用默认设置. AP 算法的聚类性能主要受参数参考度 PR 影响, 参考文献[19], 在样本相似度中值附近采用网格搜索结合交叉验证的方法确定 PR. DBSCAN 算法主要受参数样本邻域阈值 Eps 和样本邻域内样本个数阈值 MinPts 影响, 在文献[22, 24] 推荐设置值附近采用网格搜索结合交叉验证的方法分别确定 Eps 和 MinPts.

对于 KPC、F-KPC 以及  $K$ -medoids 三种算法, 分别运行 30 次后取平均结果, 其他算法运行 10 次后取平均结果. 所有算法均在 Matlab 平台上运行, 电脑配置为: 3.6 GHz 且 Intel(R) Core(TM) i7-4790 CPU, 8 G 内存, 64 位 Windows 10 操作系统.

2.2 数据集

本文在人造数据集上视觉地展示ISF-KPC算法的聚类性能,在真实数据集上验证ISF-KPC的实际聚类性能.

图3为本文采用的人造数据集,其中SD1、SD2、SD3中每一类数据对应典型的形状,意味着这3个数据集包含明显的数据风格.由于SD4中有一类数据为高斯随机分布,SD5中4类数据对应同样的形状,SD6中3类数据都为高斯随机分布,因此这3个数据集并不包含明显的数据风格.人造数据集详细配置信息如表1所示.

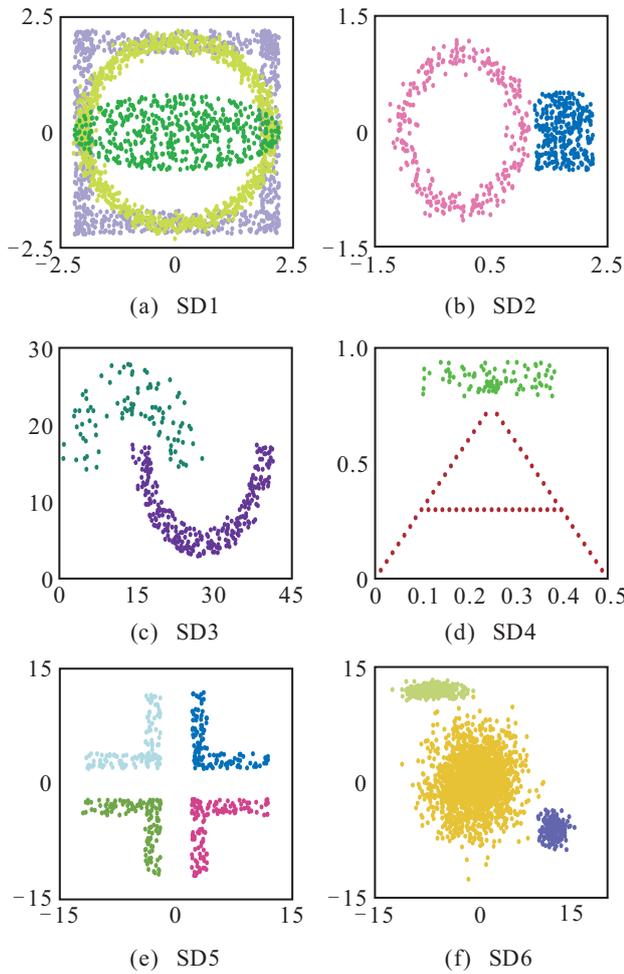


图3 人造数据集

表1 人造数据集详细配置

数据集	样本数	聚类数
SD1	2501	3
SD2	600	2
SD3	373	2
SD4	150	2
SD5	504	3
SD6	2600	3

本文在真实数据集上展示ISF-KPC独特的聚类性能,且每一个真实数据集具有明显的数据风格.如

引言所述,选取癫痫脑电信号识别、手写体识别以及元音识别作为3个案例.

对于癫痫脑电信号<sup>[5-6]</sup>,其原始信号和经核化主成分分析(kernel principal component analysis, K-PCA)降维后的特征分别如图4和图5所示.由图4和图5可见,正常人群的脑电信号(A组和B组)明显区别于患癫痫人群的脑电信号(C组、D组和E组),即使属于同一人群,所表现出的脑电信号也相互区别,如A组和B组所示.实验中分别将A组与B组,B组与D组,B组、D组与E组组成真实数据集,并分别命名为EEG-D1、EEG-D2和EEG-D3.3个真实数据集的详细配置如表2所示,每个数据集都包含不同的数据风格,且每个聚类对应的数据风格各不相同.

对于手写体识别<sup>[1,4,7-8]</sup>,实验中选取 Chinese academy of science institute of automation (CASIA) 官网公布的手写体数据集<sup>[7]</sup>,其中部分数据如图6所示.由此可见,出自每一位作者的手写体均相互区别,

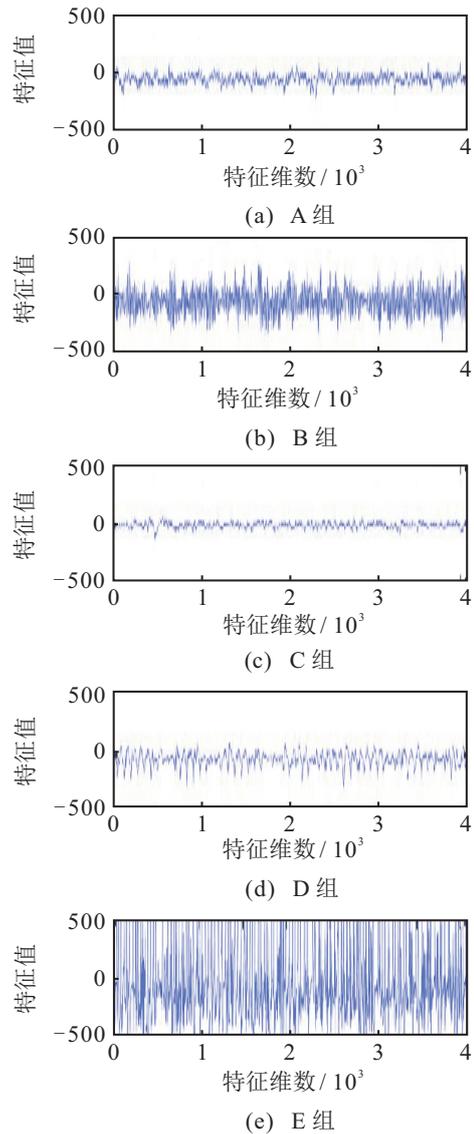


图4 EEG原始信号

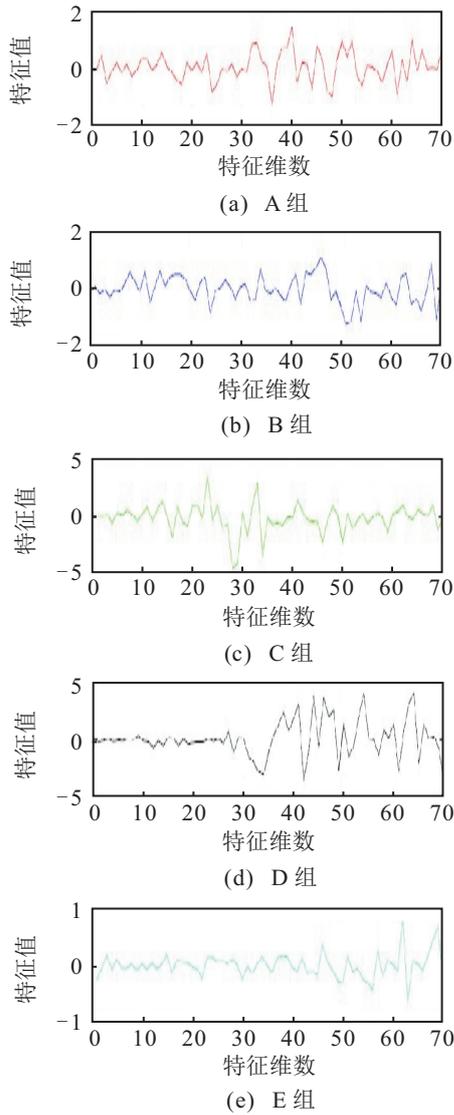


图5 经K-PCA特征降维后的EEG信号

表2 真实数据集详细配置

数据集	样本数	特征数	聚类数
EEG-D1	200	70	2
EEG-D2	200	70	2
EEG-D3	300	70	3
HWD1	4000	7	2
HWD2	4000	7	2
HWD3	6000	7	3
VD1	180	13	2
VD2	180	13	2
VD3	270	13	3

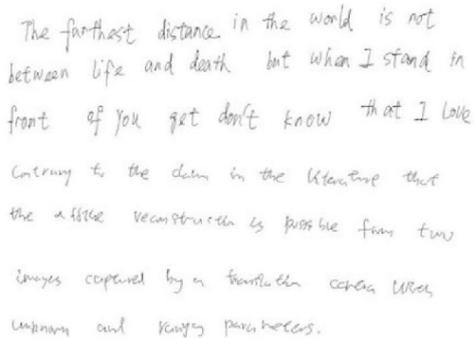


图6 手写体部分数据展示

因此相应的手写数据风格也相互区别. 实验中分别将来自以下作者的手写体数据组成真实数据集: luan 和 taojing, liuchunkai 和 lulingling, chaowenting、fuyu 和 maying, 并将这3个数据集分别命名为 HWD1、HWD2 和 HWD3. 另外, 参照文献 [1], 从每一位作者的手写体数据中随机选取 2000 个样本组成以上 3 个数据集. 3 个真实数据集的详细配置如表 2 所示.

对于元音识别<sup>[1,8]</sup>, 表 3 列出了每个具体元音展示. 由此可知, 每个元音的发音各不相同, 不管男女, 同一个元音的发音相同, 意味着同一个元音的发音具有相同的风格. 实验中分别将元音数据集 Vowel<sup>[1,8]</sup> 中标签为 3 和 5, 标签为 6 和 7, 标签为 1、4 和 9 对应的样本组成相应的真实数据集, 并分别命名为 VD1、VD2 和 VD3. 3 个真实数据集的详细配置如表 2 所示.

表3 英文中的元音展示

vowel	word	vowel	word
æ	hoard	ɑ:	hard
ʊ	hood	i:	heed
SD3	head	ɪ	hid
SD4	heard	ʌ	hud
æ	had	u:	who'd
SD6	hod		

### 2.3 结果与分析

图 7 和表 4 分别为所有对比算法在人造数据集上的聚类结果. 表 4 中: 正确率 (Acc)、F-measure 和 Rand Index (RI) 为外部类型的聚类评价指标<sup>[9-10,25]</sup>, Davies-Bouldin Index (DBI) 为内部类型的聚类评价指标<sup>[26]</sup>, “-” 代表算法参数采用默认设置, “--” 代表标准差小于  $10^{-4}$ . 最好的聚类正确率用黑体标出. 由图 7 和表 4 可得出以下结论:

1) 就聚类正确率而言, ISF-KPC 在人造数据集 SD1、SD2、SD3 和 SD5 上取得最好的聚类结果, 尤其对于包含典型数据风格的数据集 SD1、SD2 和 SD3, ISF-KPC 能够取得优越的聚类性能. 就其他内外聚类评价指标而言, ISF-KPC 至少能够取得竞争性的聚类性能.

2) 对于其他对比算法, 由于能够识别任意形状的聚类, DBSCAN 在人造数据集上能够取得较好的聚类结果.

3) 由于 ISF-KPC 考虑了数据双知识表达, 即特征增量学习和数据风格信息学习, ISF-KPC 的聚类性能明显优于仅考虑数据风格信息的 ISF-KPC\_0.

4) 将 ISF-KPC、ISF-KPC\_0 分别与 KPC 及 F-KPC 比较, 相关实验结果有力地验证了挖掘数据风格信息并用于提高聚类性能的有效性.

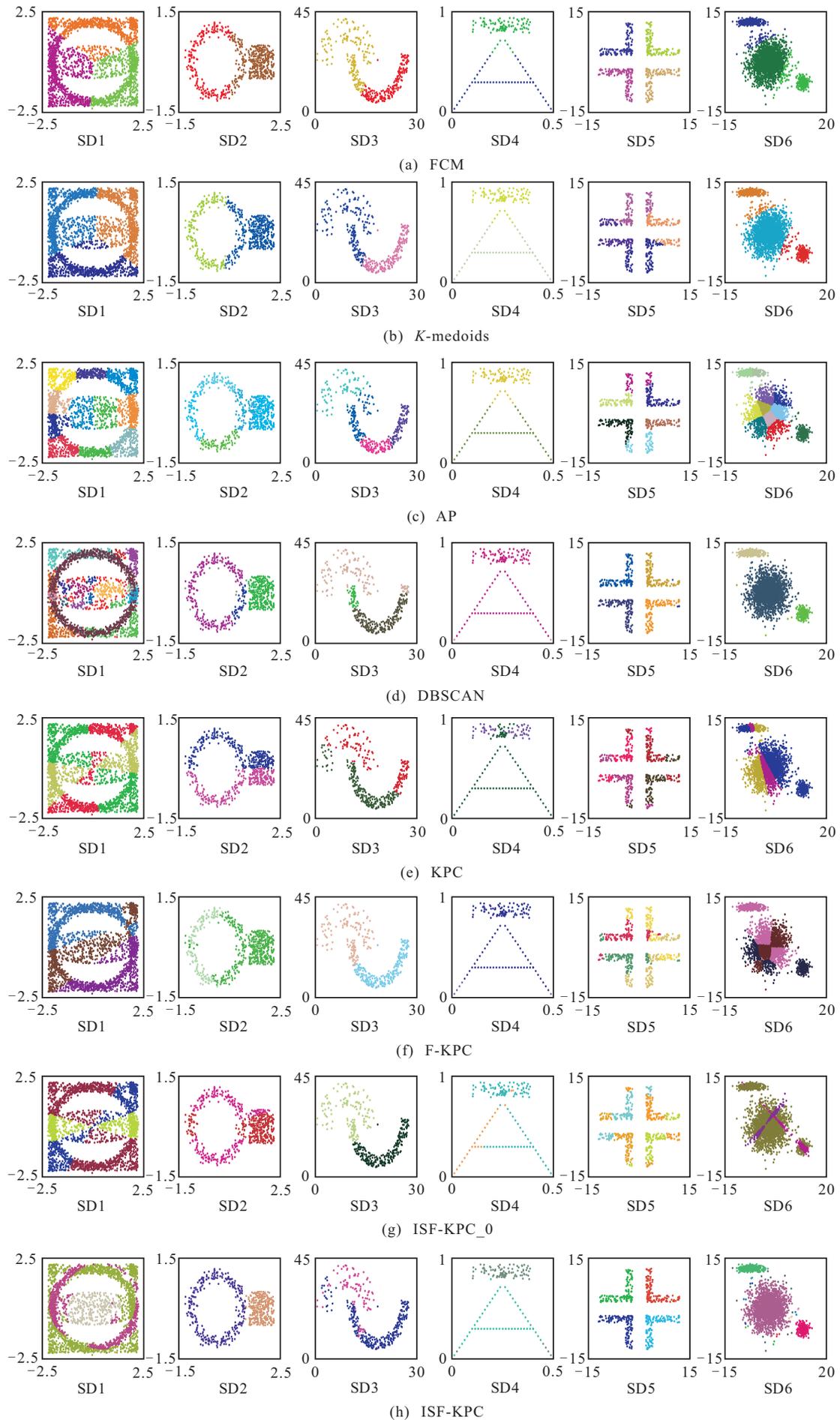


图7 对比算法在人造数据集上的聚类结果视觉展示

表 4 所有对比算法在人造数据集上的详细聚类结果

Methods	SD1					SD2				
	Acc	$F$ -measure	RI	DB	Opt	Acc	$F$ -measure	RI	DB	Opt
FCM	0.3499 (0.0077)	0.3487 (0.0010)	0.5463 (0.0013)	0.0088 (0.0001)	—	0.8483 (0)	0.7532 (—)	0.7422 (—)	0.0078 (0)	—
$K$ -medoids	0.3503 (0.0081)	0.3505 (0.0012)	0.5454 (0.0017)	0.0087 (0.0002)	$K = 3$	0.8050 (0)	0.7073 (0)	0.6855 (0)	0.0078 (0)	$K = 2$
AP	0.1887 (0)	0.2077 (—)	0.6398 (—)	0.0067 (0.0003)	PR = -82.12	0.5350 (0)	0.6103 (0)	0.7077 (—)	0.0074 (—)	PR = -27.31
DBSCAN	0.5590 (0)	0.7082 (—)	0.8369 (0)	0.0123 (0)	Eps = 0.20 MinPts = 13	0.9000 (0)	0.9100 (0)	0.9176 (0)	0.0095 (0)	Eps = 0.15 MinPts = 6
KPC	0.3886 (0)	0.3580 (0)	0.5494 (0)	0.0612 (0.0305)	$K = 3$	0.5050 (0)	0.4985 (0)	0.4992 (0)	0.0243 (0)	$K = 2$
F-KPC	0.4294 (0.0102)	0.3701 (0.0061)	0.5625 (0.0033)	0.0213 (0.0051)	$K = 3$	0.7333 (0)	0.6566 (0)	0.6082 (0)	0.0086 (0)	$K = 2$
ISF-KPC <sub>0</sub>	0.4482 (0.0245)	0.4008 (0.0487)	0.5122 (0.0534)	0.0422 (0.0520)	$\lambda = 10^2$	0.7983 (0.0638)	0.7654 (0.0464)	0.7229 (0.0376)	0.0452 (0.0210)	$\lambda = 10^{-1}$
ISF-KPC	<b>0.6062</b> (0.0712)	0.5943 (0.0493)	0.6179 (0.0901)	0.0991 (0.0210)	$\lambda = 10$ $\sigma_1 = 10^2, \sigma_2 = 10$	<b>0.9900</b> (0.0033)	0.9864 (0.0044)	0.9816 (0.0032)	0.0062 (0.0001)	$\lambda = 10^4$ $\sigma_1 = 10^2, \sigma_2 = 10$
Methods	SD3					SD4				
	Acc	$F$ -measure	RI	DB	Opt	Acc	$F$ -measure	RI	DB	Opt
FCM	0.7748 (0)	0.6858 (—)	0.6501 (0)	0.0086 (—)	—	0.9467 (—)	0.9042 (—)	0.8983 (—)	0.0053 (—)	—
$K$ -medoids	0.7668 (0)	0.6778 (0)	0.6414 (0)	0.0086 (0)	$K = 2$	0.9467 (—)	0.9042 (—)	0.8983 (—)	0.0053 (—)	$K = 2$
AP	0.4665 (0)	0.5128 (0)	0.5775 (0)	0.0078 (0.0002)	PR = -4031.25	0.9467 (0)	0.9042 (0)	0.8983 (0)	0.0053 (0)	PR = -0.79
DBSCKN	0.8284 (0)	0.7746 (0)	0.7638 (0)	0.0064 (0)	Eps = 2.15 MinPts = 25	<b>0.9733</b> (—)	0.9802 (0)	0.9800 (0)	0.0051 (0)	Eps = 0.15 MinPts = 9
KPC	0.7534 (0)	0.6705 (0)	0.6274 (0)	0.0220 (0)	$K = 2$	0.7733 (0)	0.6625 (0)	0.6471 (0)	0.0116 (0)	$K = 2$
F-KPC	0.8642 (0.0025)	0.7928 (0.0036)	0.7646 (0.0037)	0.0087 (—)	$K = 2$	0.5689 (0.0567)	0.6455 (0.0248)	0.5127 (0.0194)	0.0155 (0.0080)	$K = 2$
ISF-KPC <sub>0</sub>	0.8939 (0.0247)	0.6987 (0.0078)	0.6745 (0.0247)	0.0101 (0.0015)	$\lambda = 10$	0.8283 (0.0279)	0.8430 (0.0334)	0.8099 (0.0095)	0.0954 (0.0467)	$\lambda = 1$
ISF-KPC	<b>0.9008</b> (0.0169)	0.9046 (0.0858)	0.9141 (0.0898)	0.0301 (0.0217)	$\lambda = 1$ $\sigma_1 = 10^2, \sigma_2 = 10^2$	0.9600 (0.0851)	0.9716 (0.0841)	0.9725 (0.0893)	0.0049 (0.0098)	$\lambda = 1$ $\sigma_1 = 10^2, \sigma_2 = 10^2$
Methods	SD5					SD6				
	Acc	$F$ -measure	RI	DB	Opt	Acc	$F$ -measure	RI	DB	Opt
FCM	0.9004 (0.1614)	0.8722 (0.1985)	0.9364 (0.0988)	0.0081 (0.0003)	—	0.9465 (—)	0.9236 (—)	0.9102 (—)	0.0062 (0)	—
$K$ -medoids	0.6845 (0)	0.5656 (0.0001)	0.7839 (0.0001)	0.0082 (0.0004)	$K = 4$	0.9638 (—)	0.9477 (0)	0.9373 (—)	0.0053 (0.0006)	$K = 3$
AP	0.7579 (0)	0.7612 (—)	0.8992 (—)	0.0078 (0.0005)	PR = -1235.18	0.3008 (0)	0.2696 (0)	0.4747 (—)	0.0082 (0.0001)	PR = -1197.86
DBSCAN	0.9921 (0)	0.9921 (0)	0.9961 (0)	0.0057 (0)	Eps = 1.00 MinPts = 2	<b>0.9981</b> (0)	0.9976 (0)	0.9970 (0)	0.0027 (—)	Eps = 2.00 MinPts = 3
KPC	0.5142 (0)	0.5056 (0)	0.6839 (0)	0.0126 (0)	$K = 4$	0.3815 (0)	0.4478 (0)	0.4600 (0)	0.0447 (0)	$K = 3$
F-KPC	0.6766 (0.0449)	0.5340 (0.0581)	0.7681 (0.0285)	0.0148 (0.0007)	$K = 4$	0.4985 (0.0266)	0.5013 (0.0221)	0.4956 (0.0078)	0.0250 (0.0044)	$K = 3$
ISF-KPC <sub>0</sub>	0.7202 (0.0804)	0.5324 (0.0290)	0.7675 (0.0144)	0.1005 (0.0019)	$\lambda = 10^{-2}$	0.7226 (0.0943)	0.6678 (0.0789)	0.7124 (0.0328)	0.1078 (0.0515)	$\lambda = 1$
ISF-KPC	<b>1</b> (0)	<b>1</b> (0)	<b>1</b> (0)	0.0080 (0.0004)	$\lambda = 10^4$ $\sigma = 10^2$	0.9846 (—)	0.9812 (0.0001)	0.9771 (0.0003)	0.0035 (—)	$\lambda = 10^{-3}$ $\sigma_1 = 10^2, \sigma_2 = 10, \sigma_3 = 10$

表5 所有对比算法在真实数据集上的详细聚类结果

Methods	EEG-D1					EEG-D2				
	Acc	F-measure	RI	DB	Opt	Acc	F-measure	RI	DB	Opt
FCM	0.5225 (0.0105)	0.4990 (0.0030)	0.4987 (0.0010)	0.0584 (0.0015)	—	0.5265 (0.0195)	0.4978 (0.0031)	0.4997 (0.0028)	0.0699 (0.0019)	—
K-medoids	0.5850 (0)	0.5108 (0)	0.5120 —	0.0666 —	K = 2	0.5100 —	0.6600 (0)	0.4977 —	0.0060 (0)	K = 2
AP	0.5006 (0)	0.6644 (—)	0.4975 (—)	0.0732 (0.0006)	PR = -1 262.13	0.4925 (0)	0.6034 —	0.4867 —	0.0305 (0.0002)	PR = -2 553.33
DBSCAN	0.5750 (0)	0.6383 (0)	0.5088 (0)	0.0681 (0)	Eps = 3.00 MinPts = 15	0.5900 (0)	0.6307 (0)	0.5138 (0)	0.0337 (0)	Eps = 3.95 MinPts = 20
KPC	0.5150 (0)	0.6630 (0)	0.5302 (0)	0.0956 (0)	K = 2	0.5400 (0)	0.6329 —	0.5211 (0)	0.0974 —	K = 2
F-KPC	0.5367 (0.0455)	0.5488 (0.0605)	0.5295 (0.0642)	0.0979 (0.0044)	K = 2	0.5883 (0.0861)	0.6562 (0.0902)	0.5795 (0.1069)	0.0897 (0.0113)	K = 2
ISF-KPC <sub>0</sub>	0.5260 (0.0254)	0.4986 (0.0040)	0.5001 (0.0043)	0.1134 (0.0014)	$\lambda = 10^2$	0.5190 (0.0086)	0.4965 (0.0012)	0.4984 (0.0008)	0.1186 (0.0072)	$\lambda = 10^2$
ISF-KPC	<b>0.6325</b> (0.0100)	0.6618 (0.0008)	0.5289 (0.0010)	0.0422 (0.0085)	$\lambda = 1$ $\sigma_1 = 10^2, \sigma_2 = 10$	<b>0.6521</b> (0.0025)	0.6839 (0.0685)	0.5638 (0.0001)	0.0603 (0.0274)	$\lambda = 10^{-2}$ $\sigma_1 = 10, \sigma_2 = 10$

Methods	ECG-D3					HWD1				
	Acc	F-measure	RI	DB	Opt	Acc	F-measure	RI	DB	Opt
FCM	0.4500 —	0.4277 —	0.5265 (0)	0.0502 0.0022	—	0.5542 —	0.5122 (0)	0.5058 —	0.0072 —	—
K-medoids	0.3690 (0.0030)	0.4040 (0.0001)	0.4966 (0.0006)	0.0436 (0.0003)	K = 3	0.5659 (0.0184)	0.5169 (0.0087)	0.5092 (0.0043)	0.0072 (0.0002)	K = 2
AP	0.3564 —	0.4822 —	0.3379 (0)	0.0825 (0.0002)	PR = -2 684.71	0.2813 (0)	0.3822 —	0.6182 —	0.0048 (0.0002)	PR = $10^9$
DBSCAN	0.5600 (0)	0.5208 (0)	0.5578 (0)	0.0069 (0)	Eps = 2.55 MinPts = 5	0.5093 (0)	0.6560 (0)	0.5006 (0)	0.0001 (0)	EPS = 1.00 MinPts = 11
KPC	0.3667 (0)	0.4923 (0)	0.3626 (0)	0.0773 (0)	K = 3	0.4802 (0)	0.6348 —	0.5094 (0)	0.1057 (0)	K = 2
F-KPC	0.3522 (0.0042)	0.4878 (0.0020)	0.3589 (0.0035)	0.0208 (0.0028)	K = 3	0.6674 (0.1364)	0.6643 (0.0017)	0.5931 (0.0762)	0.0233 (0.0113)	K = 2
ISF-KPC <sub>0</sub>	0.5467 (0.0156)	0.5065 (0.0245)	0.6424 (0.0125)	0.1023 (0.0147)	$\lambda = 10^2$	0.6797 (0.1156)	0.5771 (0.0568)	0.6014 (0.0588)	0.0503 (0.0204)	$\lambda = 10^6$
ISF-KPC	<b>0.6191</b> (0.0167)	0.5542 (0.0062)	0.5645 —	0.0119 (0.0041)	$\lambda = 10^{-1}$ $\sigma_1 = 10$	<b>0.8081</b> (0.0059)	0.7028 (0.0098)	0.6899 (0.0072)	0.0314 (0.0002)	$\lambda = 10^9$ $\sigma = 10^2, \sigma_2 = 10^{-1}$

Methods	HWD2					HWD3				
	Acc	F-measure	RI	DB	Opt	Acc	F-measure	RI	DB	Opt
FCM	0.5038 (0)	0.5147 (0)	0.4999 (0)	0.0090 —	—	0.3773 —	0.3399 —	0.5562 —	0.0053 (0.0008)	—
K-medoids	0.5397 (0.0511)	0.5238 (0.0026)	0.5083 (0.0127)	0.0090 (0.0004)	K = 2	0.3845 (0.0040)	0.3440 (0.0016)	0.5570 (0.0011)	0.0055 (0.0007)	K = 3
AP	0.3472 —	0.4291 —	0.6367 —	0.0053 (0.0002)	PR = $-10^9$	0.1990 —	0.2057 —	0.6320 —	0.0066 (0.0003)	PR = $-10^9$
DBSCAN	0.5058 (0)	0.6528 (0)	0.5006 (0)	0.0003 (0)	Eps = 2.25 MinPts = 10	0.3473 (0)	0.4877 (0)	0.3683 (0)	0.0012 (0)	Eps = 3.00 MinPts = 4
KPC	0.5287 (0)	0.6361 —	0.5127 (0)	0.1104 —	K = 2	0.3862 (0)	0.4853 (0)	0.3620 —	0.0568 (0)	K = 3
F-KPC	0.6990 (0)	0.6282 (0)	0.6239 (0)	0.0284 (0)	K = 2	0.5927 (0.0055)	0.6679 (0.0091)	0.6938 (0.0100)	0.1028 (0.0074)	K = 3
ISF-KPC <sub>0</sub>	0.6591 (0.1040)	0.5795 (0.0647)	0.5722 (0.0657)	0.0522 (0.0299)	$\lambda = 10^6$	0.6100 (0.0623)	0.6616 (0.0267)	0.7109 (0.0314)	0.1053 (0.0347)	$\lambda = 10^6$
ISF-KPC	<b>0.7138</b> (0.0595)	0.6019 (0.0543)	0.5984 (0.0509)	0.0348 (0.0082)	$\lambda = 10^7$ $\sigma = 10^2$	<b>0.6185</b> (0.0499)	0.6700 (0.0523)	0.6977 (0.0718)	0.3297 (0.1213)	$\lambda = 10^7$ $\sigma = 10^3$

表5 (续)

Methods	VD1					VD2				
	Acc	$F$ -measure	RI	DB	Opt	Acc	$F$ -measure	RI	DB	Opt
FCM	0.511 1 (0)	0.495 1 --	0.497 5 (0)	0.007 0 --	—	0.500 0 (0)	0.495 5 --	0.497 2 --	0.007 6 (0)	—
$K$ -medoids	0.505 6 (0)	0.495 2 --	0.497 3 (0)	0.007 0 (0)	$K = 2$	0.500 0 (0)	0.495 5 --	0.497 2 --	0.007 8 --	$K = 2$
AP	0.516 7 (0)	0.508 5 --	0.497 8 (0)	0.007 6 --	PR = -788.07	0.505 6 (0)	0.495 2 --	0.497 3 (0)	0.007 0 --	PR = -655.85
DBSCAN	0.494 4 (0)	0.476 8 --	0.498 8 (0)	0.006 6 --	Eps = 2.50 MinPts = 20	0.561 1 (0)	0.581 0 (0)	0.504 7 (0)	0.007 2 (0)	Eps = 2.15 MinPts = 19
KPC	0.507 8 (0)	0.510 1 --	0.498 4 (0)	0.045 0 (0)	$K = 2$	0.527 8 --	0.563 6 (0)	0.500 8 (0)	0.119 5 (0)	$K = 2$
F-KPC	0.503 7 (0.022 4)	0.542 2 (0.036 2)	0.509 8 (0.008 7)	0.076 2 (0.032 1)	$K = 2$	0.530 0 (0.009 1)	0.522 9 (0.010 6)	0.532 6 (0.005 5)	0.066 4 (0.018 7)	$K = 2$
ISF-KPC <sub>0</sub>	0.530 0 (0.015 2)	0.528 6 (0.038 2)	0.499 5 (0.001 7)	0.084 8 (0.083 1)	$\lambda = 10$	0.555 6 (0.025 6)	0.524 9 (0.039 0)	0.504 7 (0.004 9)	0.062 1 (0)	$\lambda = 10$
ISF-KPC	<b>0.622 0</b> (0.013 9)	0.577 5 (0.073 1)	0.502 1 (0.002 6)	0.047 1 (0.022 6)	$\lambda = 10^{-2}$ $\sigma_1 = 10^2, \sigma_2 = 10$	<b>0.616 7</b> (0.005 6)	0.580 3 (0.004 1)	0.519 5 (0.000 7)	0.054 9 (0.036 2)	$\lambda = 10^{-2}$ $\sigma = 10^{-1}, \sigma_2 = 1$

Methods	VD3				
	Acc	$F$ -measure	RI	DB	Opt
FCM	0.355 6 (0)	0.328 3 (0)	0.554 9 --	0.009 6 (0.001 4)	—
$K$ -medoids	0.362 6 (0.016 4)	0.339 4 (0.002 1)	0.549 4 (0.004 1)	0.010 4 (0.001 1)	$K = 3$
AP	0.333 3 --	0.395 3 --	0.497 4 --	0.008 4 --	PR = -859.61
DBSCAN	0.400 0 (0)	0.438 6 (0)	0.762 1 --	0.008 1 (0)	Eps = 2.50 MinPts = 5
KPC	0.373 3 (0)	0.360 1 (0)	0.532 6 (0)	0.069 1 (0)	$K = 3$
F-KPC	0.396 3 (0.024 6)	0.368 3 (0.032 1)	0.534 6 (0.096 8)	0.053 5 (0.016 1)	$K = 3$
ISF-KPC <sub>0</sub>	0.381 5 (0.022 1)	0.352 3 (0.016 2)	0.545 7 (0.010 5)	0.110 0 (0.028 6)	$\lambda = 10$
ISF-KPC	<b>0.455 6</b> (0.050 0)	0.464 3 (0.038 3)	0.565 8 (0.074 6)	0.036 7 (0.010 8)	$\lambda = 10^{-1}$ $\sigma_1 = 10, \sigma_2 = 10, \sigma_3 = 1$

5) 在原输入特征基础上的增量学习次数只需 1 ~ 3 次, 增量学习过程简单且易于执行。

表 5 详细列出了所有对比算法在真实数据集上的聚类结果, 且每一个真实数据集均包含典型的数据风格。由表 5 可见, 就聚类正确率而言, ISF-KPC 在所有真实数据集上取得最好的聚类结果, 就其他内外部聚类评价指标而言, ISF-KPC 在多数情况下能够取得最好的聚类性能。与其他对比算法相比 (除 ISF-KPC<sub>0</sub> 之外), 由于考虑了数据风格信息, ISF-KPC 在

多数情况下至少能够取得竞争性的聚类性能。

为了进一步比较 ISF-KPC 与其他对比算法的区别, 利用文献 [27] 的统计测试方法对所有对比算法的聚类性能进行分析。

该统计测试方法主要包含  $F_F$  和 CD 两个关键值 [27]。  $F_F$  用于确定该统计测试方法的空假设 (即假设所有对比算法具有相同的聚类性能) 是否被否定; CD 用于验证对比算法之间的聚类性能是否存在显著区别。由表 5 列出的聚类正确率对所有对比算法

进行排序,如对于EEG-D2,因ISF-KPC取得了最好的聚类结果,其排名为1,即 $\text{rank} = 1$ .依此类推,可得出所有对比算法在所有真实数据集上的排序,如表6所示.根据对比算法排序值、 $F$ -分布及其自由度( $N_c - 1$ )、 $(N_c - 1)(N_d - 1)$ 可确定 $F_F$ 值.只要满足 $F_F > F((N_c - 1), (N_d - 1))$ ,该统计测试方法的空假设即被否定.其中 $N_c$ 代表所有对比算法个数, $N_d$ 代表

真实数据集个数, $F(\cdot, \cdot)$ 根据 $F$ -分布表可查出<sup>[27]</sup>.基于空假设被否定,可进一步计算CD值.其中 $\text{CD} = q_\alpha \sqrt{\frac{N_c(N_c + 1)}{6N_d}}$ ( $\alpha$ 为显著性度,其值取文献[27]推荐值 $\alpha = 0.05$ , $q_\alpha$ 可从文献[27]查出).任意两个对比算法平均排序差值的绝对值大于该CD值即表明这两个对比算法的聚类性能之间存在显著区别.

表6 所有对比算法在真实数据集上的排序

methods	FCM	$K$ -medoids	AP	DBSCAN	KPC	$F$ -KPC	ISF-KPC_0	ISF-KPC
EEG-D1	6	2	8	3	7	4	5	1
EEG-D2	5	7	8	2	4	3	6	1
EEG-D3	4	5	7	2	6	8	3	1
HWD1	5	4	8	6	7	3	2	1
HWD2	7	4	8	6	5	2	3	1
HWD3	6	5	8	7	4	3	1	1
VD1	4	6	3	8	5	7	2	1
VD2	7	7	6	2	5	4	3	1
VD3	7	6	8	2	5	3	4	1
average rank(AR)	5.67	5.11	7.11	4.22	5.33	4.11	3.33	1.00

由表5和表6可知, $F(7, 56) \approx 2.18$ 、 $F_F \approx 8.87$ ,因此该统计测试方法的空假设被否定.由给出的显著性度 $\alpha = 0.05$ 可计算 $\text{CD} \approx 3.12$ .根据统计测试结果可得出以下结论:

1) 由于ISF-KPC与其他对比算法(除ISF-KPC\_0与F-KPC外)之间的平均排序差值都大于 $\text{CD} \approx 3.12$ ,ISF-KPC与其他对比算法之间存在显著区别.另外,由于ISF-KPC考虑了双知识表达,其聚类性能优于ISF-KPC\_0和F-KPC.

2) 由于挖掘并利用了数据风格信息约束聚类行为,ISF-KPC\_0的聚类性能优于其他对比算法(除ISF-KPC外),特别地,ISF-KPC\_0与AP算法之间存在显著区别.

3) 结合1)和2),在真实数据集上的统计测试结果充分表明了数据双知识表达确实能够提高聚类算法性能,尤其基于数据风格信息的聚类算法,其聚类性能能够显著区别于典型聚类算法.

### 3 结论

由于数据集中的每一类数据都潜在或明显地具有独特的数据风格,这些风格信息显著区别于数据物理特征.如何挖掘隐藏的数据风格信息并用于聚类分析是一个值得研究的课题.本文利用风格矩阵迭代地捕捉数据集中每个聚类的数据风格信息,这些风格矩阵相互区别.另外,沿着高斯型SVM将输入特征

从低维空间映射到高维空间的思想,本文在原输入特征空间中增量式地对原输入特征进行扩维.通过增量学习和数据风格矩阵对数据进行双知识表达以强化聚类性能.人造数据集尤其真实数据集上的实验结果验证了ISF-KPC中双知识表达的有效性.真实数据集上的统计测试结果表明ISF-KPC与典型聚类方法之间存在显著区别.未来的研究会重点关注如何将ISF-KPC扩展到回归模型,进一步研究本文双知识表达并将ISF-KPC推广到深度学习.

### 参考文献(References)

- [1] Huang K Z, Jiang H C, Zhang X Y. Field support vector machines[J]. IEEE Transactions on Emerging Topics in Computational Intelligence, 2017, 1(6): 454-463.
- [2] Jiang S H, Shao M, Jia C C, et al. Learning consensus representation for weak style classification[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2018, 40(12): 2906-2919.
- [3] Veeramachaneni S, Nagy G. Analytical results on style-constrained Bayesian classification of pattern fields[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2007, 29(7): 1280-1285.
- [4] Sarkar P, Nagy G. Style consistent classification of isogenous patterns[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2005, 27(1): 88-98.
- [5] Jiang Y Z, Deng Z H, Wang S T, et al. Recognition of epileptic EEG signals using a novel multiview TSK fuzzy system[J]. IEEE Transactions on Fuzzy Systems, 2017,

- 25(1): 3-20.
- [6] Xie L X, Deng Z H, Wang S T, et al. Generalized hidden-mapping transductive transfer learning for recognition of epileptic electroencephalogram signals[J]. IEEE Transactions on Cybernetics, DOI: 10.1109/TCYB.2018.2821764.
- [7] Chinese Academy of Sciences Institute of Automation (CASIA). Handwriting database[DB/OL]. [2018-10-07]. <http://biometrics.idealtest.org/>.
- [8] Zhang X Y, Huang K Z, Liu C L. Pattern field classification with style normalized transformation[C]. Proceedings of the 22th International Joint Conference on Artificial Intelligence. Spain: IEEE, 2011: 1621-1626.
- [9] 陈爱国, 王士同. 基于多代表的大规模数据模糊聚类算法[J]. 控制与决策, 2016, 31(12): 2122-2130. (Chen A G, Wang S T. Fuzzy clustering algorithm based on multiple medoids for large-scale data[J]. Control and Decision, 2016, 31(12): 2122-2130.)
- [10] 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.)
- [11] Bradley P S, Mangasarian O L.  $k$ -plane clustering[J]. Journal of Global Optimization, 2000, 16(1): 23-32.
- [12] Shao Y H, Bai L, Wang Z, et al. Proximal plane clustering via eigenvalues[J]. Procedia Computer Science, 2013, 17: 41-47.
- [13] Jayadeva, Khemchandani R, Chandra S, et al. Twin support vector machines for pattern classification[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2007, 29(5): 905-910.
- [14] Wang Z, Shao Y H, Bai L, et al. Twin support vector machine for clustering[J]. IEEE Transactions on Neural Networks and Learning Systems, 2015, 26(10): 2583-2588.
- [15] Park H S, Jun C H. A simple and fast algorithm for  $K$ -medoids clustering[J]. Expert Systems with Applications, 2009, 36(2): 3336-3341.
- [16] Bai X, Chen Z, Zhang Y, et al. Infrared ship target segmentation based on spatial information improved FCM[J]. IEEE Transactions on Cybernetics, 2016, 46(12): 3259-3271.
- [17] Anand S, Mittal S, Tuzel O, et al. Semi-supervised kernel mean shift clustering[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2014, 36(6): 1201-1215.
- [18] Fang X Z, Wong W K, Teng S H, et al. Flexible affinity matrix learning for unsupervised and semisupervised classification[J]. IEEE Transactions on Neural Networks and Learning Systems, DOI: 10.1109/TNNLS.2018.2861839.
- [19] Frey B J, Dueck D. Clustering by passing messages between data points[J]. Science, 2007, 315: 972-976.
- [20] Arzeno N M, Vikalo H. Semi-supervised affinity propagation with soft instance-level constraints[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2015, 37(5): 1041-1052.
- [21] Sakai T, Tamura K, Kitakami H. Density-based adaptive spatial clustering algorithm for identifying local high-density areas in georeferenced documents[C]. Proceedings IEEE International Conference System, Man, and Cybernetics. New York: IEEE, 2014: 513-518.
- [22] Ester M, Kriegel H P, Sander J, et al. A density-based algorithm for discovering clusters in large spatial databases with noise[C]. Proceedings of ACM Conference on Knowledge Discovery and Data Mining. New York: ACM, 1996: 226-231.
- [23] Jiang Y Z, Deng Z H, Wang S T, et al. Realizing two-view TSK fuzzy classification system by using collaborative learning[J]. IEEE Transactions on Systems, Man, and Cybernetics: Systems, 2017, 47(1): 145-160.
- [24] Hou J, Gao H J, Li X L. DSets-DBSCAN: A parameter-free clustering algorithm[J]. IEEE Transactions on Image Processing, 2016, 25(7): 3182-3193.
- [25] Wang Y T, Chen L H, Mei J P. Incremental fuzzy clustering with multiple medoids for large data[J]. IEEE Transactions on Fuzzy Systems, 2014, 22(6): 1557-1568.
- [26] Davies D L, Bouldin D W. A cluster separation measure[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 1979, PAMI-1(2): 224-227.
- [27] Demsar J. Statistical comparisons of classifiers over multiple data sets[J]. Journal of Machine Learning Research, 2006, 7: 1-30.

## 作者简介

顾苏杭(1989—), 男, 博士生, 从事人工智能与模式识别、机器学习的研究, E-mail: gusuhang09@163.com;

王士同(1964—), 男, 教授, 博士生导师, 从事人工智能与模式识别、机器学习、深度学习等研究, E-mail: wxwangst@yahoo.com.cn.

(责任编辑: 郑晓蕾)