

## 基于模糊 Fisher 准则的自适应降维模糊聚类算法

支晓斌<sup>①</sup> 范九伦<sup>②</sup>

<sup>①</sup>(西安电子科技大学电子工程学院 西安 710071)

<sup>②</sup>(西安邮电学院信息与控制系 西安 710121)

**摘要:** 该文指出曹苏群等人提出的基于模糊 Fisher 准则(FFC)的半模糊聚类算法(FFC-SFCA)中的一个推导错误, 结合模糊紧性和分离性(FCS)聚类算法提出新的聚类算法: FFC-FCS。FFC-FCS 充分利用 FFC 的特征提取和降维特性, 交替运行原始数据空间中 FFC 和投影空间中的 FCS, 通过对降维数据的聚类实现对原始数据的聚类。FFC-FCS 不仅对低维数据具有优异的分类性能而且对高维数据也表现出一定的分类优势。实验结果表明, FFC-FCS 的性能明显优于原有的 FCS 算法, FFC-SFCA 算法以及经典的模糊 C-均值(FCM)算法。

**关键词:** 模糊散布矩阵; 模糊 Fisher 准则; 最优投影矢量; FCS 聚类

中图分类号: TP181

文献标识码: A

文章编号: 1009-5896(2009)11-2653-06

## Fuzzy Fisher Criterion Based Adaptive Dimension Reduction Fuzzy Clustering Algorithm

Zhi Xiao-bin<sup>①</sup> Fan Jiu-lun<sup>②</sup>

<sup>①</sup>(School of Electronic Engineering, Xidian University, Xi'an 710071, China)

<sup>②</sup>(Department of Information and Control, Xi'an Institute of Post and Telecommunications, Xi'an 710121, China)

**Abstract:** The derivation mistake in Cao's Fuzzy Fisher Criterion (FFC) based Semi-Fuzzy Clustering Algorithm (FFC-SFCA) is pointed out. Combining Fuzzy Compactness and Separation (FCS) clustering algorithm, a new clustering algorithm, FFC-FCS, is proposed in this paper. FFC-FCS make full use of the feature extraction and dimension reduction characteristics of FFC, alternately running FFC in the original data space and FCS in the projection space, clustering the original data is accomplished by clustering the dimension reduction data. FFC-FCS not only shows excellent capability of classifying low dimensional data but also has a certain grade classification advantage with respect to high dimensional data. The experimental results show that FFC-FCS has super performance over original FCS, FFC-SFCA and classical Fuzzy C-Means(FCM).

**Key words:** Fuzzy scatter matrix; Fuzzy Fisher Criterion(FFC); Optimal projection vector; FCS clustering

### 1 引言

聚类分析是多元统计分析的方法之一, 也是无监督模式识别的一个重要分支。聚类方法将相似的数据点分为一类, 而将不相似的数据点分为不同的类。在众多的聚类算法中, FCM 聚类算法<sup>[1]</sup>是最为著名的聚类算法之一, 由于它在一定程度上克服了传统硬 C-均值聚类算法对初值敏感, 抗噪性差等缺点, 得到了学者们的广泛关注<sup>[2-4]</sup>。然而 FCM 是一种基于紧性致度量的算法, 它只考虑了聚类的类内紧致性而没有考虑类间分离性<sup>[5]</sup>。为了进一步提高 FCM 的性能, Wu 等人提出了模糊散布矩阵的概念<sup>[5]</sup>, 在此基础上通过修改聚类有效性指标 FS(c)<sup>[6]</sup>提出了模糊紧性和分离性(FCS)聚类算法。FCS 将 FCM 作为特例, 是 FCM 的广义形式。由于 FCS 同

时考虑了类内紧致性和类间分离性, 聚类性能得到了进一步的提高。然而 FCS 的缺点是需手动选取参数, 自适应性差。

Fisher 判别分析(FDA)是监督模式识别的一种重要的特征提取和数据降维方法, FDA 经过寻找最具有判别性的投影方向, 使得数据经投影处理后, 不仅维数得到降低, 而且获得了最大类内紧致性和类间分离性。在此基础上的分类也就变得更加高效和容易。但是 FDA 是一种监督方法, 需要知道已有数据的类标号。2002 年, Clausi 提出了 KIF(K-means Iterative Fisher)方法<sup>[7]</sup>, 将 FDA 与硬 C-均值聚类算法结合。该算法通过 K-means 初始化与迭代 FDA 组合实现。但该方法的抗噪性差。鉴于此, 曹苏群等人在模糊散布矩阵概念的基础上, 提出了模糊 Fisher 准则(FFC)<sup>[8]</sup>, 并以 FFC 为目标函数, 提出了一种新的聚类算法: FFC-SFCA。该算法通过

对 FFC 迭代优化, 不仅能够实现聚类而且同时得到最优投影矢量。

本文指出曹苏群等人提出的 FFC-SFCA 中聚类中心的表达式是错误的, 这种错误会导致算法不收敛。在此基础上我们给出新的基于 FFC 的模糊聚类算法: FFC-FCS。FFC-FCS 算法由 FFC 嵌套低维空间中的 FCS 构成。FFC-FCS 充分利用 FFC 的特征提取和降维特性, 交替运行原始数据空间中 FFC 和投影空间中的 FCS, 通过对降维数据的聚类实现对原始数据的聚类。与 FCS 不同, FFC-FCS 中 FCS 的参数就是最优投影矢量对应的特征值, 在算法中是不断迭代寻优的, 因此 FFC-FCS 可以看作是参数自适应选取的 FCS 算法。因为 FFC-FCS 中的 FCS 是在原始数据所在空间经 FFC 降维后得到的投影空间中进行的, 所以 FFC-FCS 不仅可以获得优异的分类性能并且可以实现对高维数据的聚类。实验表明, FFC-FCS 总体性能明显优于原有的 FCS 算法, FFC-SFCA 算法以及经典的 FCM 算法。

## 2 FCS 和 FFC-SFCA 算法简介

### 2.1 FCS

鉴于 FCM 仅考虑了类内紧致性, 而没有考虑类间分离性, 为了能够同时描述聚类的类内紧致性和类间分离性, Wu 等人提出了模糊总散布矩阵  $\mathbf{S}_{\text{FT}}$ , 模糊类内散布矩阵  $\mathbf{S}_{\text{FW}}$  和模糊类间散布矩阵  $\mathbf{S}_{\text{FB}}$  的概念<sup>[5]</sup>。它们分别定义为

$$\mathbf{S}_{\text{FT}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^{\text{T}} \quad (1)$$

$$\mathbf{S}_{\text{FW}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{x}_i - \mathbf{v}_j)(\mathbf{x}_i - \mathbf{v}_j)^{\text{T}} \quad (2)$$

$$\mathbf{S}_{\text{FB}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}})(\mathbf{v}_j - \bar{\mathbf{x}})^{\text{T}} \quad (3)$$

其中  $\mathbf{v}_j = \sum_{i=1}^n u_{ij}^m \mathbf{x}_i / \left( \sum_{i=1}^n u_{ij}^m \right)$  称为模糊样本均值,  $u_{ij}$

$\in [0, 1]$ ,  $\sum_{j=1}^c u_{ij} = 1$ ,  $m > 1$ ; 并且  $\mathbf{S}_{\text{FT}}$ ,  $\mathbf{S}_{\text{FW}}$  和  $\mathbf{S}_{\text{FB}}$

满足  $\mathbf{S}_{\text{FT}} = \mathbf{S}_{\text{FW}} + \mathbf{S}_{\text{FB}}$ 。通过修改聚类有效性指标

$$\text{FS}(c) = \text{tr}(\mathbf{S}_{\text{FW}}) - \text{tr}(\mathbf{S}_{\text{FB}}) = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{x}_i - \mathbf{v}_j\|^2 -$$

$$\sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{v}_j - \bar{\mathbf{x}}\|^2 \quad [6], \text{ Wu 等人提出了 FCS 聚类算}$$

法。FCS 的目标函数如下:

$$J_{\text{FCS}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{x}_i - \mathbf{v}_j\|^2 - \sum_{j=1}^c \sum_{i=1}^n \eta_j u_{ij}^m \|\mathbf{v}_j - \bar{\mathbf{x}}\|^2 \quad (4)$$

其中  $u_{ij} \in [0, 1]$ ,  $\sum_{j=1}^c u_{ij} = 1$ ,  $m > 1$ ,  $\eta_j \geq 0$ 。当  $\eta_j = 0$

时,  $J_{\text{FCS}} = J_{\text{FCM}}$ , 当  $\eta_j = 1$  时,  $J_{\text{FCS}} = \text{FS}(c)$ 。式

(4)的前半部分  $\sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{x}_i - \mathbf{v}_j\|^2 = \text{tr}(\mathbf{S}_{\text{FW}})$  衡量

了类内紧致性, 后半部分  $\sum_{j=1}^c \sum_{i=1}^n \eta_j u_{ij}^m \|\mathbf{v}_j - \bar{\mathbf{x}}\|^2$  衡量

了类间分离性。定义 Lagrange 函数。

$$L_{\text{FCS}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{x}_i - \mathbf{v}_j\|^2 - \sum_{j=1}^c \sum_{i=1}^n \eta_j u_{ij}^m \|\mathbf{v}_j - \bar{\mathbf{x}}\|^2 + \sum_{i=1}^n \alpha_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \quad (5)$$

用  $L_{\text{FCS}}$  对  $\mathbf{v}_j$  和  $u_{ij}$  分别求偏导数, 并令偏导数为零, 得到一对方程:

$$\mathbf{v}_j = \frac{\sum_{i=1}^n u_{ij}^m \mathbf{x}_i - \sum_{i=1}^n \eta_j u_{ij}^m \bar{\mathbf{x}}}{\sum_{i=1}^n u_{ij}^m - \sum_{i=1}^n \eta_j u_{ij}^m} \quad (6)$$

和

$$u_{ij} = \frac{\left( \|\mathbf{x}_i - \mathbf{v}_j\|^2 - \eta_j \|\mathbf{v}_j - \bar{\mathbf{x}}\|^2 \right)^{\frac{1}{m-1}}}{\sum_{k=1}^c \left( \|\mathbf{x}_i - \mathbf{v}_k\|^2 - \eta_j \|\mathbf{v}_k - \bar{\mathbf{x}}\|^2 \right)^{\frac{1}{m-1}}} \quad (7)$$

鉴于由式(7)确定的  $u_{ij}$  可能出现负值, 为了使  $u_{ij}$  的值在  $[0, 1]$  区间内, 文献[5]对  $u_{ij}$  做了如下的修正:

如果  $\|\mathbf{x}_i - \mathbf{v}_j\|^2 \leq \eta_j \|\mathbf{v}_j - \bar{\mathbf{x}}\|^2$ , 则  $u_{ij} = 1$ ,

且对其它的  $j' \neq j$ ,  $u_{ij'} = 0$  (8)

另外, Wu 等人给出了用下面的式子确定参数  $\eta_j$  的方法:

$$\eta_j = \frac{(\beta/4) \min_{j' \neq j} \|\mathbf{v}_j - \mathbf{v}_{j'}\|^2}{\max_k \|\mathbf{v}_k - \bar{\mathbf{x}}\|^2}, \quad 0 \leq \beta \leq 1 \quad (9)$$

FCS 通过交替优化式(6)和式(7) 最小化  $L_{\text{FCS}}$ , 从而最小化  $J_{\text{FCS}}$ , 具体算法请参见文献[5]。FCS 算法的缺陷是算法性能依赖于参数  $\beta$  的选取, 而  $\beta$  需要手动选取。

### 2.2 FFC-SFCA

曹苏群等人在 Wu 等人提出的模糊散布矩阵概念的基础上, 提出了模糊 Fisher 准则<sup>[8]</sup>。记  $\mathbf{y} = \boldsymbol{\omega}^{\text{T}} \mathbf{x}$ , 在该投影空间, 各类样本均值向量记为  $\tilde{\mathbf{v}}_j$ , 有  $\tilde{\mathbf{v}}_j = \boldsymbol{\omega}^{\text{T}} \mathbf{v}_j$ 。模糊类内散布矩阵  $\tilde{\mathbf{S}}_{\text{FW}}$  定义为

$$\tilde{\mathbf{S}}_{\text{FW}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{y}_i - \tilde{\mathbf{v}}_j)(\mathbf{y}_i - \tilde{\mathbf{v}}_j)^{\text{T}} = \boldsymbol{\omega}^{\text{T}} \mathbf{S}_{\text{FW}} \boldsymbol{\omega} \quad (10)$$

模糊类间散布矩阵  $\tilde{\mathbf{S}}_{\text{FB}}$  定义为

$$\tilde{\mathbf{S}}_{\text{FB}} = \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\tilde{\mathbf{v}}_j - \bar{\mathbf{y}})(\tilde{\mathbf{v}}_j - \bar{\mathbf{y}})^{\text{T}} = \boldsymbol{\omega}^{\text{T}} \mathbf{S}_{\text{FB}} \boldsymbol{\omega} \quad (11)$$

定义模糊 Fisher 准则(FCC)函数:

$$J_{\text{FFC}} = \frac{\tilde{\mathbf{S}}_{\text{FB}}}{\tilde{\mathbf{S}}_{\text{FW}}} = \frac{\boldsymbol{\omega}^{\text{T}} \mathbf{S}_{\text{FB}} \boldsymbol{\omega}}{\boldsymbol{\omega}^{\text{T}} \mathbf{S}_{\text{FW}} \boldsymbol{\omega}} \quad (12)$$

据此, 曹苏群等人提出了基于 FFC 的半模糊聚类算法: FFC-SFCA. FFC-SFCA 以 FFC 为目标函数, 即  $J_{\text{FFC}} = \frac{\omega^T \mathbf{S}_{\text{FB}} \omega}{\omega^T \mathbf{S}_{\text{FW}} \omega}$ . 定义 Lagrange 函数为

$$L_{\text{FFC}} = \omega^T \mathbf{S}_{\text{FB}} \omega - \lambda \omega^T \mathbf{S}_{\text{FW}} \omega + \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \quad (13)$$

将  $L_{\text{FFC}}$  分别对  $\omega, \mathbf{v}_j$  和  $u_{ij}$  求偏导数, 并令其为零, 得到  $L_{\text{FFC}}$  取极大值必须满足的方程为

$$\mathbf{S}_{\text{FW}}^{-1} \mathbf{S}_{\text{FB}} \omega = \lambda \omega \quad (14)$$

$$\mathbf{v}_j = \frac{\sum_{i=1}^n u_{ij}^m (\mathbf{x}_i - (1/\lambda) \bar{\mathbf{x}})}{\sum_{i=1}^n u_{ij}^m (1 - 1/\lambda)} \quad (15)$$

$$u_{ij} = (\omega^T (\mathbf{x}_i - \mathbf{v}_j) (\mathbf{x}_i - \mathbf{v}_j)^T \omega - (1/\lambda) \omega^T (\mathbf{v}_j - \bar{\mathbf{x}}) \cdot (\mathbf{v}_j - \bar{\mathbf{x}})^T \omega)^{-\frac{1}{m-1}} \left/ \left[ \sum_{k=1}^c (\omega^T (\mathbf{x}_i - \mathbf{v}_k) (\mathbf{x}_i - \mathbf{v}_k)^T \omega - (1/\lambda) \omega^T (\mathbf{v}_k - \bar{\mathbf{x}}) (\mathbf{v}_k - \bar{\mathbf{x}})^T \omega)^{-\frac{1}{m-1}} \right] \right. \quad (16)$$

由式(16)确定的  $u_{ij}$  可能出现负值, 为了使  $u_{ij}$  的值在  $[0,1]$  区间内, 文献[8]对  $u_{ij}$  做了如下的修正: 如果  $\omega^T (\mathbf{x}_i - \mathbf{v}_j) (\mathbf{x}_i - \mathbf{v}_j)^T \omega < (1/\lambda) \omega^T (\mathbf{v}_j - \bar{\mathbf{x}}) (\mathbf{v}_j - \bar{\mathbf{x}})^T \omega$ , 则

$$u_{ij} = 1, \text{ 且对其它的 } j' \neq j, u_{ij'} = 0 \quad (17)$$

FFC-SFCA 以硬 C-均值聚类算法作为初始化, 通过交替优化式(14), 式(15)和式(16)最小化  $L_{\text{FFC}}$ , 从而最小化  $J_{\text{FFC}}$ , 具体算法参见文献[8].

需要强调指出的是用式(13)对  $\mathbf{v}_j$  求偏导数不能导出如式(15)中各类中心  $\mathbf{v}_j$  的表达式. 事实上,

$$\begin{aligned} L_{\text{FFC}} &= \omega^T \mathbf{S}_{\text{FB}} \omega - \lambda \omega^T \mathbf{S}_{\text{FW}} \omega + \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \\ &= \omega^T \left( \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}}) (\mathbf{v}_j - \bar{\mathbf{x}})^T \right) \omega \\ &\quad - \lambda \omega^T \left( \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{x}_i - \mathbf{v}_j) (\mathbf{x}_i - \mathbf{v}_j)^T \right) \omega \\ &\quad + \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \end{aligned}$$

用  $L_{\text{FFC}}$  对  $\mathbf{v}_j$  求偏导数, 并令其等于零可得如下方程:

$$2\omega^T \sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}}) + 2\lambda \omega^T \sum_{i=1}^n u_{ij}^m (\mathbf{x}_i - \mathbf{v}_j) = 0$$

整理得

$$\omega^T \sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}} + \lambda (\mathbf{x}_i - \mathbf{v}_j)) = 0 \quad (18)$$

我们无法从上式反解出  $\mathbf{v}_j$ . 观察式(18)可知, 文献[8]

的错误在于直接在上式中将  $\omega^T$  约去, 得方程  $\sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}} + \lambda (\mathbf{x}_i - \mathbf{v}_j)) = 0$ , 进而解出关于  $\mathbf{v}_j$  的表达式式(15). 这种错误导致 FFC-SFCA 在迭代求解过程中可能是不收敛的(在实验过程中我们发现常常会出现这样的现象: FFC-SFCA 在迭代过程中总是迭代到设置的最大迭代次数时才停止迭代, 见第 4 节的表 1 和表 3).

实验中, FFC-SFCA 表现出一定的分类能力, 其原因在于当  $\sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}} + \lambda (\mathbf{x}_i - \mathbf{v}_j)) = 0$  时,  $\omega^T \sum_{i=1}^n u_{ij}^m (\mathbf{v}_j - \bar{\mathbf{x}} + \lambda (\mathbf{x}_i - \mathbf{v}_j))$  也必然为零. 即 FFC-SFCA 在一个人为缩小的区域内搜寻最优解, 这样可能只会得到一个近似最优解, 而不是最优解.

### 3 FCM-FCS 算法

鉴于文献[8]的错误, 本文重新考虑基于模糊 Fisher 准则的聚类算法. 由  $\tilde{\mathbf{S}}_{\text{FW}}$  和  $\tilde{\mathbf{S}}_{\text{FB}}$  的定义和式(13)可知

$$\begin{aligned} L_{\text{FFC}} &= \omega^T \mathbf{S}_{\text{FB}} \omega - \lambda \omega^T \mathbf{S}_{\text{FW}} \omega + \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \\ &= \tilde{\mathbf{S}}_{\text{FB}} - \lambda \tilde{\mathbf{S}}_{\text{FW}} + \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \\ &= \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\tilde{\mathbf{v}}_j - \bar{\mathbf{y}}) (\tilde{\mathbf{v}}_j - \bar{\mathbf{y}})^T \\ &\quad - \lambda \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m (\mathbf{y}_i - \tilde{\mathbf{v}}_j) (\mathbf{y}_i - \tilde{\mathbf{v}}_j)^T \\ &\quad + \sum_{i=1}^n \lambda_i \left( \sum_{j=1}^c u_{ij} - 1 \right) = -\lambda \left( \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{y}_i - \tilde{\mathbf{v}}_j\|^2 \right. \\ &\quad \left. - \sum_{j=1}^c \sum_{i=1}^n \frac{1}{\lambda} u_{ij}^m \|\tilde{\mathbf{v}}_j - \bar{\mathbf{y}}\|^2 + \sum_{i=1}^n \left( -\frac{\lambda_i}{\lambda} \right) \left( \sum_{j=1}^c u_{ij} - 1 \right) \right) \end{aligned}$$

其中  $\tilde{\mathbf{v}}_j = \omega^T \mathbf{v}_j$ ,  $\bar{\mathbf{y}} = \omega^T \bar{\mathbf{x}}$ ,  $\mathbf{y}_i = \omega^T \mathbf{x}_i$ .

令  $\lambda' = 1/\lambda$ ,  $\lambda'_i = -\lambda_i/\lambda$ , 则有

$$\begin{aligned} L_{\text{FFC}} &= -\lambda \left( \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{y}_i - \tilde{\mathbf{v}}_j\|^2 - \sum_{j=1}^c \sum_{i=1}^n \lambda' u_{ij}^m \|\tilde{\mathbf{v}}_j - \bar{\mathbf{y}}\|^2 \right. \\ &\quad \left. + \sum_{i=1}^n \lambda'_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \right) \end{aligned}$$

令

$$\begin{aligned} L'_{\text{FFC}} &= \sum_{j=1}^c \sum_{i=1}^n u_{ij}^m \|\mathbf{y}_i - \tilde{\mathbf{v}}_j\|^2 - \sum_{j=1}^c \sum_{i=1}^n \lambda' u_{ij}^m \|\tilde{\mathbf{v}}_j - \bar{\mathbf{y}}\|^2 \\ &\quad + \sum_{i=1}^n \lambda'_i \left( \sum_{j=1}^c u_{ij} - 1 \right) \quad (19) \end{aligned}$$

则当  $\lambda > 0$  时, 最大化  $L_{\text{FFC}}$  等价于最小化  $L'_{\text{FFC}}$ 。用  $L'_{\text{FFC}}$  分别对  $\tilde{\mathbf{v}}_j$  和  $u_{ij}$  求偏导数, 并令偏导数为零, 可得如下的一对方程:

$$\tilde{\mathbf{v}}_j = \frac{\sum_{i=1}^n u_{ij}^m \mathbf{y}_i - \sum_{i=1}^n \lambda' u_{ij}^m \bar{\mathbf{y}}}{\sum_{i=1}^n u_{ij}^m - \sum_{i=1}^n \lambda' u_{ij}^m} \quad (20)$$

$$u_{ij} = \frac{(\|\mathbf{y}_i - \tilde{\mathbf{v}}_j\|^2 - \lambda' \|\tilde{\mathbf{v}}_j - \bar{\mathbf{y}}\|^2)^{-\frac{1}{m-1}}}{\sum_{k=1}^c (\|\mathbf{y}_i - \tilde{\mathbf{v}}_k\|^2 - \lambda' \|\tilde{\mathbf{v}}_k - \bar{\mathbf{y}}\|^2)^{-\frac{1}{m-1}}} \quad (21)$$

则当  $\lambda > 0$  且值固定时(此时  $\lambda'$  自然也为定值), 最小化  $L'_{\text{FFC}}$  可以由交替迭代式(20)和式(21)实现。式(19), 式(20)和式(21)与式(5), 式(6)和式(7)的对比可知:  $L'_{\text{FFC}}$  与 FCS 的目标函数  $L_{\text{FCS}}$  是一致的, 只是 FCS 算法的参数  $\eta_j$  被这里的  $\lambda'$  取代。

综上所述有如下结论: 当最优投影矢量  $\boldsymbol{\omega}$  及其对应的特征值  $\lambda$  (对应于  $\mathbf{S}_{\text{FW}}^{-1} \mathbf{S}_{\text{FB}}$  的最大特征值, 是大于零的)固定时, 交替迭代式(20)和式(21)最小化  $L'_{\text{FFC}}$  就是在投影空间中进行的以  $L'_{\text{FFC}}$  为目标函数的 FCS 算法。又因为最大化  $L_{\text{FFC}}$  等价于最小化  $L'_{\text{FFC}}$ , 所以当最优投影矢量  $\boldsymbol{\omega}$  及其对应的特征值  $\lambda$  固定时, 交替迭代式(20)和式(21)可实现在投影空间中进行的以  $L_{\text{FFC}}$  为目标函数的 FCS 算法。另外, 既然 FFC 可以求得最优投影矢量, 可以设想在投影空间中进行聚类。这样既容易分类, 又由于数据得到降维, 使得聚类算法的效率更高。因此, 本文考虑利用 FFC 来求得最优投影矢量, 在投影空间进行 FCS 的聚类算法。

因为要利用 FFC 来求得最优投影矢量, 所以要计算模糊类内散布矩阵  $\mathbf{S}_{\text{FW}}$  和模糊类间散布矩阵  $\mathbf{S}_{\text{FB}}$ , 由模糊类内散布矩阵  $\mathbf{S}_{\text{FW}}$  和模糊类间散布矩阵  $\mathbf{S}_{\text{FB}}$  的定义可知, 要计算模糊类内散布矩阵  $\mathbf{S}_{\text{FW}}$  和模糊类间散布矩阵  $\mathbf{S}_{\text{FB}}$  需要知道各类的模糊样本均值  $\mathbf{v}_j$  和隶属度  $u_{ij}$ , 事实上只要知道隶属度  $u_{ij}$ , 我们可以由  $\mathbf{v}_j = \sum_{i=1}^n u_{ij}^m \mathbf{x}_i / \left( \sum_{i=1}^n u_{ij}^m \right)$  求得各类的模糊样本均值  $\mathbf{v}_j$ 。而隶属度  $u_{ij}$  可以由投影空间中的 FCS 聚类算法获得, 这样就形成一个迭代优化的过程。这种基于 FFC 的求得最优投影矢量和在投影空间中的 FCS 的交替聚类算法, 记为 FFC-FCS。FFC-FCS 可以看作是参数自适应选取的 FCS 算法, 参数由最优投影矢量  $\boldsymbol{\omega}$  对应的特征值  $\lambda$  的倒数给出。FFC-FCS 算法在运行过程中的初始隶属度  $u_{ij}$  由 FCM 给出, 在迭代过程中  $u_{ij}$  由投影空间中 FCS 计算出的结果传递出来。整个算法流程如图 1 所示。

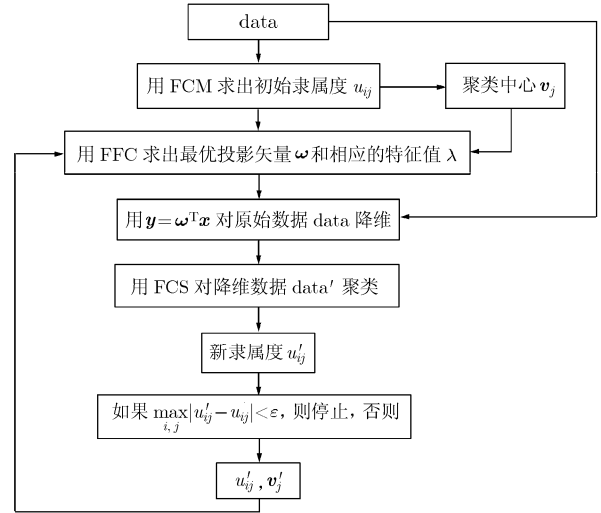


图 1 FFC-FCS 算法流程图

算法步骤:

- (1) 运行 FCM 算法给出初始隶属度  $u_{ij}$  和各类模糊样本均值  $\mathbf{v}_j$  (即聚类中心);
- (2) 利用(1)求得的隶属度  $u_{ij}$  和各类中心  $\mathbf{v}_j$  和 FFC 求出矩阵  $\mathbf{S}_{\text{FW}}^{-1} \mathbf{S}_{\text{FB}}$  的最大特征值  $\lambda$ , 并取  $\boldsymbol{\omega}$  为矩阵  $\mathbf{S}_{\text{FW}}^{-1} \mathbf{S}_{\text{FB}}$  属于  $\lambda$  的模为 1 的特征向量(即最优投影矢量), 并用  $\mathbf{y} = \boldsymbol{\omega}^T \mathbf{x}$  对数据降维;
- (3) 在降维的投影空间中执行以  $1/\lambda$  为参数的 FCS 并更新隶属度  $u'_{ij}$ ;
- (4) 如果  $\max_{i,j} |u'_{ij} - u_{ij}| < \epsilon$ , 停止; 否则, 更新各类中心  $\mathbf{v}'_j$  并返回到(2)。

## 4 实验结果及分析

为了验证本文提出的 FFC-FCS 算法的有效性, 本节用 FCM, FCS, FFC-SFCA 和 FFC-FCS 4 个算法分别对 1 个人造数据和两个标准数据 Iris 和 Glass<sup>[9]</sup>进行仿真实验。在本文的实验中, 为避免算法陷入局部最优, 算法在随机初始化后运行 100 次取最优分类结果, 并取相应的迭代次数和 CPU 运行时间。最大迭代次数设为 200, 停止阈值设为  $10^{-5}$ 。4 个算法中的参数  $m$  都取为 2, FCS 算法中的参数  $\eta_j$  取为 0.5。

### 实验 1 人造数据

图 2 所示的长条型数据来自文献[10]中的数据“fourlines”, 选其中长度接近的两类数据, 为了增大聚类的难度, 将两类数据间的距离变小, 新数据记为“twolines”。4 个算法对“twolines”的分类结果如图 2 和图 3 所示。其中, “•”和“+”分别表示两类数据, “o”表示聚类中心。图 2 中的 3 幅图 2(a), 2(b)和 2(c)分别是 FCM, FCS 和 FFC-SFCA 对“twolines”的分类结果。图 3 中的 3 幅图 3(a),

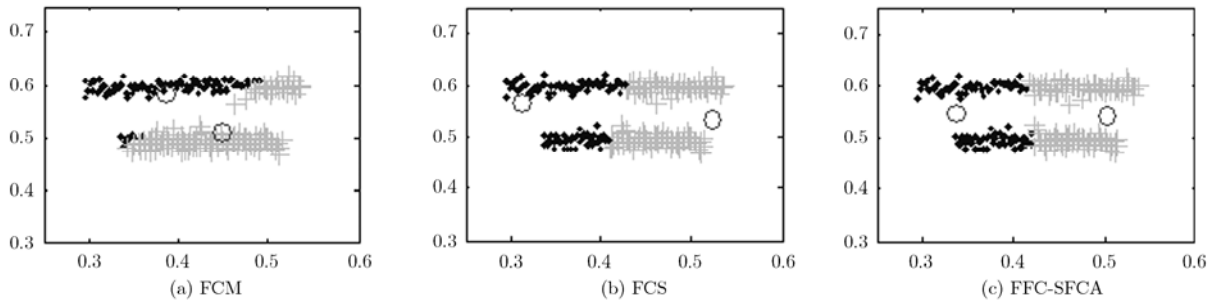


图 2 FCM, FCS 和 FFC-SFCA 对“twolines”的分类结果

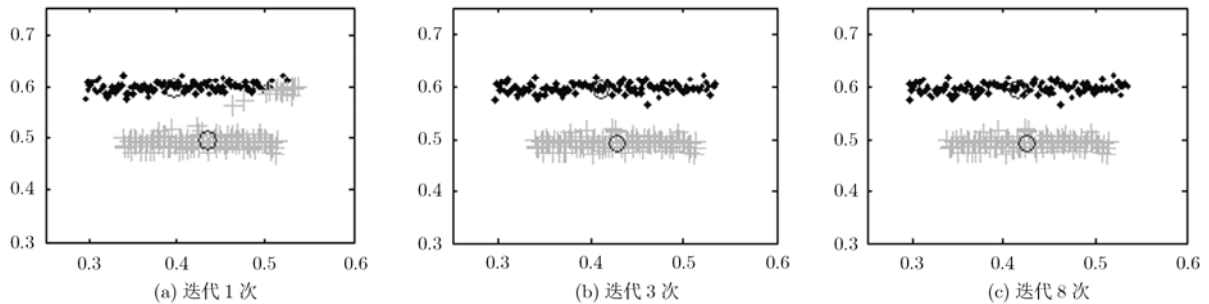


图 3 FFC-FCS 不同迭代次数时对“twolines”的分类结果

3(b)和 3(c)分别是 FFC-FCS 对“twolines”聚类迭代过程中迭代次数为 1 次, 3 次和 8 次时分类结果。由图 2 可知, FCM, FCS 和 FFC-SFCA 都不能对“twolines”正确分类。由图 3 可知, 本文提出的 FFC-FCS 算法可以在迭代过程中自适应地优化投影方向, 最终找到最优投影方向并对“twolines”正确分类。由表 1 可知, FFC-FCS 迭代次数和运行时间比 FCM 和 FCS 略多一些。但少于 FFC-SFCA。FFC-SFCA 的迭代次数达到最大迭代次数, 导致用时也较多。

表 1 4 个算法对“twolines”分类时的迭代次数和运行时间

	迭代次数	运行时间(s)
FCM	67	0.234000
FCS	5	0.187000
FFC-SFCA	200	8.875000
FCC-FCS	8	0.647000

### 实验 2 标准数据

本节选取的两个标准数据 Iris 和 Glass 均来自于 UCI repository of machine learning databases<sup>[9]</sup>, 数据的特性如表 2 所示。本节用这两个数据对 4 个聚类算法的分类性能进行比较, 用错分数, 迭代次数和运行时间 3 个指标综合衡量聚类算法的优劣。结果如表 3 和表 4 所示。其中“×”表示该算法不能运行。

表 2 数据描述

	样本数	维数	类数
Iris	150	4	3
Glass	214	9	7

表 3 4 个算法对 IRIS 数据分类结果

	错分数	迭代次数	运行时间(s)
FCM	16	24	0.047000
FCS	16	28	0.188000
FFC-SFCA	5	200	10.531000
FCC-FCS	2	14	0.547000

表 4 4 个算法对 Glass 数据分类结果

	错分数	迭代次数	运行时间(s)
FCM	109	200	0.594000
FCS	109	200	7.500000
FFC-SFCA	×	×	×
FCC-FCS	16	17	9.438000

由表 3 可知, FFC-FCS 对 IRIS 的分类精度最高, 要明显优于 FCM 和 FCS。FFC-SFCA 也达到很高的分类精度, 但是 FFC-SFCA 的迭代次数达到最大迭代次数, 导致用时也较多。

由表 4 可知, FFC-FCS 对 Glass 的分类精度最

高,要明显优于 FCM 和 FCS,FFC-FCS 的用时相对多一些。FFC-SFCA 对 Glass 分类时算法不能运行。

综上所述,本文提出的 FFC-FCS 虽然用时稍微多一些,但是在分类精度上表现出优异的性能,FFC-FCS 算法总体上要优于 FCM, FCS 和 FFC-SFCA。

## 5 结论

本文首先指出曹苏群等人提出的 FFC-SFCA 聚类算法中的一个推导错误,在分析该算法错误根源的基础上提出一个新的基于 FFC 的聚类算法:FFC-FCS。FFC-FCS 充分利用 FFC 的特征提取和降维优势,由原始数据空间的 FFC 和投影空间中的 FCS 聚类交替进行实现,可以在迭代过程中不断优化投影方向,同时完成最优投影矢量的寻找和对原始数据的聚类,FFC-FCS 可以看作是参数自适应选取的 FCS。实验表明,本文提出的 FFC-FCS 的算法性能在总体上要优于 FCM, FCS 和 FFC-SFCA。本文提出的 FFC-FCS 本质上是一种线性聚类算法,如何通过“核”技巧构造核版本的 FFC-FCS 聚类算法是本文下一步的工作。

## 参 考 文 献

- [1] Bezdek J C. Pattern Recognition with Fuzzy Objective Function Algorithms[M]. New York, Plenum Press, 1981: 95-107.
- [2] Yang Miin-Shen and Wu Kuo-Lung, *et al.* Alpha-cut implemented fuzzy clustering algorithms and switching regressions[J]. *IEEE Transactions on Systems, Man, and Cybernetics, Part B*, 2008, 38(3): 588-603.
- [3] Cai Wei-ling, Chen Song-can, and Zhang Dao-qiang. Fast and robust fuzzy c-means clustering algorithms incorporating local information for image segmentation[J]. *Pattern Recognition*, 2007, 40(3): 825-838.
- [4] Xing Hong-jie and Hu Bao-gang. An adaptive fuzzy c-means clustering-based mixtures of experts model for unlabeled data classification [J]. *Neurocomputing*, 2008, 71(4): 1008-1021.
- [5] Wu Kuo-lung, Yu Jian, and Yang Miin-Shen. A novel fuzzy clustering algorithm based on a fuzzy scatter matrix with optimality test [J]. *Pattern Recognition Letters*, 2005, 26(5): 639-652.
- [6] Yoshiki Fukuyama and Michio Sugeno. A new method of choosing the number of clusters for the fuzzy c-means method [C]. Proceedings of the 5th Fuzzy System Symposium, in Japanese, 1989: 247-250.
- [7] Clausi David. K-means Iterative Fisher (KIF) unsupervised clustering algorithm applied image texture segmentation [J]. *Pattern Recognition*, 2002, 35(9): 1959-1972.
- [8] 曹苏群, 王士同, 陈晓峰等. 基于模糊 Fisher 准则的半模糊聚类算法[J]. 电子与信息学报, 2008, 30(9): 2162-2165.  
Cao Su-qun, Wang Shi-tong, and Chen Xiao-feng, *et al.* Fuzzy fisher criterion based semi-fuzzy clustering algorithm[J]. *Journal of Electronics & Information Technology*, 2008, 30(9): 2162-2165.
- [9] Blake C L and Merz C J. UCI repository of machine learning databases, Irvine. CA: University of California, Department of Information and Computer Science, <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998, 7.
- [10] 王玲, 薄列峰, 焦李成. 密度敏感的谱聚类[J]. 电子学报, 2007, 35(8): 1577-1581.  
Wang Ling, Bo Lie-feng, and Jiao Li-cheng. Density-sensitive spectral clustering [J]. *Acta Electronica Sinica*, 2007, 35(8): 1577-1581.

支晓斌: 男, 1976 年生, 讲师, 研究方向为模式识别、模糊集理论及其应用。

范九伦: 男, 1964 年生, 教授, 博士生导师, 研究方向为模糊集理论、模式识别与图像处理、智能信息处理。