

基于聚类优化的非负矩阵分解方法及其应用

栗茂林¹ 梁霖² 陈元明² 徐光华² 何康康²

1.西安交通大学工程坊,西安,710049

2.西安交通大学机械工程学院,西安,710049

摘要:针对不断增加的机电系统运行状态信息,传统的特征提取和选择方法已无法满足需求。根据非负矩阵分解典型算法的特点,基于非负矩阵分解的聚类特性,提出了一种面向故障诊断的分解方法。通过分类能力和迭代效率的对比分析,选择了相关性约束和稀疏性约束的改进型交替最小二乘迭代算法,确定了低维嵌入维数及迭代初始化方法,在UCI测试数据集和TEP系统的特征选择应用中验证了该方法的有效性。

关键词:非负矩阵分解;聚类;迭代算法;特征选择

中图分类号:TH16

DOI:10.3969/j.issn.1004-132X.2018.06.013

开放科学(资源服务)标识码(OSID):



Non-negative Matrix Factorization Based on Clustering and Its Application

LI Maolin¹ LIANG Lin² CHEN Yuanming² XU Guanghua² HE Kangkang²

1.Engineering Workshop, Xi'an Jiaotong University, Xi'an, 710049

2.School of Mechanical Engineering, Xi'an Jiaotong University, Xi'an, 710049

Abstract: With the increasing complexity of electromechanical system state informations, traditional feature extraction and selection methods were unable to meet the needs. According to the characteristics of conventional non-negative matrix factorization (NMF) algorithm, a NMF method for monitoring and fault diagnosis was proposed based on the clustering property of NMF. By comparing classification accuracy and iteration efficiency, an improved alternating least square iterative algorithm with sparsity and correlation constraints was selected, and the low-dimensional embedded dimension and iterative initialization method were also determined. Experimental results to UCI test datasets and fault diagnosis of Tennessee-Eastman process (TEP) systems show that this approach is more effective to extract the fault features, and enhance the failure pattern capabilities.

Key words: non-negative matrix factorization (NMF); clustering; iterative algorithm; feature selection

0 引言

随着信息获取技术的不断进步,描述系统状态的信息量越来越大,导致特征空间维数不断增加,从而引发维数灾难等问题。在传统的特征约简方法中,主分量分析、独立成分分析以及矢量量化等方法要求信号平稳、满足高斯条件,限制了应用范围。

作为一种新兴的多元数据处理方法,非负矩阵分解(non-negative matrix factorization, NMF)可在高维空间中获得原始数据的局部特征,其纯加性的表达方式符合“局部构成整体”的认知规律,成为信号处理、模式识别等领域的热门工具。张培林等^[1]采用NMF技术对发动机故障信号进行特征参数提取,获得了更高的分类精度。李兵

等^[2]采用二维NMF技术来提取时频分布矩阵特征参数,得到了较好的诊断效果。除了特征提取外,NMF在复杂工业过程和机电设备的监测诊断中也有应用^[3-4]。

根据NMF的问题模型,可将其求解归结为一个优化问题,即通过目标函数来刻画它的逼近程度。实际应用中,一般采用交替迭代方法获得局部最优解。常用的乘性迭代算法、梯度下降算法和交替最小二乘算法是基于不同思路提出的,而对于设备诊断的特征提取,需要选择出合适的分解模型。

研究表明,NMF进行数据约简的目的是估计出原数据中的结构,而其分解的基向量与K均值聚类中的“类”有相通之处,即具有软聚类特性。因此,本文从分类性能和迭代效率角度出发,基于NMF的聚类特性,将交替非负最小二乘算法用

收稿日期:2017-01-05

基金项目:国家自然科学基金资助项目(51575438)

于故障诊断,并通过基矩阵的聚类性优化出约束参数与嵌入维数,最后,通过测试数据和特征选择实例应用验证了其有效性。

1 非负矩阵分解的迭代准则

1.1 NMF 原理

NMF 是一种非负性约束下的矩阵分解方法,具体可描述如下:给定一个非负矩阵 $V \in \mathbf{R}^{m \times n}$ (m 为样本数, n 为特征个数),将矩阵 V 分解成非负的基矩阵 $W \in \mathbf{R}^{m \times r}$ 和系数矩阵 $H \in \mathbf{R}^{r \times n}$ (通常情况下,要求 $r < n$) 的乘积,使得 $V = WH$ 成立。

非负矩阵分解可以通过优化迭代来解决,即通过目标函数来刻画逼近程度,并在非负约束条件下进行迭代求解。距离目标函数中,基于欧氏距离目标函数的优化问题应用较广泛^[5],其表达式如下:

$$\min_{W, H} \frac{1}{2} \|V - WH\|_F^2 \quad (1)$$

式中, $\|\cdot\|_F$ 表示 2 范数。

1.2 迭代准则

在目标函数(式(1))中,当 W 和 H 同时为变量时,求解最小化问题是非凸的,因此采用交替迭代更新 W 和 H 以获得局部解。常用的迭代算法有乘性迭代(multiplicative iterative, MI)、Lin's 投影梯度(Lin's projected gradient, LPG)及交替最小二乘法(alternating least squares, ALS)。其中,MI 算法是针对欧氏距离目标函数的优化问题提出的;LPG 算法用投影梯度法进行更新迭代,采用 Armijo 规则来搜索每次迭代步长;ALS 算法根据库恩塔克一阶最优性条件,采用直接估计“驻点”的固定点方法获得局部解^[6]。在这三种典型迭代算法中,ALS 算法形式更简单,理论上分解结果优于 MI 算法,在处理高维数据时的效率高于 LPG 算法,缺点是对噪声敏感,但可以通过增加约束条件来提高优化效果。标准 ALS 算法的更新规则为

$$\left. \begin{aligned} W &\leftarrow [VH^T(HH^T)^{-1}]_+ \\ H &\leftarrow [(W^T W)^{-1}W^T V]_+ \end{aligned} \right\} \quad (2)$$

式中, $[\cdot]_+$ 表示非负矩阵。

2 面向聚类的非负矩阵分解算法

2.1 NMF 的聚类特性

NMF 数据约简的目的就是估计出原始数据的本质结构,即获得基矩阵 W 。事实上,基矩阵 W 也可理解为模式识别理论中的“模式”,基向量空间就形成了 K 均值聚类中的“类”^[7]。因此,NMF 就是要挖掘出原始特征中的本质结构,发现其中

的“类”,将相似的局部特征聚集为一类,然后用提取出的“类”代表原始特征集。

NMF 分解过程中没有正交性约束的要求,即具有软聚类性,有助于克服一些硬聚类所遇到的问题。同时,NMF 不受样本数据的限制,无需预先确定聚类数,有效提高了其应用能力。

2.2 施加约束的改进型交替最小二乘法

为了获得更稳定、优质的分解结果,需要对 W 和 H 施加其他约束条件,可同时对 W 和 H 或只对其中一个施加约束。从故障诊断的角度来说,希望基矩阵 W 的聚类效果尽可能好,相关性尽可能小,也希望权矩阵 H 具有一定的稀疏度。因此,对基矩阵 W 施加相关性约束,对权矩阵 H 施加稀疏性约束,最优化问题转化为^[8]

$$\min(\frac{1}{2} \|V - WH\|_F^2 + \alpha_w J_w(W) + \alpha_H J_H(H)) \quad (3)$$

式中, α_w 为相关系数, $\alpha_w \geq 0$; $J_w(W)$ 为 W 的迹, $J_w(W) = \text{tr}(WBW^T)$; 矩阵 B 为 k 阶全 1 矩阵; α_H 为稀疏化系数; $J_H(H)$ 为 H 的迹, $J_H(H) = \text{tr}(HBH^T)$ 。

约束条件下,改进型 ALS 的迭代更新准则为

$$W \leftarrow [VH^T(HH^T + \alpha_w B)^{-1}]_+ \quad (4)$$

$$H \leftarrow [(W^T W)^{-1}(W^T V - \alpha_H B)]_+ \quad (5)$$

诊断应用中,约束条件可根据情况进行选择,可同时施加两个约束,也可只施加其一。约束参数的选取可通过基矩阵 W 的聚类效果进行确定。

2.3 初始化方法的选择

迭代更新方式对 W 、 H 的初始值很敏感,常用的初始化方法^[9-11]中,随机初始化法、随机 Acot 初始化法以及模糊 C 均值初始化法的初始值不唯一,分解结果不稳定,而基于奇异值分解的初始化方法可以获得唯一的初始值,使得分解结果稳定,具有可复现性。

2.4 嵌入维数的选择

实际应用中,往往很难获得先验知识,因此无法预先确定嵌入维数 r 。 r 过大或过小,都会直接影响分解结果中基矩阵结构表达的准确性。因此,在处理高维特征空间时,在经验值 $C \sim \sqrt{n}$ (C 为样本类别数)附近进行筛选,即利用 KNN 分类算法对分解的基矩阵 W 进行统计,根据基矩阵的聚类效果来选择最优 r 值。虽然需要计算不同 r 下的聚类性能,但可以尽量减小因 r 选择不佳对算法造成的影响。

3 实验分析

3.1 测试数据集

选取 4 个常用的 UCI 数据集(电离层数据集、

垃圾邮件分类数据集、钢板缺陷数据集以及肺癌数据集)对迭代算法的分解结果进行比较分析,详情见表 1。

表 1 测试数据集的描述

Tab.1 Description of test datasets

数据集	样本数	特征数	数据类型	类别数	嵌入维数
电离层数据集	351	34	连续值	2	5
垃圾邮件分类数据集	4 601	57	连续值	2	4
钢板缺陷数据集	1 941	27	连续值	7	4
肺癌数据集	32	56	离散值	3	3

3.2 迭代算法的效果对比

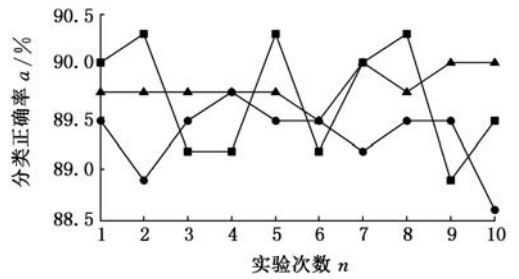
为了验证标准 ALS 的效果,将其与 MI 和 LPG 算法进行对比,其中,嵌入维数 r 的选取见表 1。由于随机初始化的多次处理有可能包含最优解,因此,3 种算法均采用随机初值,统计 10 次的迭代结果。图 1 所示为不同迭代算法在 4 个数据集上应用的分类正确率。

由图 1 可知,MI 算法的分类正确率总体上落后于其他 2 种算法,仅仅是针对肺癌数据集的分类率优于 LPG 算法;MI 算法不稳定且易收敛至局部点,如图 1a 中第 10 次和图 1c 中第 3 次实验结果,分类正确率都处于最低点。LPG 算法在三者中的稳定性最好,波动较小。标准 ALS 算法可以获得比 MI 算法和 LPG 算法更优的解,如图 1a 中第 2 次和图 1d 中的所有结果,但波动性比较大,主要是对噪声比较敏感。此外,即使是同一种迭代方法,每次实验结果都有较大的差距,说明随机初始值对 NMF 影响较大。标准 ALS 算法在分类性能方面表现良好,而通过施加约束条件还有提高分解效果的可能。另外,从数据集的分解效果来看,垃圾邮件分类数据的 3 种分解效果基本相同,这与垃圾邮件分类较明确有关。

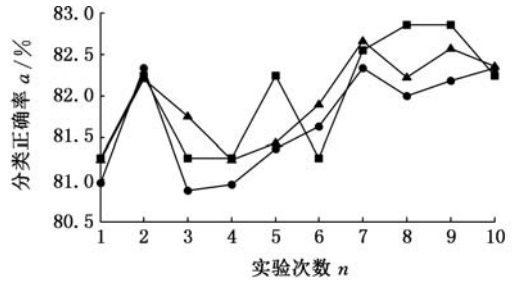
3.3 初始化方法的影响

采用电离层和垃圾邮件分类数据集对初始化方法进行测试,目标函数为欧氏距离,迭代规则为标准 ALS,嵌入维数 r 分别为 5 和 4,评判标准不变。不同的初始化方法得到的基矩阵 W 的分类结果如图 2 所示。

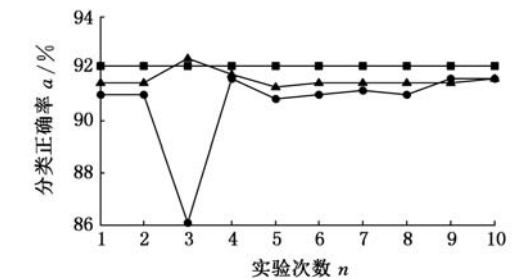
对于电离层和垃圾邮件分类两种数据集来说,最优和最差的分解结果皆由随机法和随机 Acol 法获得,这两种方法的波动性也最大;模糊 C 均值法与奇异值分解法虽然最终分类率不如上述两种方法,但其波动性很小,且与最优值差距不大。这四种方法当中,奇异值分解法的结果最稳定,平均分类率也较高,且具有可重复性,是最为理想的初始化方法。



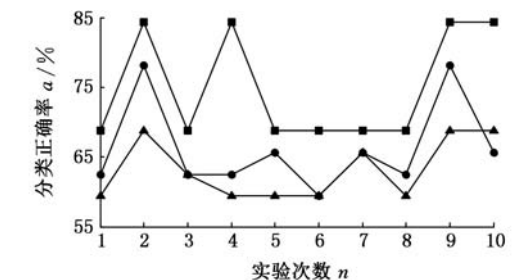
(a) 电离层数据集的分类正确率



(b) 垃圾邮件分类数据集的分类正确率



(c) 钢板缺陷数据集的分类正确率



(d) 肺癌数据集的分类正确率

—●— MI —▲— LPG —■— ALS

图 1 三种迭代算法的分类正确率

Fig.1 Accuracy of three iterative algorithms

3.4 约束条件的影响

相关性约束可以减少基矩阵 W 中的冗余结构,有利于控制基向量的冗余性。图 3 所示为相关性约束对电离层数据集的影响结果,其中,初始化和评判标准不变。显然,施加约束下的分类正确率略有提高。当分类结果较差(图 3 中第 10 次实验)时,施加相关性约束后的效果较为明显。当分类结果已较为理想(图 3 中第 7 次实验)时,约束效果反而变差,究其原因,主要是当已较为准确地提取出数据的本质结构时,施加约束反而无法迭代出更好的结果,为此,在后续应用中不施加相关性约束条件。

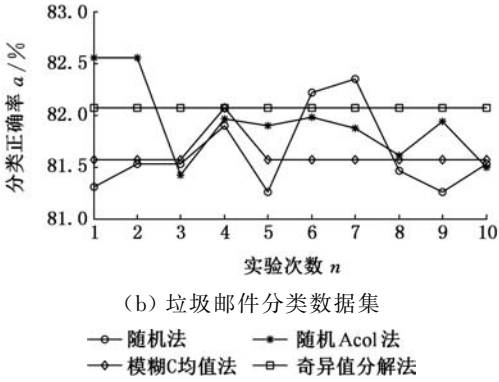
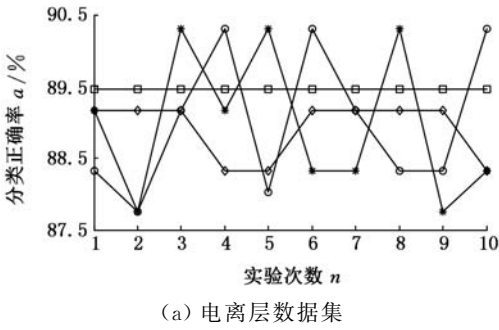


图2 两种数据集的初始化方法对比结果
Fig.2 Comparison with the initial methods for two datasets

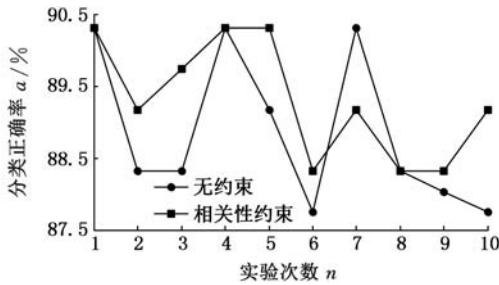


图3 相关性约束对电离层数据集的影响
Fig.3 Effects of the Ionosphere dataset with the correlation constraints

特征约简的重点是获取稀疏的权矩阵,即希望用极少维度的向量表征数据的突出特征。因此,对钢板缺陷数据集进行稀疏性约束,其实验结果如图4所示。由图4可知,在10次随机初始优化搜索中,在施加稀疏性约束的情况下,NMF分解得到基矩阵的聚类效果基本都得到了改善,即分类正确率得到了提高。这表明了稀疏约束对提高特征约简效果的有效性。

稀疏度约束能够提高分解质量,但稀疏度系数的不同对分类效果也有一定的影响。图5所示为电离层数据集在不同稀疏度系数下的分类正确率变化情况,其中,嵌入维数 $r=5$,稀疏度系数在 $0 \sim 0.5$ 内变化,KNN分类器的邻域为5。曲线变化情况表明,当 $\alpha_H=0.4$ 时,分解得到的基矩阵 W 的聚类效果最好,对应着较大的分类正确率。参数选择不佳,则导致约束下的分类效果变差,因此

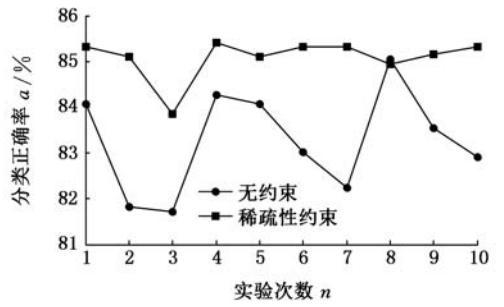


图4 稀疏性约束对钢板缺陷数据集的影响
Fig.4 Effects of the Steel plates faults dataset with the sparsity constraints

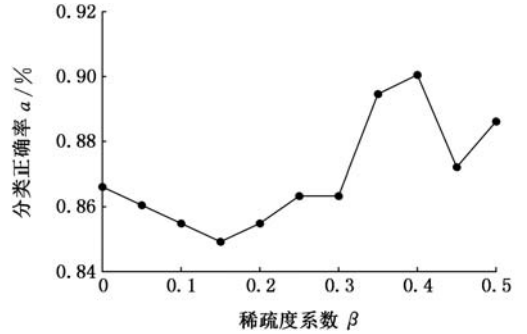


图5 不同稀疏度参数对电离层数据集的影响
Fig.5 Effects of the Ionosphere dataset with the different sparsitys

在应用中需要根据分类情况来优选稀疏度系数。

4 在特征选择中的应用

4.1 实例对象数据

采用由伊斯曼化学公司创建的田纳西-伊斯曼过程(Tennessee-Eastman process,TEP)进行实例分析。TEP过程包含了1种正常状态IDV(0)和21种故障状态(IDV(1)~IDV(21)),涉及22个过程测量变量(XMEAS(1)~XMEAS(22))、19个成分测量值(XMEAS(23)~XMEAS(41))和12个控制变量(XMV(1)~XMV(12))。22个过程测量变量包括进料、反应器压力以及汽提器的流量等信息;19个成分测量值为反应过程中各种气体成分测量值,都包含高斯噪声;12个控制变量描述了进料量、分离器罐液流量以及冷却水流量等信息。其中,搅拌速度变量恒定不变,在应用中不包含该项,即一共52个观测变量。

从TEP中选择IDV(2)、IDV(3)、IDV(4)、IDV(5)四类故障数据(每类故障480个样本)的52个特征进行特征子集的选择。因此,待分解数据为 1920×52 的高维矩阵 V 。

4.2 约束参数和嵌入维数的选择

特征选择的目的是从原始52维特征中选出对故障分类敏感的特征。因此,首先在原始特征

中,通过特征间的相关系数剔除掉矩阵 V 中相关性较高的冗余特征,将矩阵特征维数从 52 维降到 16 维,则待分解矩阵 X 为 1920×16 维矩阵。其次,评估不同参数下分解矩阵 X 得到的基矩阵 W 的分类性能,将其中分类性能好的基矩阵作为特征选择的基础。在评估中,低维嵌入维数 r 在 $3 \sim 5$ 内变化,稀疏度约束 α_H 以 0.05 的步长在 $0 \sim 0.5$ 内变化,通过分解基矩阵的聚类情况确定出优化的维数和稀疏度。

图 6 所示为不同参数下基矩阵的分类率曲线,其中,KNN 分类器的近邻数为 5。由图 6 可知,在 $r=5, \alpha_H=0.5$ 时分解得到的基矩阵 W 的分类性能最好。根据各基向量在分类性能上的互补性和组合性,剔除掉基矩阵 W 中的第 2 和第 4 冗余基向量,保留第 1、3、5 基向量,使得彼此间具有结构上的互补性,组合在一起则具有良好的分类性,其投影空间中的样本分布如图 7 所示。

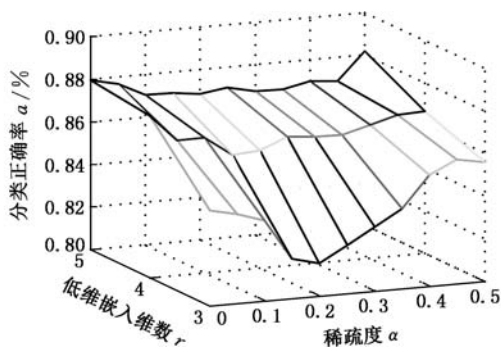


图 6 不同参数下基矩阵的分类正确率

Fig.6 Accuracy curve of basis matrix with different parameters

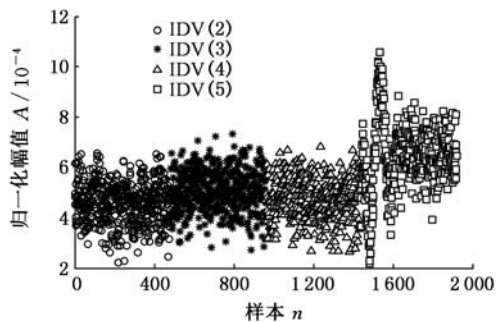
4.3 特征选择的效果

根据优化参数确定的权矩阵 H 分布能够进行特征选择,即在最优的基矩阵 W 对应的系数矩阵 H 中找出对应行幅值最大的元素,该元素所在的列即为有效特征。

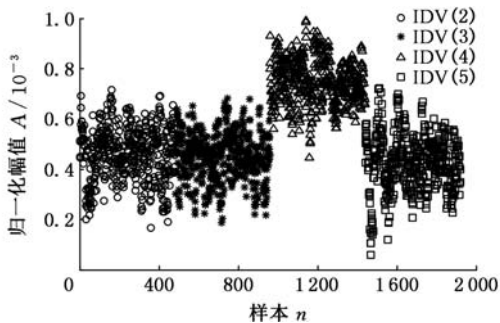
在优化参数确定的系数矩阵 H 中,找出对应行幅值最大的元素,该元素所在的列依次为特征 16、15 和 2,对应原始集合特征子集为 $\{52, 51, 10\}$ 。其中,特征 52、51 和 10 分别为控制变量 XMV(11)、XMV(10) 和过程测量变量 XMEAS(10)。该子集的分类率达到 89.74%。

图 8 所示为 $\{52, 51, 10\}$ 所张成的三维空间投影,由样本分布可知,除 IDV(3) 故障样本 (“*”) 和 IDV(5) 故障样本 (“□”) 之间有一定的重合外,其他故障样本间都区分得较为明显,显然是一个较理想的特征子集。

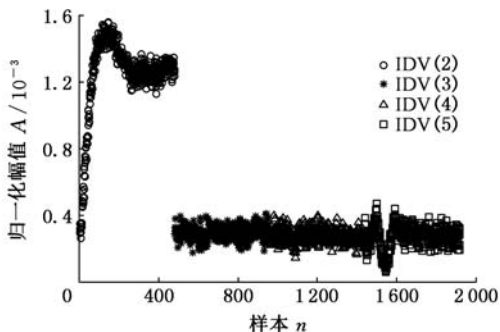
若不考虑稀疏度的约束条件,在 $r=5$ 的情况



(a) 第 1 基向量空间中的样本分布



(b) 第 3 基向量空间中的样本分布



(c) 第 5 基向量空间中的样本分布

图 7 三个基向量空间中的样本分布

Fig.7 Distribution of three projection vectors

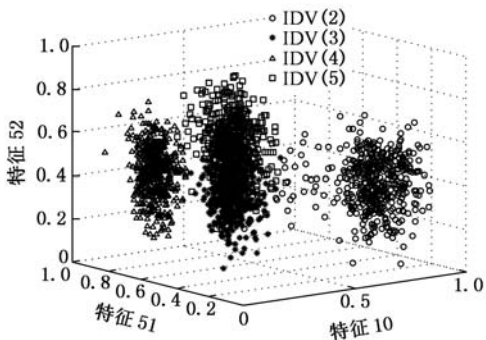


图 8 子集 $\{52, 51, 10\}$ 的三维样本分布

Fig.8 Distribution of the feature subset $\{52, 51, 10\}$

下,稀疏度和相关系数均取零,选出的特征依次为 12、14 和 15,对应原始集合中的分类特征子集为 $\{41, 51, 52\}$,该子集的分类率只能达到 81.87%。由此可见,不考虑约束条件的话,识别效果会有一定程度的下降。

基于主成分分析的特征选择结果为 $\{10, 28, 47\}$,样本分布的正确率仅为 65.62%,显然提取的

特征子集更不理想。

5 结语

结合 NMF 的聚类性能,在 NMF 的迭代求解优化问题基础上,结合相关性约束和稀疏性约束,提出了一种面向 NMF 聚类性的迭代优化算法,通过基矩阵空间中的分类性能选择出具有最优聚类性的嵌入维数及相关参数。通过测试数据验证了最小二乘迭代算法的有效性以及约束条件的效果,通过复杂机电系统的应用实例验证了特征选择的有效性。

参考文献:

[1] 张培林,王怀光,张磊,等.非负矩阵分解在发动机故障特征提取中的应用[J].振动工程学报,2013,26(6):944-950.
ZHANG Peilin, WANG Huaiguang, ZHANG Lei, et al. Feature Extraction for Engine Fault Diagnosis by Utilizing Adaptive Multi-scale Morphological Gradient and Non-negative Matrix Factorization[J]. Journal of Vibration Engineering, 2013, 26(6): 944-950.

[2] 李兵,米双山,刘鹏远,等.二维非负矩阵分解在齿轮故障诊断中的应用[J].振动测试与诊断,2012,32(5):836-840.
LI Bing, MI Shuangshan, LIU Pengyuan, et al. Application of Two-dimensional Non-negative Factorization for Gear Fault Diagnosis[J]. Journal of Vibration, Measurement & Diagnosis, 2012, 32(5): 836-840.

[3] 陈霖,邹金慧.基于 NMF-SVM 的复杂化工过程故障诊断[J].计算机与应用化学,2014,31(8):1015-1018.
CHEN Lin, ZOU Jinhui. Fault Diagnosis of Complex Chemical Process Based on NMF-SVM [J]. Computers and Applied Chemistry, 2014, 31(8): 1015-1018.

[4] 唐曦凌,梁霖,高慧中,等.结合连续小波变换和多约束非负矩阵分解的故障特征提取方法[J].振动与冲击,2013,32(19):7-11.
TANG Xiling, LIANG Lin, GAO Huizhong, et al. Fault Feature Extraction Method Combining Con-

tinuous Wavelet Transformation with Multi-constraint Nonnegative Matrix Factorization[J]. Journal of Vibration and Shock, 2013, 32(19): 7-11.

[5] WANG Yuxiong, ZHANG Yujin. Nonnegative Matrix Factorization: a Comprehensive Review [J]. IEEE Transactions on Knowledge and Data Engineering, 2013, 25(6): 1336-1353.

[6] BERRY M W, BROWNE M, LANGVILLE A N, et al. Algorithms and Applications for Approximate Nonnegative Matrix Factorization[J]. Computational Statistics & Data Analysis, 2007, 52(1): 155-173.

[7] LI T, DING C. The Relationships among Various Nonnegative Matrix Factorization Methods for Clustering[C]//Sixth International Conference on Data Mining. Hong Kong, 2006: 362-371.

[8] CICHOCKI A, ZDUNEK R, PHAN A H, et al. Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation [M]. Hoboken: Wiley, 2009: 203-213.

[9] YU Shaohui, ZHANG Yujun, LIU Wenqing, et al. A Novel Initialization Method for Nonnegative Matrix Factorization and Its Application in Component Recognition with Three-dimensional Fluorescence Spectra[J]. Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy, 2012, 86: 315-319.

[10] XUE Yun, TONG C S, CHEN Ying, et al. Clustering-based Initialization for Non-negative Matrix Factorization[J]. Applied Mathematics and Computation, 2008, 205(2): 525-536.

[11] JANECEK A, TAN Y. Using Population Based Algorithms for Initializing Nonnegative Matrix Factorization[J]. Lecture Notes in Computer Science, 2011, 6729: 307-316.

(编辑 张 洋)

作者简介:栗茂林,女,1978年生,讲师、博士。研究方向为智能信息处理技术及故障诊断。发表论文20余篇。E-mail: maolinli@xjtu.edu.cn。梁霖(通信作者),男,1973年生,副教授、博士生导师。研究方向为机械设备故障诊断技术。E-mail: lianglin@xjtu.edu.cn。