

基于 W2DPCA-FCM 的近红外显微图像分割

杨秀坤¹ 钟明亮^{1,2} 景晓军² 岳新启¹

(¹ 哈尔滨工程大学信息与通信工程学院, 黑龙江 哈尔滨 150001)
² 北京邮电大学信息与通信工程学院, 北京 100876)

摘要 特征提取与聚类分析相结合的图像分割方法可以用于近红外显微图像化学信息的快速提取。针对基于主成分分析(PCA)的特征提取运算较为复杂的缺点,提出了一种加权二维主成分分析(W2DPCA)光谱特征提取方法,与模糊 C 均值(FCM)算法相结合用于近红外显微图像化学分布信息提取。通过片剂的近红外显微图像的仿真实验,验证了 W2DPCA-FCM 方法的可行性和有效性。实验结果表明,W2DPCA-FCM 方法可以减少计算时间、提高聚类精度,是一种有效的红外显微图像分析方法。

关键词 成像系统;近红外显微成像;加权二维主成分分析;模糊 C 均值;图像分割

中图分类号 O657.33; TP391.41 **文献标识码** A **doi:** 10.3788/AOS201333.0811002

Near-Infrared Microscopic Image Segmentation Based on W2DPCA-FCM

Yang Xiukun¹ Zhong Mingliang^{1,2} Jing Xiaojun² Yue Xinqi¹

¹ College of Information and Communication Engineering, Harbin Engineering University,
Harbin, Heilongjiang 150001, China

² School of Information and Communication Engineering, Beijing University of Posts and Telecommunications,
Beijing 100876, China

Abstract Segmentation of near-infrared (NIR) microscopic image by feature extraction and clustering analysis methods can be used for efficient extraction of chemical information. Due to the high computational complexity of principal component analysis (PCA) in extracting features, we propose a weighted two-dimensional PCA (W2DPCA) spectral feature extraction scheme in this paper, which is combined with fuzzy C-mean (FCM) algorithm to extract the chemical information of NIR microscopic image. The feasibility and effectiveness of the proposed algorithm are verified by simulation experiments performed on NIR microscopic image of tablets. Experimental results show that W2DPCA-FCM is an effective infrared microscopy image analysis method since it can reduce the computation time and improve the clustering accuracy.

Key words imaging systems; near-infrared microscopic imaging; weighed two-dimensional principal component analysis; fuzzy C-mean; image segmentation

OCIS codes 110.2960; 100.5010; 300.6340

1 引 言

红外显微成像技术是一种将红外光谱技术和显微技术相结合的快速、直接、绿色的微区分析技术,它能够在不破坏样品原始结构的前提下深入微观视野,探测生物组织、高分子聚合物等样品的化学组分,具有较高的空间分辨率和光谱分辨率。其中近

红外(NIR)显微成像技术能够无损地提供足够的空间分布信息,获得片剂表面重要的参数^[1],可以用于辨别药品真伪^[2]、检测成分均匀性^[3]以及提取成分分布信息^[4]等,因而在制药领域得到了广泛的应用。

近年来人们将红外显微成像技术与多元分析技术[主成分分析(PCA)^[5-6]、多元曲线分辨-交替最

收稿日期: 2013-02-27; 收到修改稿日期: 2013-03-26

基金项目: 国家自然科学基金(61143008)、黑龙江省自然科学基金(ZD200915)

作者简介: 杨秀坤(1971—),女,博士,教授,主要从事图像处理、模式识别和多光谱检测等方面的研究。

E-mail: yangxiukun@hrbeu.edu.cn

本文电子版彩色效果请详见中国光学期刊网 www.opticsjournal.net

小二乘法(MCR-ALS)^[7-8]、聚类分析^[9]等]相结合,用于红外显微图像化学信息提取。其中聚类分析可以按像素光谱到中心光谱距离的不同将红外显微图像划分为不同区域,使得同一区域的光谱尽可能相似,不同区域的光谱尽可能不同,从而实现红外显微图像化学分布信息的提取。但直接对原始光谱数据进行聚类不仅计算量较大,而且容易受到噪声的干扰^[10],因此,需要对原始光谱数据进行特征提取,以降低数据维度,提高聚类精度。PCA是一种常用的特征提取方法,它无需先验知识即可降低数据维数,揭示隐藏复杂数据背后的简单结构,并且可以用少量的载荷光谱以及得分图表示原始光谱数据,然而PCA需要先将原始二维(2D)矩阵转换为一维向量才能构造协方差矩阵,增加了计算复杂度。因此本文提出了一种基于加权二维主成分分析-模糊C均值(W2DPCA-FCM)的近红外显微图像分割方法,它可以通过列(或行)光谱矩阵来构造广义协方差矩阵,实现光谱数据特征的快速提取和噪声去除;可以通过对输入模糊C均值分析的特征向量进行加权,突出不同特征向量对聚类的不同贡献,提高近红外显微图像的分割精度。

2 基于 W2DPCA-FCM 的近红外显微图像分割

2.1 基于 W2DPCA 的特征提取

基于谱方向的二维 PCA 可以利用列(或行)光谱矩阵来构造广义协方差矩阵,使构造的协方差矩阵包含不同通道分布图之间的相关性^[11]。设三维的红外显微图像为 $\mathbf{B}_{m \times n \times k}$,它包含有 k 个通道,每个通道的分布图的大小为 $m \times n$,定义行光谱矩阵构造的协方差矩阵为

$$\mathbf{C}_l = \frac{1}{m} \sum_{i=1}^m (\mathbf{A}_i - \bar{\mathbf{A}})^T (\mathbf{A}_i - \bar{\mathbf{A}}), \quad (1)$$

式中 \mathbf{A}_i 为 $n \times k$ 的行光谱, $\bar{\mathbf{A}} = \frac{1}{m} \sum_{i=1}^m \mathbf{A}_i$ 为均值行光谱。对 \mathbf{C}_l 进行特征分解,选择前 l 个较大的特征值 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_l$ 及其对应的特征向量 $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l$, 构成主分量矩阵 $\mathbf{W}_{k \times l}$ 。将红外显微图像 $\mathbf{B}_{m \times n \times k}$ 投影到特征矩阵 $\mathbf{W}_{k \times l}$ 上,得到三维投影特征矩阵 $\mathbf{Z}_{m \times n \times l}$ 。得分图由 $m \times n$ 维投影特征矩阵 $\mathbf{Z}_j (j=1, 2, \dots, l)$ 构成,可以用来表示物质的分布情况;载荷光谱由主分量矩阵 \mathbf{W}^T 的行向量构成,它与得分图代表物质的光谱具有相近的特征峰峰位和峰形。由于不同的投影特征向量对聚类的贡献是不同的,通

常特征值大的特征向量包含的信息较多,对聚类贡献较大,因此利用特征值对投影特征向量进行加权,并通过调整因子来获得适当的加权值,以突出不同投影特征向量对聚类的贡献。加权投影特征矩阵为 $\mathbf{Z}' = [\mathbf{Z}'_1, \mathbf{Z}'_2, \dots, \mathbf{Z}'_l] = [\lambda_1^{\alpha} \mathbf{Z}_1, \lambda_2^{\alpha} \mathbf{Z}_2, \dots, \lambda_l^{\alpha} \mathbf{Z}_l] = [\lambda_1^{\alpha} \mathbf{B} \mathbf{w}_1, \lambda_2^{\alpha} \mathbf{B} \mathbf{w}_2, \dots, \lambda_l^{\alpha} \mathbf{B} \mathbf{w}_l] = [\mathbf{B}(\lambda_1^{\alpha} \mathbf{w}_1), \mathbf{B}(\lambda_2^{\alpha} \mathbf{w}_2), \dots, \mathbf{B}(\lambda_l^{\alpha} \mathbf{w}_l)] = \mathbf{B} \mathbf{W} \Lambda^{\alpha}$, (2) 其中 $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_l)$ 是前 l 个较大的特征值, $\alpha \in [0, 1]$ 为权值的调整因子。 α 的大小直接影响投影特征向量对聚类的贡献, $\alpha = 0$ 时,每个投影特征向量加权值相等; $\alpha = 1$ 时,直接以特征值作为加权值; α 过大或过小都得不到好的分割效果,故选择 $\alpha = 0.2$ 。将加权矩阵 \mathbf{Z}' 拉伸为二维矩阵 $\mathbf{M}_{m \times l}$ 即为 W2DPCA 提取的特征。

2.2 FCM 聚类算法

FCM 聚类算法是基于目标函数的非线性迭代优化方法,它将每个样本以不同的隶属度分配到某一类。设 $X = \{x_1, x_2, \dots, x_n\} \in R^m$ 为待聚类的数据集,FCM 通过最小化目标函数 J 来将数据集 X 分为 c 类:

$$J = \sum_{k=1}^c \sum_{i=1}^n \mu_{ik}^g \|x_i - m_k\|, \quad (3)$$

式中 m_k 为第 k 类的聚类中心, $\sum_{k=1}^c \mu_{ik} = 1, \forall i = 1, 2, \dots, n, \mu_{ik}$ 表示 x_i 对 m_k 的隶属度, g 为模糊加权指数。为了使目标函数 J 最小化,利用拉格朗日算法可以得到

$$\mu_{ik} = \sum_{k=1}^c \left(\frac{\|x_i^{(k)} - m_j\|}{\|x_i^{(k)} - m_k\|} \right)^{-2/(g-1)}, \quad (4)$$

$$m_k = \frac{\sum_{i=1}^n \mu_{ik}^g x_i}{\sum_{i=1}^n \mu_{ik}^g}. \quad (5)$$

FCM 算法的步骤如下。

步骤 1: 初始化各输入参数,给定聚类簇个数 c ,设定迭代停止阈值 ξ 及模糊加权指数 g ,初始化聚类中心 $m_k, k = 1, 2, \dots, c$,设置迭代计数器 $a = 1$;

步骤 2: 按(3)式和(4)式计算并更新规划矩阵和聚类中心;

步骤 3: 重复步骤 2,直到聚类中心变得稳定,算法停止并输出隶属度矩阵和聚类中心。

2.3 近红外显微图像分割

经过 W2DPCA 特征提取,可以用少量的得分代替全光谱来表示每个像素,将提取的特征输入

FCM 可以实现近红外显微图像化学分布信息的提取。具体步骤如下。

步骤 1: 采用 W2DPCA 提取近红外显微图像特征;

步骤 2: 采用 FCM 对降维后的数据进行聚类分析, 得到隶属度矩阵;

步骤 3: 将隶属度向量恢复为二维图像矩阵, 即可得到相应成分的分布信息;

步骤 4: 按隶属度对原始光谱数据进行加权平均即可得到相应的中心光谱。

3 结果与讨论

3.1 实验数据和仿真环境

实验数据由美国 Perkin Elmer 公司提供, 通过傅里叶变换近红外成像系统 Spotlight 400 采集得到的包

含活性成分(AI)、乳糖、碳酸钙(CaCO₃)三种成分的药片近红外显微图像, 成像面积为 4825 μm×4800 μm, 空间分辨为 25 μm×25 μm, 共有 37056 pixel。光谱分辨力 8 cm⁻¹, 波数范围 7800~4000 cm⁻¹, 共有 476 个通道。采用 PE 公司的 HyperView 软件进行一阶导数预处理和最小二乘拟合(LSF)分析, 其余的仿真实验主要基于 Matlab7. 11。所有程序均在 Celeron(R) Dual-Core CPU T3000 @ 1. 80 GHz、250 G 硬盘和 1. 98 GB 内存的某品牌笔记本电脑上运行。

3.2 特征提取

分别采用 PCA、2DPCA 方法对药片近红外显微图像进行光谱分解。图 1 和图 2 分别是 PCA、2DPCA 分解得到的前五帧得分图, 图中颜色由深到浅表示浓度由低到高。前五帧载荷光谱与三种成分参考光谱的相关系数如表 1 所示。

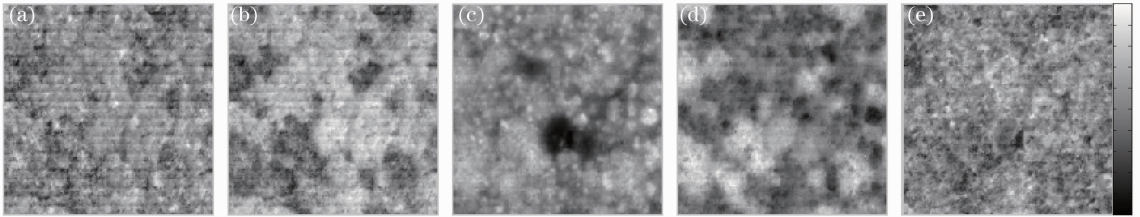


图 1 PCA 前五帧得分图。(a)第 1 帧; (b)第 2 帧; (c)第 3 帧; (d)第 4 帧; (e)第 5 帧

Fig. 1 Score images of PCA. (a) Score 1; (b) score 2; (c) score 3; (d) score 4; (e) score 5

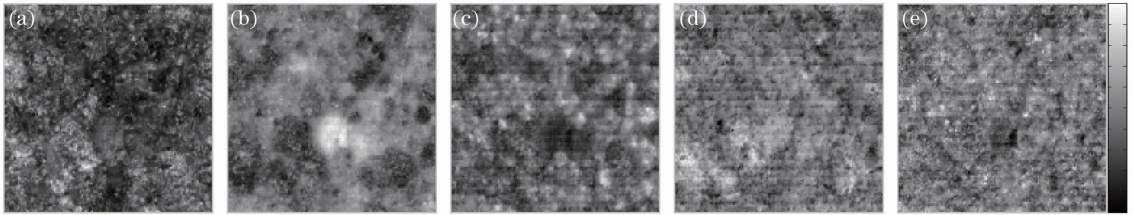


图 2 2DPCA 前五帧得分图。(a)第 1 帧; (b)第 2 帧; (c)第 3 帧; (d)第 4 帧; (e)第 5 帧

Fig. 2 Score images of 2DPCA. (a) Score 1; (b) score 2; (c) score 3; (d) score 4; (e) score 5

表 1 载荷光谱与参考光谱的相关系数

Table 1 Correlation coefficients between loading spectra and reference spectra

	PCA			2DPCA		
	Lactose	CaCO ₃	AI	Lactose	CaCO ₃	AI
Factor1	97. 31	92. 11	95. 80	-6. 82	27. 96	-13. 40
Factor2	80. 35	87. 97	83. 28	12. 19	52. 02	31. 24
Factor3	47. 96	-1. 33	23. 49	77. 48	44. 82	60. 12
Factor4	20. 52	36. 70	29. 86	54. 82	63. 54	59. 50
Factor5	-13. 42	-20. 95	-25. 13	-23. 21	-31. 66	-35. 81

由表 1 可以看出 PCA 分解得到的前四帧载荷光谱与成分参考光谱的相关性较高, 而 2DPCA 则是第 2~4 帧载荷光谱与成分参考光谱的相关性高, 因此选择 PCA 的前四帧得分图作为聚类的输入, 对 2DPCA 的第 2~4 帧得分图加权后作为聚类的

输入。

3.3 聚类分析

分别将原始光谱数据、PCA 以及 W2DPCA 提取的特征进行 FCM 聚类分析, 可以得到相应的成分分布图及估计浓度直方图。将三种成分分布图

合成 RGB 图像^[12]可以直观呈现三种成分的分布情况。RGB 合成图像的红色、绿色和蓝色通道与活性

成分、乳糖和碳酸钙相对应,结果如图 3 所示。

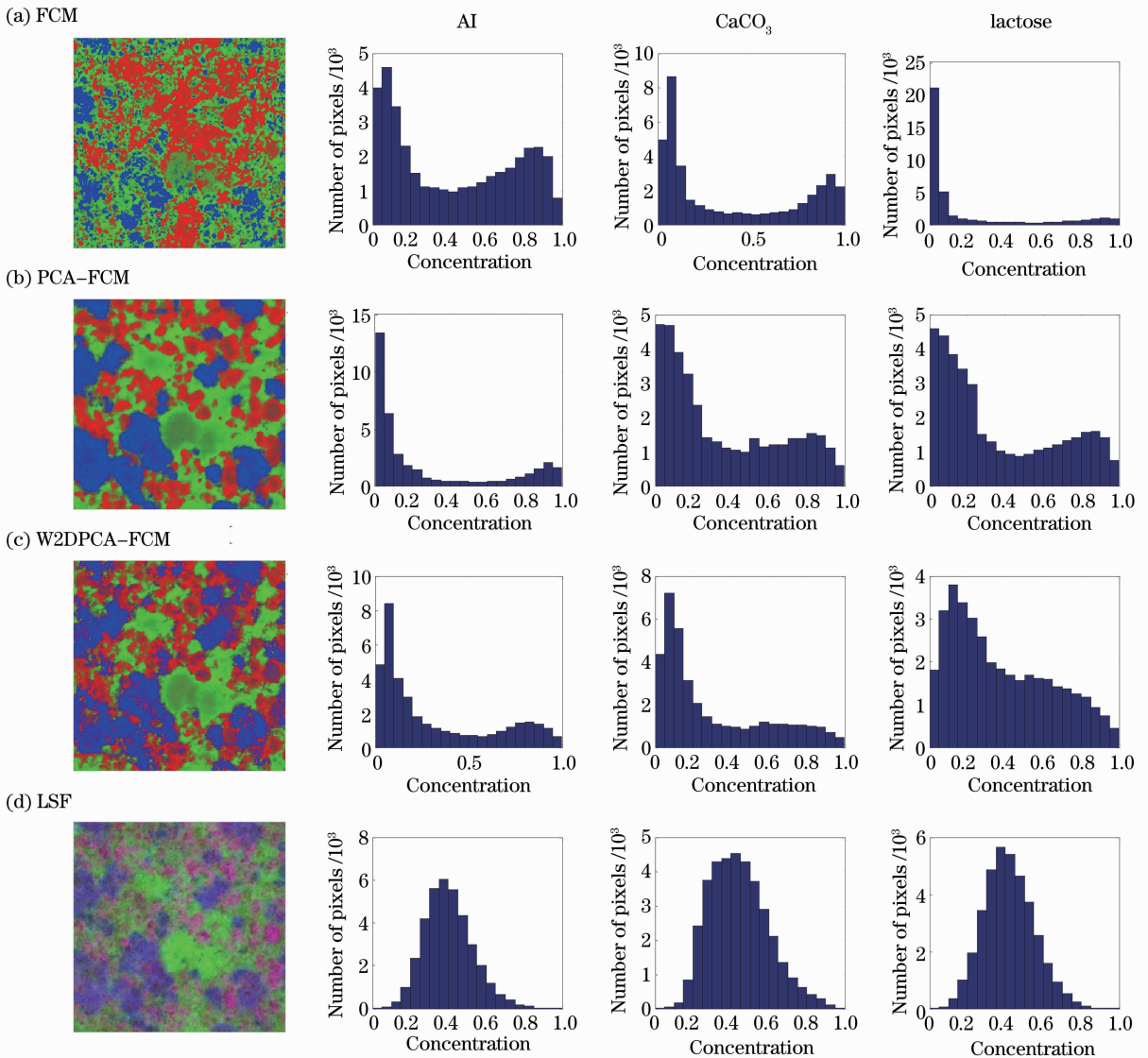


图 3 四种方法的 RGB 合成图以及估计浓度直方图

Fig. 3 RGB composite images and estimated concentration histograms of four algorithms

由于 LSF 分析基于参考光谱以及适当的预处理,分析的结果较为可靠,因而将其作为参考用于比较三种聚类方法。图 3(d)为采用 HyperView 软件对原始光谱数据进行一阶导数和 LSF 分析的结果。图 3(d)的成分分布图多为两种或三种颜色混合的像素,三种成分直方图的估计浓度(摩尔分数,下同)均集中于 0.3~0.4 之间,由此说明三种成分的实际分布均较为均匀。比较图 3(a)与图 3(d),活性成分的估计浓度直方图较为相似,但二者的成分分布图存在明显的差异;图 3(b)与图 3(c)的分布图较为相似,但图 3(c)的红色区域略大于图 3(b)的红色区域,与图 3(d)的红色区域的分布更为相似;图 3(b)

与图 3(d)的碳酸钙成分估计浓度直方图较为相似;图 3(c)与图 3(d)的活性成分和乳糖成分的估计浓度直方图较为相似。由以上分析可以判断,FCM 分割的结果较差,PCA-FCM 和 W2DPCA-FCM 分割的结果较好,且 W2DPCA-FCM 得到的成分分布图更接近于实际的成分分布情况。

通过计算三种方法得到的估计浓度、相对于 LSF(乳糖、碳酸钙和活性成分的浓度分别是 35.03%、31.76%和 33.21%)的浓度误差、中心光谱与对应参考光谱的相关系数、加权平均相关系数以及 20 次实验的平均运算时间,可以更好地进行性能比较,结果如表 2 所示。

表 2 性能比较
Table 2 Performance comparison

Method	Ingredient	Estimated concentration /%	Concentration error /%	Correlation coefficient /%	Average correlation coefficient /%	Time /s
FCM	Lactose	18.49	47.22	97.38	94.88	207.54
	CaCO ₃	39.02	22.86	92.67		
	AI	42.49	27.94	95.83		
PCA-FCM	Lactose	36.53	4.28	97.63	95.49	48.48
	CaCO ₃	35.94	13.16	93.63		
	AI	27.53	17.10	95.09		
W2DPCA-FCM	Lactose	38.22	9.11	97.53	95.68	8.15
	CaCO ₃	30.12	5.16	93.73		
	AI	31.66	4.67	95.31		

由表 2 中可以看出:PCA-FCM 算法估计的乳糖浓度与定量分析的结果最为接近,且对应中心光谱与乳糖相关性最高;W2DPCA-FCM 算法估计的碳酸钙以及活性成分的浓度与定量分析的结果最为接近,且碳酸钙中心光谱与碳酸钙参考光谱的相关性最高;FCM 算法估计活性成分浓度与定量分析的结果的误差最大,但对应的中心光谱与活性成分的相关性最高,这是因为在样本光谱存在干扰时,相关系数并不总是能够很好地区分光谱之间的差异性,相关系数最高不一定表示聚类的结果最好^[13]。因此从估计浓度上判断 W2DPCA-FCM 得到的活性成分的分布更符合药片活性成分的实际分布。三种算法中,W2DPCA-FCM 的加权平均相关系数最高,总浓度误差相对最小,由此判断 W2DPCA-FCM 算法的分割结果优于 FCM 和 PCA-FCM 的分割结果。

三种算法中,PCA-FCM 的运算时间小于 FCM 的运算时间,分割精度优于 FCM 的分割精度。W2DPCA-FCM 的运算时间远小于 FCM 和 PCA-FCM 的运算时间,分割精度优于 FCM 和 PCA-FCM 的分割精度。这是因为原始数据包含海量的光谱数据,经过 PCA 和 2DPCA 特征提取,可以有效压缩数据,减少运算复杂度,同时去除一定的噪声有利于提高分割精度,因此 PCA-FCM 和 W2DPCA-FCM 的运算时间小于 FCM 的运算时间,PCA-FCM 和 W2DPCA-FCM 的分割精度优于 FCM 的分割精度。PCA 的协方差矩阵包含的是任意两个光谱之间的相关信息,2DPCA 的协方差矩阵中包含的是不同通道分布图之间的相关性,由于实验数据的光谱数远大于通道数,因此 2DPCA 协方差的复杂度远小于 PCA 协方差复杂度,W2DPCA-FCM 运算时间要小于 PCA-FCM 的运算时间。

W2DPCA 对提取的特征进行加权能更好地突出不同特征对聚类的不同贡献,因此 W2DPCA-FCM 的分割精度优于 PCA-FCM 的分割精度。

4 结 论

将 W2DPCA 特征提取与 FCM 聚类方法相结合,提出一种基于 W2DPCA-FCM 的近红外显微图像分割方法。与 PCA 方法相比,W2DPCA 能够直接利用原始光谱数据构造的协方差矩阵,有利于提高运算效率;同时对提取的特征进行加权能更好地突出不同特征对聚类的不同贡献,有利于提高聚类精度。因此采用该方法可以快速有效地提取红外显微图像的化学信息。

参 考 文 献

- 1 A Palou, J Cruz, M Blanco, *et al.*. Determination of drug, excipients and coating distribution in pharmaceutical tablets using NIR-CI [J]. *J Pharmaceutical Analysis*, 2012, 2(2): 90-97.
- 2 M B Lopes, J C Wolff, J M Bioucas-Dias, *et al.*. Near-infrared hyperspectral unmixing based on a minimum volume criterion for fast and accurate chemometric characterization of counterfeit tablets [J]. *Anal Chem*, 2010, 82(4): 1462-1469.
- 3 J Cruz, M Blanco. Content uniformity studies in tablets by NIR-CI [J]. *J Pharmaceutical and Biomedical Analysis*, 2011, 56(2): 408-412.
- 4 J M Amigo, C Ravn. Direct quantification and distribution assessment of major and minor components in pharmaceutical tablets by NIR-chemical imaging [J]. *Eur J Pharmaceutical Sciences*, 2009, 37(2): 76-82.
- 5 S Serranti, D Cesare, F Marini, *et al.*. Classification of oat and groat kernels using NIR hyperspectral imaging [J]. *Talanta*, 2013, 103: 276-284.
- 6 Yang Xiukun, Zhong Mingliang, Jing Xiaojun, *et al.*. FTIR microscopic image analysis based on principal component analysis-2nd derivative spectral imaging [J]. *Acta Optica Sinica*, 2012, 32(7): 0711002.
- 7 杨秀坤, 钟明亮, 景晓军, 等. 基于主成分分析-二阶导数光谱成像的红外显微图像分析[J]. *光学学报*, 2012, 32(7): 0711002.
- 7 S Piqueras, L Duponchel, R Tauler, *et al.*. Resolution and

- segmentation of hyperspectral biomedical images by multivariate curve resolution-alternating least squares [J]. *Analytica Chimica Acta*, 2011, 705(1-2): 182–192.
- 8 B Vajna, A Farkas, H Pataki, *et al.*. Testing the performance of pure spectrum resolution from Raman hyperspectral images of differently manufactured pharmaceutical tablets [J]. *Analytica Chimica Acta*, 2012, 712: 45–55.
- 9 N Alajlan, Y Bazi, F Melgani, *et al.*. Fusion of supervised and unsupervised learning for improved classification of hyperspectral images [J]. *Information Sciences*, 2012, 217: 39–55.
- 10 J M Amigo, J Cruz, M Bautista, *et al.*. Study of pharmaceutical samples by NIR chemical-image and multivariate analysis [J]. *Trends in Analytical Chemistry*, 2008, 27(8): 696–713.
- 11 Yang Xiukun, Zhong Mingliang, Jing Xiaojun, *et al.*. FTIR microscopic imaging analysis based on 2DPCA [J]. *J Communications*, 2012, 33(9): 147–151.
杨秀坤, 钟明亮, 景晓军 等. 基于 2DPCA 的红外显微图像分析 [J]. *通信学报*, 2012, 33(9): 147–151.
- 12 L Zhang, M J Henson, S S Sekulic. Multivariate data analysis for Raman imaging of a model pharmaceutical tablet [J]. *Analytica Chimica Acta*, 2005, 545(2): 262–278.
- 13 P R Griffiths, L Shao. Self-weighted correlation coefficients and their application to measure spectral similarity [J]. *Appl Spectrosc*, 2009, 63(8): 916–919.

栏目编辑: 李文喆