

# 基于高光谱图像技术的苹果粉质化 LLE-SVM 分类

赵桂林, 朱启兵\*, 黄敏

江南大学通信与控制工程学院, 江苏 无锡 214122

**摘要** 苹果粉质化程度是衡量其内部品质的一个重要因素, 采用了高光谱散射图像技术进行苹果粉质化的无损检测。针对高光谱散射图像数据量大的特点, 提出了局部线性嵌入(local linear embedded, LLE)和支持向量机(support vector machine, SVM)相结合的用于检测苹果粉质化的新分类方法。LLE 是一种通过局部线性关系的联合来揭示全局非线性结构的非线性降维方法, 能有效计算高维输入数据在低维空间的嵌入流形。对降维后的高光谱数据采用 SVM 进行分类。将 LLE-SVM 分类方法与传统的 SVM 分类方法比较, 仿真结果表明, 对高光谱数据而言, 用 LLE-SVM 得到的训练精度高于单纯使用 SVM 的训练精度; 降维前后, 分类器的测试精度变化不大, 波动范围不超过 5%。LLE-SVM 为高光谱散射图像技术进行苹果粉质化无损检测提供了一个有效的分类方法。

**关键词** 粉质化; 高光谱散射图像技术; 局部线性嵌入; 非线性降维; 支持向量机

中图分类号: O657.3 文献标识码: A DOI: 10.3964/j.issn.1000-0593(2010)10-2739-05

## 引言

苹果的粉质化是指苹果非正常软化、汁液减少和果肉质地发棉等一系列生理失调现象, 是影响苹果等级的重要口感参数<sup>[1]</sup>。如何利用无损检测方法代替传统的破坏性检测是苹果粉质化检测研究的一个重要趋势。当前, 一种能集光谱检测和图像检测优点的新技术——高光谱图像技术正好能满足水果产品检测技术发展的需要。国内外许多学者利用高光谱散射图像技术开展了对苹果、柑橘、梨等进行无损检测的研究工作, 并取得了较好的研究成果<sup>[2-6]</sup>。

高光谱数据一般含有几十甚至几百个波段, 高光谱图像具有高空间分辨率和时间分辨率、图谱合一等特点, 能为各项应用提供更详细的观测信息。然而, 高光谱特征和分类研究中主要存在以下两个难点: 一是高维使得计算速度受到很大影响, 训练样本的不足也会导致不好的分类结果, 即所谓的“维数灾难”或者 Hughes 现象; 二是波段间的强相关性增加了冗余性, 如果不能有效处理, 会对结果产生一定的影响。降维是一种有效的消除噪声并提取有用信息的方法<sup>[7]</sup>, 如何能够既有效降低特征空间的维数, 同时又要尽可能多的保留原始数据所包含的信息, 相关学者进行了大量的研究, 提出了连续投影、遗传算法、主成分分析等降维方法<sup>[8-10]</sup>。

流形学习算法是近年来出现的一类非线性降维方法。其基本思想是: 高维观测空间中的点由少数独立变量的共同作用在观测空间形成一个流形, 如果能有效的展开观测空间卷曲的流形或发现内在的主要变量, 就可以对该数据集进行降维。LLE 是流形学习算法的一个典型代表, 它试图保持数据的局部几何特征, 就本质上说, 它是将流形上的近邻点映射到低维空间的近邻点。具有较高的计算效率、较少的自由参数、成本函数的非迭代全局最优、实现容易等特点<sup>[11]</sup>。

结合 LLE 和 SVM, 本文提出一种基于高光谱散射图像的苹果粉质化分类方法: 先用 LLE 对苹果的高光谱图像数据作非线性降维, 再利用 SVM 进行分类。结果表明, 用 LLE-SVM 方法建立的分类器的性能优于 SVM。

## 1 实验材料与方法

### 1.1 实验材料

实验中 580 个“Red Delicious”样本由两部分组成, 180 个于 2008 年 10 月份采摘于美国密歇根州立大学农业实验站 (Michigan State University, MSU), 实验之前保存在可以控制的储藏条件下 (0°C, 2% O<sub>2</sub> 和 3% CO<sub>2</sub>); 余下的 400 个从商店购买 (commercial packinghouse, CP)。所有样本分为两组进行储藏, 第一组 240 个样本 (180 个来源于 CP, 60 个来

收稿日期: 2009-11-26, 修订日期: 2010-03-02

基金项目: 国家自然科学基金项目 (60805014) 和中央高校基本科研业务费专项资金项目 (JUSRP20913) 资助

作者简介: 赵桂林, 女, 1986 年生, 江南大学通信与控制工程学院硕士研究生 e-mail: lin861103@sina.com

\* 通讯联系人 e-mail: zhuqib@163.com

源于 MSU) 储藏在 4 °C 的冷藏室里面; 为了加速苹果的粉质化过程, 第二组 340 个样本(220 来源于 CP, 120 来源于 MSU) 储藏在相对湿度为 95%, 温度为 20 °C 的储藏室中。实验前, 所有的样本需要放到室温条件下至少 15 h。苹果的损坏棉质化标准值(压缩硬度及汁液含量) 测试参见文献[13], 实验过程中苹果的切割圆柱直径为 18 mm, 长度为 16 mm。

## 1.2 实验装置

高光谱图像数据是利用基于光谱仪的高光谱图像系统采集得到的<sup>[13]</sup>。它是由图像光谱仪(ImSpector V10, Spectral Imaging Ltd, Oulu, Finland)、高光谱摄像头(Model C4880 21-24A, Hamamatsu Photonics Systems, Bridgewater, NJ, USA), 250 W 的光纤卤素灯和一套输送装置等部件组成。高光谱摄像头的光谱范围为 400~1 000 nm, 光谱分辨率为 4.54 nm, 实际实验数据记录时, 取近似值 5 nm, 空间分辨率为 0.20 mm。

## 1.3 数据采集与处理平台

在高光谱图像数据采集前, 预先确定好高光谱摄像头的曝光时间以保证图像清晰; 确定好输送装置的速度以避免图像尺寸和空间分辨率失真。对于每个待测试的苹果而言, 在 200 ms 的曝光时间内获得 10 幅图像, 然后将这些图像平均, 将平均后的图像保存起来作为以后的研究用。

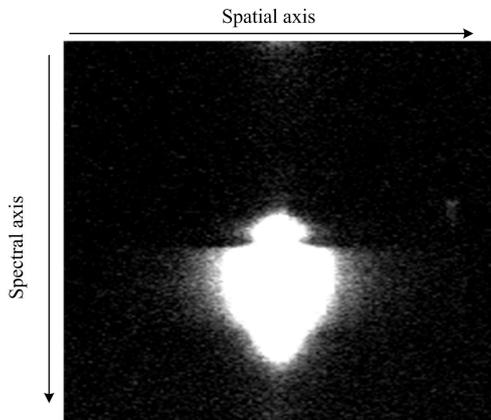


Fig 1 Hyperspectral scattering image of a "Red Delicious" apple

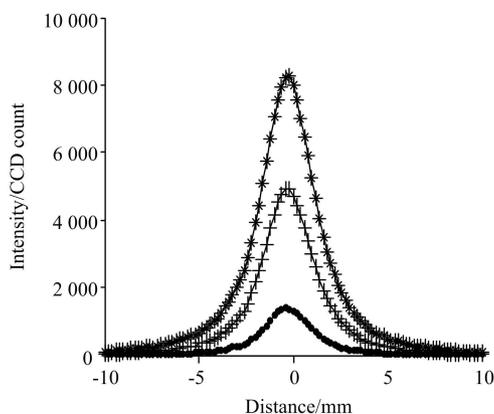


Fig 2 Raw spatial scattering profiles at three different wavelengths

1: 650 nm; 2: 700 nm; 3: 750 nm

图 1 是一副原始的苹果高光谱散射图像, 水平坐标代表空间位置, 垂直坐标代表光谱波段信息。每一副散射图像由特定波长下的 100 多幅空间散射图像组成。图 2 中的曲线分别为 650, 700 和 750 nm 下不同空间位置 CCD 的光谱强度。

实验中, 样本的光谱范围取为 600~1 000 nm, 每个样本含 81 个光谱波段信息, 也就是每隔 5 nm 取一个光谱波段值。

## 2 数据分析

### 2.1 图像特征提取

图像特征提取之前, 要对图像进行预处理, 包括图像读取、系统校正、距离校正、果形校正以及光源波动影响校正[参考 Jianwei Qin 的博士论文(2007, Department of Biosystems and Agricultural Engineering, Michigan State University)]. 由图 2 可知, 曲线的空间位置距离为 20 mm, 且沿纵轴对称, 所以在处理之前先对散射图像的两边取平均, 再利用平均值法进行特征提取<sup>[14]</sup>, 对于每个样本共得到 81 个特征描述。

### 2.2 LLE

刘建军等<sup>[15]</sup>主要是应用局部的线性来逼近全局的非线性, 它试图保持数据的局部几何特征, 就本质而言, 它是将流形上的近邻点映射到低维空间的近邻点。

LLE 算法的具体步骤如下:

Step1: 寻找每个样本点的  $k$  个近邻点。即对高维空间的每个样本点  $x_i$  ( $i = 1, 2, \dots, n$ ) 计算它和其他  $n-1$  个样本点之间的距离, 并且选择前  $k$  个与  $x_i$  ( $i = 1, 2, \dots, n$ ) 最近的点作为其近邻点, 这里的  $k$  是一个预先给定的值。距离的计算常采用欧式距离, 即  $d_{ij} = |x_i - x_j|$ 。

Step2: 由每个样本点的近邻点计算出该样本点的局部重建权值矩阵。即已知每个  $x_i$  ( $i = 1, 2, \dots, n$ ) 和它的  $K$  个近邻点, 计算该点和它的每个近邻之间的权重  $w_j^{(i)}$ , 即最小化成本函数(cost function)

$$\varepsilon_i(W) = \sum_{i=1}^n \left| x_i - \sum_{j=1}^k w_j^{(i)} x_j \right|^2 \quad (1)$$

其中, 权重  $w_j^{(i)}$  服从特定的对称性, 且  $\sum_{j=1}^k w_j^{(i)} = 1$ , 如果  $x_j$  ( $j = 1, 2, \dots, n$ ) 不是  $x_i$  ( $i = 1, 2, \dots, n$ ) 的近邻, 则  $w_j^{(i)} = 0$ 。

Step3: 由该样本点的局部重建权值矩阵和其近邻点计算出该样本点的输出值。降维的目的是在低维空间中尽量保持高维空间中的局部线性结构, 权重  $w_j^{(i)}$  代表着局部信息, 则固定权重  $w_j^{(i)}$ , 使下面的损失函数最小化

$$\varepsilon_r(Y) = \sum_{j=1}^n \left| y_j - \sum_{i=1}^k w_i^{(j)} y_i \right|^2 \quad (2)$$

要求  $\sum_{i=1}^n y_i = 0$  且  $\frac{1}{n} \sum_{i=1}^n y_i y_i^T = 1$ , 使  $\varepsilon_r(Y)$  对平移、旋转和伸缩变化都具有不变性。使  $\varepsilon_r(Y)$  最小化的解为矩阵  $M = (I - W)^T / (I - W)$  的最小几个特征值所对应的特征向量构成的矩阵  $Y$ , 将  $M$  的特征值从小到大排列, 舍去第一个几乎接近

0 的特征值, 取第 2~ d+ 1 之间的特征值所对应的特征向量作为输出结果。

### 2.3 SVM

设样本集  $(x_i, y_i), i = 1, 2, \dots, n, x \in R^m$  为输入的样本特征向量,  $y \in \{+1, -1\}$  为分类类别, 则 SVM 所构成的决策函数为

$$f(x) = \text{sgn} \left\{ \sum_{i=1}^n a_i y_i K(x_i \cdot x) + b \right\} \quad (3)$$

其中  $\text{sgn}(\cdot)$  为符号函数,  $K(\cdot)$  为满足 Mercer 条件的核函数,  $a_i$  为拉格朗日乘子, 其值可由下列最优化问题解得

$$\begin{aligned} \max L = & \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i,j} a_i a_j y_i y_j K(x_i, x_j) \\ \text{s.t.} & \begin{cases} 0 \leq a_i \leq C \\ \sum a_i y_i = 0 \end{cases} \end{aligned} \quad (4)$$

式中  $C$  为惩罚因子, 用于调节学习机的置信范围和风险的比例。

$b$  为阈值, 其值可由下式确定

$$b = y_j - \sum_{i=1}^L y_i a_i K(x_i, x_j), j \in \{j | C > a_j > 0\} \quad (5)$$

这样, 对于任意样本  $x$ , 可以根据  $f(x)$  的符号来判断样本所属类别。

满足 Mercer 条件的核函数有很多种, 在本文的应用中, 选择 RBF 核函数  $K(x, x') = \exp(-\|x - x'\|^2/\gamma^2)$ 。  $\gamma$  为 RBF 核函数的宽度, 用于描述样本数据在高维特征空间中分布的复杂程度。

## 3 结果与讨论

由高光谱图像原理可知, 苹果图像上每个像素都存在不同波长下的光谱信息。图 3 为 10 个粉质化(mealy) 苹果在 600~ 1 000 nm 范围内的光谱曲线, 图 4 是 10 个非粉质化(nonmealy) 苹果在 600~ 1 000 nm 范围内的光谱曲线。由图可以看出, nonmealy 的相对反射强度比 mealy 的相对反射强度大。

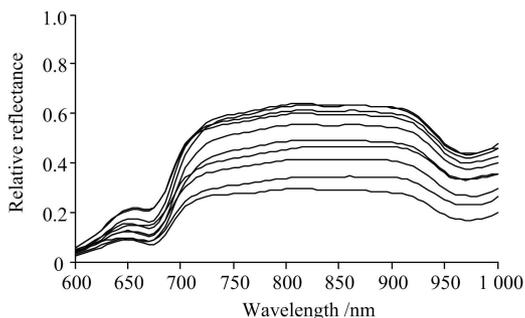


Fig 3 Relative mean reflectance for 10 "mealy" apples

### 3.1 数据集分类结果比较

对 580 个实验样本, 采用随机选择样本的方式, 选取 480 个样本作为分类器的训练样本, 剩余的 100 个作为测试样本。分别用 LLE-SVM 和 SVM 算法进行分类。表 1 给出了

降维前后分类精度的比较。在进行计算的过程中, 有 4 个可调参数, LLE 中的  $k=12, d=50$ ; SVM 中  $\gamma=20, C=200$ 。

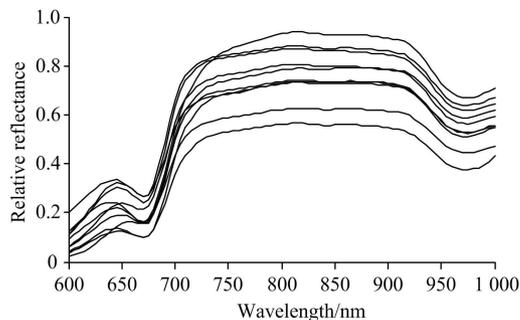


Fig 4 Relative mean reflectance for 10 "nonmealy" apples

Table 1 Comparison of classification accuracy of different intervals

训练样本个数	测试样本个数	SVM		LLE-SVM	
		测试精度/%	训练精度/%	测试精度/%	训练精度/%
480	100	75	85.29	79	94.12

由表 1 可以看出, 无论是训练精度还是测试精度, LLE-SVM 都要高于 SVM。

### 3.2 RBF 参数对分类结果的影响

为了更有力的说明 LLE-SVM 的优越性, 考察 SVM 参数  $C, \gamma$  对分类结果的影响, 同样取训练样本数为 480, 测试样本数为 100, 固定参数  $k=12, d=50$ 。实验结果如表 2 所示。

Table 2 Comparison of the classification results by different( $C, \gamma, C$ )

$\gamma$	$C$	SVM		LLE-SVM	
		测试精度/%	训练精度/%	测试精度/%	训练精度/%
20	300	72	84.03	75	91.60
20	500	72	82.56	76	89.29
20	1000	74	80.88	79	84.45
10	200	72	83.61	76	92.02
5	200	74	82.56	77	90.13
1	200	77	80.88	79	82.98

由表 2 可以看出, 经 LLE 降维后, 对于不同的核参数  $\gamma$  和  $C$ , LLE-SVM 的训练和测试精度敏感性虽有提高, 但是始终较单纯的 SVM 高。

### 3.3 LLE 参数对分类结果的影响

LLE 算法有 2 个可调的参数  $k$  和  $d$ , 其中  $k$  为邻域参数,  $d$  为样本本真维数。在不同参数下, LLE-SVM 得到不同的训练和测试精度。对数据集, 分别取  $k$  从 5 到 20,  $d$  从 5 到 80 作 LLE 降维, 再用 SVM 进行分类, 固定参数  $\gamma=20, C=200$ 。经初步尝试, 对于本例, 当  $k$  的取值范围为 10~ 20,  $d$  的取值范围为 20~ 50 时, 有比较理想的结果。

## 4 结 论

正常苹果和粉质化苹果在 600~1 000 nm 范围内的光谱反射强度有着较大的差异, 因此, 利用这一差异可以达到苹果粉质化分类的目的。

LLE 能够对高光谱进行非线性降维, 提取有用信息, 结合 SVM 可实现高光谱数据的分类。实验结果表明: LLE-SVM 算法在训练和测试精度方面都要好于单纯的 SVM 方

法。本研究为将流形学习算法引入高光谱数据分类提供了一种有益尝试。

当然除此之外, 也出现了许多新的问题。比如本文中只考虑了 RBF 参数对实验结果的影响, LLE 参数  $k$ ,  $d$  对分类结果的影响只得到一个比较理想的范围, 最优  $k$  和  $d$  参数的确定将是作者下一步研究需要解决的问题。

致谢: 论文作者对美国农业部 Postharvest Engineering Laboratory 的 Dr. Lu 在实验工作中的指导深表感谢。

## 参 考 文 献

- [ 1 ] Noh, Peng Y, Lu R, et al. Transactions of the ASABE, 2007, 50(3): 963.
- [ 2 ] Lu Renfu, Huang Min, Qin Jianwei. Sensing for Agriculture and Food Quality and Safety, 2009, 7315: 291.
- [ 3 ] Qin J, Lu R. Postharv. Bio. Techn., 2008, 49(3): 357.
- [ 4 ] ZOU Xiaobo, ZHAO Jiewen(邹小波, 赵杰文). Agricultural Non-Destructive Testing Technology and Data Analysis Methods(农产品无损检测技术与数据分析方法). Beijing: China Light Industry Press(北京: 中国轻工业出版社), 2008. 92.
- [ 5 ] XUE Long, LI Jing, LIU Muzhuo(薛龙, 黎静, 刘木华). Grain and Oil Processing(粮油加工), 2009, (4): 137.
- [ 6 ] CAI Jianrong, WANG Jianhe, CHEN Quansheng, et al(蔡健荣, 王建黑, 陈全胜, 等). Transactions of the Chinese Society of Agricultural Engineering(农业工程学报), 2009, 25(1): 128.
- [ 7 ] HOU Weiguang, DING Mingyue(侯文广, 丁明跃). Acta Electronica Sinica(电子学报), 2009, 37(11): 2580.
- [ 8 ] WANG Weijun, ZHANG Junying, YANG Liying(王文俊, 张军英, 杨利英). Journal of Sichuan University(四川大学学报), 2009, 41(6): 155.
- [ 9 ] GUAN Xiaoying, HU Xiaomin, ZHANG Jun(关晓颖, 胡晓敏, 张军). Computer Engineering and Design(计算机工程与设计), 2009, 30(1): 161.
- [ 10 ] ZHAO Chunhui, SONG Xiaoyue(赵春晖, 宋晓玥). Journal of Natural Science of Heilongjiang University(黑龙江大学自然科学学报), 2009, 26(5): 687.
- [ 11 ] ZHU Rong, YAO Min(朱容, 姚敏). Journal of Zhejiang University(Science)(浙江大学学报·科学版), 2009, 10(12): 1721.
- [ 12 ] Barreiro P, Ruiz Cabello J, Ortiz C, et al. Magnetic Resonance Imaging, 1998, 17(2): 275.
- [ 13 ] Qin J, Lu R, Peng Y. Transactions of the ASABE, 2009, 52(2): 500.
- [ 14 ] Lu R. Food Quality & Safety, 2007, 1(1): 24.
- [ 15 ] LIU Jianjun, XIA Shengping, YU Weixian(刘建军, 夏胜平, 郁文贤). Systems Engineering and Electronics(系统工程与电子技术), 2009, 31(2): 469.

# LLE SVM Classification of Apple Mealiness Based on Hyperspectral Scattering Image

ZHAO Guǎnlin, ZHU Qǎnbing\*, HUANG Min

School of Communication and Control Engineering, Jiangnan University, Wuxi 214122, China

**Abstract** Apple mealiness degree is an important factor for its internal quality. hyperspectral scattering, as a promising technique, was investigated for noninvasive measurement of apple mealiness. In the present paper, a locally linear embedding (LLE) coupled with support vector machine (SVM) was proposed to achieve classification because of large number of image data. LLE is a nonlinear lowering dimension method, which reveals the structure of the global nonlinearity by the local linear joint. This method can effectively calculate high dimensional input data embedded in a low dimensional space manifold. The dimension reduction of hyperspectral data was classified by SVM. Comparing the LLE-SVM classification method with the traditional SVM classification, the results indicated that the training accuracy obtained with the LLE-SVM was higher than that just with SVM; and the testing accuracy of the classifier changed a little before and after dimensionality reduction, and the range of fluctuation was less than 5%. It is expected that LLE-SVM method would provide an effective classification method for apple mealiness nondestructive detection using hyperspectral scattering image technique.

**Keywords** Mealiness; Hyperspectral scattering image; Locally linear embedding; Nonlinear dimensionality reduction; Support vector machine

(Received Nov. 26, 2009; accepted Mar. 2, 2010)

\* Corresponding author

欢迎订阅 欢迎投稿 欢迎刊登广告

## 《分析实验室》技术期刊

国内统一刊号: CN11-2017/TF

国际标准刊号: ISSN 1000-0720

国际 CODEN 码: FENSE4

邮发代号: 82-431

国外代号: M848

广告经营许可证: 京西工商广字第 0441 号

《分析实验室》是中文核心期刊, 月刊, 大 16 开, 128 页, 国内外公开发行。

《分析实验室》1982 年创刊, 目前已成为我国著名的分析化学专业刊物。影响遍及冶金、地质、石油化工、环保、药物、食品、农业、商品检验和海关等社会各行业及各学科领域。《分析实验室》以突出创新性和实用性为办刊宗旨, 作者来自全国各行业的生产、科研第一线; 已被列为全国中文核心期刊、中国科技论文统计用期刊、美国“CA 千种表”中我国化学化工类核心期刊、中国学术期刊(光盘版)和中国期刊网全文数据库等国内外多家检索数据库、文摘收录, 影响因子连续多年列化学类前列。本刊常设“研究报告”、“研究简报”、“仪器装置与设备”等栏目。“定期评述”栏目系统发布特邀知名专家学者撰写的国内外分析化学各领域的综合评述, 连续跟踪学术发展前沿。“国际会议”栏目每期介绍影响广泛的分析化学领域国际学术会议。

2011 年《分析实验室》每期定价 18 元, 全年 12 期, 216 元。

全国各地邮局征订, 邮发代号 82-431。漏订的读者可直接与编辑部联系。

编辑部地址: 北京新街口外大街 2 号 邮编: 100088

电话: 010-82013328 e-mail: analysislab@263.net; anrinfo@263.net