

文章编号: 2095-2163(2020)09-0086-05

中图分类号: TP391.4

文献标志码: A

基于电子舌和 DWT-PSO-LSSVM 模型的普洱茶存储年限快速检测

荆晓语, 缪楠, 杨正伟, 李庆盛, 张鑫, 王志强

(山东理工大学 计算机科学与技术学院, 山东 淄博 255049)

摘要: 普洱茶存储年限检测方式存在操作繁琐、分析过程复杂等问题, 为实现对不同存储年限的普洱茶的客观、快速的评价, 采用电子舌结合离散小波变换-粒子群优化最小二乘支持向量机(DWT-PSO-LSSVM)模型对5种不同存储时间的普洱茶样本进行定性分析。针对电子舌输出信号特点, 采用离散小波变换(DWT)作为特征提取方法对输出信号预处理。在此基础上, 采用粒子群优化最小二乘支持向量机(PSO-LSSVM)对不同存储年限的普洱茶进行分类鉴别。实验表明, 与传统机器学习模型相比, DWT-PSO-LSSVM模型对普洱茶存储年限的分类效果更优, 其精确率(Precision)、召回率(Recall)和F1-Score分别达到94.8%、94%和0.936。结果证实, DWT-PSO-LSSVM结合电子舌适合于对普洱茶存储年限进行快速检测, 且具有较高的分类准确性。

关键词: 普洱茶; 电子舌; 离散小波变换; 粒子群优化算法; 最小二乘支持向量机; 快速检测

Fast detection of storage year of Pu'er tea based on electronic tongue and DWT-PSO-LSSVM model

JING Xiaoyu, MIAO Nan, YANG Zhengwei, LI Qingsheng, ZHANG Xin, WANG Zhiqiang

(School of Computer Science and Technology, Shandong University of Technology, Zibo Shandong 255049, China)

[Abstract] The detection method of storage year of Pu'er tea has the problems of complicated operation and analysis process. In order to realize fast and objective evaluation of pu'er tea with different storage year, electronic tongue combined with DWT-PSO-LSSVM detection model are adopted to conduct qualitative analysis of tea samples with 5 different storage years. According to the characteristics of the output signal of electronic tongue, discrete wavelet transform is used as the feature extraction method for preprocessing. On this basis, particle swarm optimization-least squares support vector machine (PSO-LSSVM) is used to classify and identify the storage year of Pu'er tea. The experiment shows that compared with traditional machine learning models, DWT-PSO-LSSVM model possesses a better classification effect for pu'er tea, in which the Precision, Recall and F1-Score are 4.8%, 94% and 0.936, respectively. The results show that DT-PSO-LSSVM combined with electronic tongue are suitable for fast detection of storage year of Pu'er tea and has high classification accuracy.

[Key words] Pu'er tea; electronic tongue; discrete wavelet transform; particle swarm optimization algorithm; least square support vector machines; rapid detection

0 引言

普洱熟茶是以云南大叶种晒青毛茶为原料经特殊发酵后加工而成的^[1]。普洱茶具有抗氧化、降血脂、降血糖、抑菌、助消化、醒酒、解毒等作用^[2]。随着贮存时间的增加, 普洱茶的内部会发生复杂的化学变化, 使得普洱茶的风味和口感有很大提升^[3]。近年来, 受经济利益驱动, 市场上常会出现普洱茶产品以新替旧、以次充好等现象, 严重损害了消费者的权益和普洱市场的声誉。传统普洱茶存储年限鉴别方法主要有感官分析法和理化分析法。感官分析法受人为因素影响较大, 结果的客观性容易受到干扰^[4]; 理化分析法目前主要采用傅里叶红外光谱、

表面增强拉曼光谱、高效液相色谱法、近红外光谱等技术对普洱茶中的酚类、醇类、酸类、酯类等成分进行分析, 但是其检测仪器成本高、体积大、分析过程繁琐、耗时耗力。

电子舌是一种新型的现代化分析仪器, 应用多元统计分析 with 电极阵列相结合的方法, 分析检测复杂的普洱茶溶液样本, 这种方法优势在于操作便捷、容易携带、成本低、检测速度快以及再现性良好。近几年来, 电子舌在食品的定性或定量分析过程中取得了一定的效果, 现已在茶叶、水、酒、肉等食品中展开应用。电子舌系统的关键技术是模式识别, 其适用性直接影响到检测结果的准确性, 主要包含特征

基金项目: 山东省自然科学基金(ZR2019MF024); 赛尔网络下一代互联网技术创新项目(NGII20170314); 教育部科技发展中心产学研创新基金(2018A02010)。

作者简介: 荆晓语(1994-), 女, 硕士研究生, 主要研究方向: 机器学习、模式识别。

通讯作者: 王志强 Email: wzq@sdut.edu.cn

收稿日期: 2020-07-10

提取与分类识别两个步骤。特征提取是从原始的电子舌信号中提取出最重要的信息,从而减少后续数据分析的复杂性。截面积法、特征点法、傅里叶变换和主成分分析法为当前采用的特征提取方法,但这些方法只能提取信号中有限的表征参数,样本的整体信息无法全面的反映出来^[5-6];分类识别是基于所提取的特征信息对样本分类或识别的方法。神经网络、支持向量机、随机森林、极限学习机为现常用的分类识别方法,但因工作参数优化受限的问题,很难使分类效果达到最优。

基于前期对普洱茶的研究,以5种储存年限的普洱茶为实验材料,使用实验室自主研发的电子舌系统为平台对其进行识别。针对电子舌响应信号的特点,本文的数据预处理采用离散小波变换(discrete wavelet transform, DWT)对电子舌信号进行分析,对不同普洱茶样本分类识别是通过构建粒子群优化最小二乘支持向量机(Particle swarm optimization-least squares support vector machine, PSO-LSSVM)模型,从而建立了一种对普洱茶存储年限快速、准确鉴别的新方法。

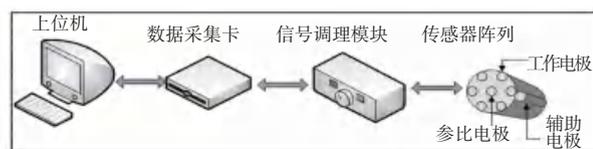
1 材料与方法

1.1 实验材料及样本处理

实验材料来自于勐海茶厂出产的普洱熟茶,出厂时间分别为2012年、2014年、2016年、2018年和2020年5个不同的年份。准确称量5g茶叶,将茶样放在200ml沸腾蒸馏水冲泡5min,茶溶液经滤纸过滤,冷却至室温(25±2℃)后采用伏安电子舌进行数据采集。每个样本采集完成后,用AL₂O₃粉末对传感器阵列进行打磨,放入超声波清洗仪中清洗。共采集到250个样本信号,其中2012年、2014年、2016年、2018年和2020年的样品数量均为50个,用于模型训练和测试的样本数量之比为8:2。

1.2 电子舌系统

实验设备采用本实验室自行开发的大幅脉冲电子舌系统、传感器、信号调理电路、数据采集卡和基于LabVIEW的上位机软件组成的系统。本文自主研发的系统结构如图1所示。工作过程为:由上位机软件控制数据采集卡产生大幅方波脉冲激励信号,该信号通过信号调理模块施加至传感器阵列,在激励信号的激发下,样本溶液在传感器工作电极表面发生电化学反应,而且产生弱电流响应信号。通过信号调理模块I/V对电流信号进行转换、放大和滤波,然后送入数据采集卡进行模/数转换,最后送入上位机进行信号处理和模式识别分析^[7]。



(a) 结构示意图

(a) Structure diagram



(b) 实物图

(b) Physical picture

图1 电子舌系统结构图

Fig. 1 Electronic tongue system structure

1.3 数据分析方法

1.3.1 离散小波变换

小波变换被称为“数字显微镜”,其可以对信号进行多分辨率、时频域、自适应等局部化的分析^[8]。由于其分辨能力强、压缩效果好、有效信息保存完整,已被成功应用在激光诱导击穿光谱、医学影像分析等领域中。离散小波变换(DWT)是一种在尺度和位移上离散化处理的小波变换。本研究根据电子舌响应信号的特点采用离散小波变换对其进行Mallat小波分解,公式(1)如下:

$$\begin{aligned} P_k^j &= \sum_{m \in Z} \alpha_{m-2k} P_m^{j-1} \\ r_k^j &= \sum_{m \in Z} \beta_{m-2k} P_m^{j-1} \end{aligned} \quad (1)$$

其中, P_k^j 为原始信号经小波分解后的第 j 层低频分量,即近似系数; r_k^j 为原始信号小波分解后的第 j 层高频分量,即细节系数; α 为低通滤波器; β 为高通滤波器; α_{m-2k} 为低通滤波器系数; β_{m-2k} 为高通滤波器系数; j 为分解层数, P_m^{j-1} 为原始信号小波分解后的第 $j-1$ 层低频分量,即近似系数。

电子舌数据经小波分解后,低频分量被完整的留了下来,而高频分量则被剔除。原有信号中有效信息被保留下来,一方面可以获得特征提取的作用,另一方面起到数据压缩的作用。

DWT对电子舌数据的压缩效果主要取决于小波基函数和分解层数的选择。特征信息的提取效果是由小波基函数决定的,最终的数据规模是由分解层数决定的。为了尽量减少特征信息丢失,寻求最优的压缩效果,本研究采用波形相似系数对离散小波变换效果进行可视化评价,再选择最合适的分解

层数和小波基函数,波形相似系数公式(2):

$$f_c = \frac{A \cap B}{A \cup B} = \frac{\sum_{i=1}^N (\max\langle a_i, b_i \rangle - |a_i - b_i|)}{\sum_{i=1}^N (\min\langle a_i, b_i \rangle + |a_i - b_i|)} \quad (2)$$

其中, A 为电子舌原始信号数据; B 为经离散小波变换后的压缩重构信号数据; a_i 为 A 的第 i 个数据; b_i 为 B 的第 i 个数据; f_c 为波形相似系数 ($0 \sim 1$), 反映了原始信号与压缩重构信号之间的相似性, f_c 越大表明两波形之间的相似性越高。

1.3.2 PSO-LSSVM

最小二乘支持向量机 (LSSVM) 是在 SVM 的基础上衍生的一类算法。其将 SVM 中的二次规划问题转化为线性方程组的求解问题, 从而提高了算法的求解速度和泛化能力^[9]。LSSVM 的性能主要取决于正则化参数和核函数宽度的选择, 本文采用粒子群优化算法 (PSO) 对 LSSVM 的参数进行全局寻优, 从而提高 LSSVM 算法的分类准确率, 其流程如图 2 所示。

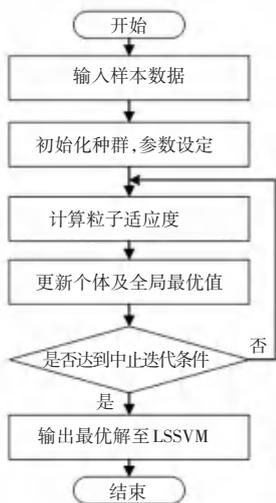


图 2 PSO 优化 LSSVM 算法的基本流程

Fig. 2 Basic flow of PSO optimization LSSVM algorithm

PSO-LSSVM 算法步骤如下:

对于学习样本:

$$T = \{(x_i, y_i) \mid i = 1, 2, \dots, n\}$$

(1) 假设最优回归函数, 式(3):

$$y = \omega^T \cdot x + b. \quad (3)$$

其中, ω 为权向量, b 为偏置量。

(2) 将回归问题转化为求解最优问题, 式(4):

$$\min J(\omega, \xi) = \frac{\|\omega\|^2}{2} + \frac{1}{2} \gamma \sum_{i=1}^n \xi_i^2. \quad (4)$$

其约束条件为式(5):

$$y_i = \omega^T \cdot x_i + b + \xi_i, \quad i = 1, 2, \dots, n. \quad (5)$$

其中, ξ_i 为松弛因子, γ 为正则化参数。

(3) 引入拉格朗日函数, 式(6):

$$L(\omega, b, \xi, \alpha) = J(\omega, \xi) - \sum_{i=1}^n \alpha_i (\omega^T \cdot x_i + b + \xi_i - y_i). \quad (6)$$

其中, α_i 为拉格朗日因子。

(4) 选用高斯径向基函数 (RBF) 作为核函数, 式(7):

$$K(x, x_i) = e^{-\frac{\|x - x_i\|^2}{2\sigma^2}}. \quad (7)$$

其中, σ^2 为核函数宽度。

(5) 采用 PSO 算法对正则化参数 γ 和核函数宽度 σ^2 进行寻优, 式(8):

$$\begin{cases} v_{id}^{k+1} = v_{id}^k + c_1 r_1^k (p_{id}^k - x_{id}^k) + c_2 r_2^k (p_{gd}^k - x_{id}^k) \\ x_{id}^{k+1} = x_{id}^k + v_{id}^{k+1}, \quad d = 1, 2, 3, \dots \end{cases} \quad (8)$$

其中, i 代表第 i 个粒子; d 为粒子的维度; k 为迭代次数; v_{id}^k 为粒子 i 的速度; x_{id}^k 为粒子 i 的位置; p_{id}^k 为粒子 i 的最优位置; p_{gd}^k 为群落的全局最优位置; c_1 和 c_2 为学习因子; r_1 和 r_2 为 $[0, 1]$ 范围内随机数。

(6) 确定 LSSVM 回归函数, 式(9):

$$y = \sum_{i=1}^n \alpha_i K(x, x_i) + b. \quad (9)$$

2 试验与分析

2.1 电子舌响应信号

图 3 为 8 个工作电极在大幅脉冲信号激励下单个普洱茶样本的响应信号。从图 3 可以看出, 在相同样本中, 8 个工作电极的响应信号之间存在着较大差异。为了避免实验数据的复杂性, 应选择尽可能少并能反映样本整体信息的电极。经过多次实验验证分析, 玻、碳、镍、钨、铂、钨、钛、金、银 8 个电极可以基本反映普洱茶样本的整体信息。因此本实验采用电子舌系统单次检测, 从普洱茶样本中获得 8 000 个数据点。

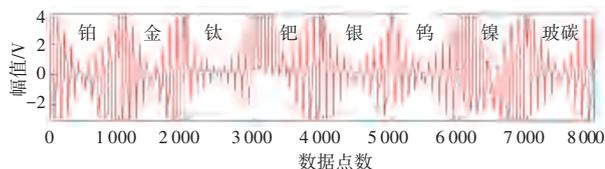


图 3 工作电极阵响应信号

Fig. 3 Response signal of working electrode array

2.2 模式识别处理

2.2.1 DWT-PSO-LSSVM 的建立

分别采用 Symlets、Daubechies、Haar、Coiflets 小

波函数对普洱茶电子舌响应信号执行4至7层分解。其重构信号与原信号的波形系数 f_c 变化情况,如图4所示。

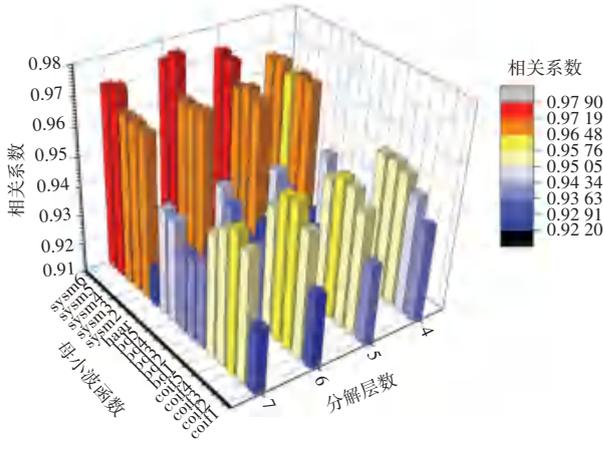


图4 不同小波基和分解层数下的波形相似系数统计

Fig. 4 Waveform similarity coefficient statistics under different wavelet bases and decomposition layers

从图4可以看出,不同的小波基对应的波形系数存在较大不同。整体而言,随着分解层数不断地增加, f_c 变化均呈现下降的趋势,说明在小波变换过程中,随着分解层数增加,能够减少更多的冗余信息,数据量会大大减少。通过分析,最后选择sym6

小波基函数,分解层数为7层。经过DWT预处理后,信号原始数据由8000个点压缩至73个,极大减少了后续模式识别处理的难度。

采用PSO算法对LSSVM模型的参数进行优化,PSO算法参数设置为: $c_1 = 1.5, c_2 = 1.7$,初始粒子数为20,最大迭代次数为200。通过对训练数据的学习及PSO算法对LSSVM模型参数的优化调整,得到LSSVM模型对普洱茶年限预测的最优参数为 $\sigma = 81.4352, c = 43.0277$ 。

2.2.2 模型性能对比分析

基于DWT方法提取电子舌特征信息,利用测试集数据分别建立基于BPNN、LSSVM和PSO-LSSVM模型,结果如图5所示。其中横坐标为预测类别,纵坐标为目标类别,0~8分别代表5种不同普洱茶的存储年限。从图5中可以看出,LSSVM对普洱茶年限鉴别的正确分类样本数大于BPNN模型,而PSO-LSSVM的正确分类的样本数大于LSSVM。说明DWT-PSO-LSSVM模型对普洱茶具有良好的区别辨识能力,可以对不同年限的普洱茶进行有效鉴别。

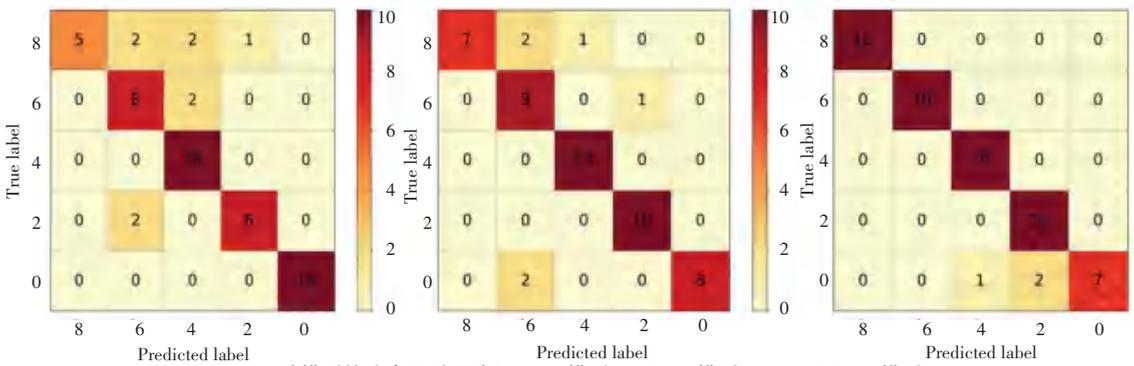


图5 3种模型的分类型混淆矩阵(BPNN模型、LSSVM模型、PSO-LSSVM模型)

Fig. 5 Classification confusion matrix of three models (BPNN model, LSSVM model, PSO-LSSVM model)

对三种模型采用精确率(Precision)、召回率(Recall)和F1-Score参数对比分析,结果见表1,各参数计算公式(10)、(11)和(12):

$$Precision = \frac{T_p}{T_p + F_p}, \quad (10)$$

$$Recall = \frac{T_p}{T_p + F_N}, \quad (11)$$

$$F1 - Score = \frac{2 \times Precision \times Recall}{Precision + Recall}. \quad (12)$$

其中, T_p 表示真实的正样本数量; F_p 表示真实的负样本数量; F_N 表示虚假的负样本数量。

从表1中可以看出,PSO-LSSVM模型的

Precision、Recall和F1-Score分别为94.8%、94%、0.936,均高于BPNN和LSSVM模型。

表1 不同模式识别模型的性能对比

Tab. 1 Performance comparison of different pattern recognition models

模型	Precision/%	Recall/%	F1-Score
BPNN	85.4	82	0.814
LSSVM	90.2	88	0.878
PSO-LSSVM	94.8	94	0.936

3 结束语

采用实验室研究的电子舌系统对不同存储时间

(下转第94页)