

文章编号: 1001 - 9014(2009)06 - 0423 - 05

基于无信息变量消除法和连续投影算法的 可见 近红外光谱技术白虾种分类方法研究

吴迪¹, 吴洪喜^{2, 3*}, 蔡景波^{2, 3}, 黄振华^{2, 3}, 何勇^{1*}

(1. 浙江大学 生物系统工程与食品科学学院, 浙江 杭州 310029;

2 浙江省海洋水产养殖研究所, 浙江 温州 325005;

3 浙江省近岸水域生物资源开发与保护重点实验室, 浙江 温州 325005)

摘要:应用无信息变量消除法结合连续投影算法对可见 近红外光谱区进行有效波长的选择,选择后的波长作为输入变量建立最小二乘支持向量机模型,对白虾属中三种典型种,脊尾白虾、秀丽白虾和东方白虾进行鉴别分类.实验采用 Kennard-Stone 算法选取 150 个样本作为建模集,50 个样本作为预测集,通过 UVE-SPA 优选了数值分别为 392、431、517、551、595、627、676、734、760、861、943 和 1018 nm 的 12 个波长为 LS-SVM 的输入变量,建立了白虾种分类模型.该模型对 50 个预测集样本检验的准确率达到 92.00%.结果表明,采用可见 近红外光谱对白虾种进行鉴别是可行的,UVE-SPA 能够有效地进行波长选择,使 LS-SVM 模型获得最优的分类结果.

关键词:可见 近红外光谱;无信息变量消除;连续投影算法;最小二乘支持向量机

中图分类号:O657.33 文献标识码:A

CLASSIFYING THE SPECIES OF EXOPALAEON BY USING VISIBLE AND NEAR INFRARED SPECTRA WITH UNINFORMATIVE VARIABLE ELIMINATION AND SUCCESSIVE PROJECTIONS ALGORITHM

WU Di¹, WU Hong-Xi^{2, 3*}, CAI Jing-Bo^{2, 3}, HUANG Zhen-Hua^{2, 3}, HE Yong^{1*}

(1. School of Biosystems Engineering and Food Science, Zhejiang University, Hangzhou 310029, China;

2 Zhejiang Mariculture Research Institute, Wenzhou 325005, China;

3 Zhejiang Key Lab of Exploitation and Preservation of Coastal Bio-resource, Wenzhou 325005, China)

Abstract: Using visible-near infrared spectra to classify different species of exopalaemon was studied. Successive projections algorithm (SPA) combined with uninformative variable elimination (UVE) were used to select effective wavelengths from visible and near infrared (Vis-NIR) bands. The selected effective wavelengths were set as inputs of least square-support vector machine (LS-SVM) for the classification of three typical exopalaemon species, namely, *E. carincauda*, *E. modestus* and *E. orientis*. Kennard-Stone algorithm was used to select 150 samples for calibration and the remaining 50 samples for prediction. Twelve effective wavelengths were selected by UVE-SPA, and they were 392, 431, 517, 551, 595, 627, 676, 734, 760, 861, 943 and 1018 nm. The correct rate is 92.00% for classifying samples in prediction set by LS-SVM model based on these twelve effective wavelengths. The overall results demonstrate that it is feasible to utilize Vis-NIR spectroscopy to classify different species of exopalaemon, and UVE-SPA can extract the most effective wavelengths to build the LS-SVM model with an optimal classification result.

Key words: visible-near infrared spectroscopy; uninformative variable elimination (UVE); successive projections algorithm (SPA); least square-support vector machine (LS-SVM)

收稿日期: 2008 - 10 - 22, 修回日期: 2009 - 05 - 22

Received date: 2008 - 10 - 22, revised date: 2009 - 05 - 22

基金项目: 国家科技支撑项目 (2006BAD10A04); 国家高技术研究发展计划 (863 计划) 项目 (2006AA10Z234); 浙江省自然科学基金项目 (Y506152); 浙江省台州市重大科技招标项目 (20071ZB02)

作者简介: 吴迪 (1984-), 男, 浙江杭州人, 博士生, 主要从事数字农业和农产品光谱检测技术研究

* 通讯联系人: 何勇: yhe@zju.edu.cn; 吴洪喜: whxchina@126.com

引言

白虾属是十足目真虾次目长臂虾科下的一个较小的属。目前,我国已鉴定出 4 个种,分别为脊尾白虾 (*E. carinicauda*)、秀丽白虾 (*E. modestus*) 东方白虾 (*E. orientis*, 俗称“绿籽虾”)和安氏白虾 (*E. anandalei*), 皆为经济虾类。该属种类体长一般不超过 90mm, 甲壳薄而透明, 微带蓝褐或红色点, 死后个体呈白色, 因此不易区别, 尤其是均生活在海水中的脊尾白虾和东方白虾更易混淆。在浙江, 若在脊尾白虾养殖塘混入了东方白虾, 东方白虾极易形成优势种群, 使养殖效益明显降低, 给白虾围塘养殖业带来重大的经济损失, 引起了当地有关部门的高度重视。

可见近红外光谱 (Visible and near infrared, Vis-NIR) 分析技术, 是一种快速、无损、低成本、无污染的分析方法, 它可以对物质的品质、种类、成分等进行定性和定量分析^[1,2], 但在虾种分析中的应用还较少。由于可见近红外光谱信息重叠严重, 利用全波段进行建模分析时, 光谱中的大量冗余信息及噪声等使模型的性能受到影响。如何在纷繁复杂的光谱信息中提取出有用信息, 提高模型校正的速度和建模效率是本文研究的重点内容。

常用的波长选择方法, 如相关系数法^[3]、载荷值法^[4]、回归系数法^[4]等大多根据主观经验进行阈值选择, 而退火算法^[5]和遗传算法^[6]的搜寻过程非常耗时, 且不稳定。无信息变量消除算法 (uninformative variable elimination, UVE) 是基于偏最小二乘 (partial least squares, PLS) 的回归系数建立的波长选择算法, 用于消除不提供信息的变量, 减少模型变量, 降低模型的复杂性^[7]。但是有时 UVE 得到的变量数依然较多, 因此仍需从 UVE 得到的变量中进一步选择有效变量。连续投影算法 (successive projections algorithm, SPA) 是一种新的变量提取方法^[8,9], 它能够利用向量的投影分析, 寻找含有最低限度的冗余信息的变量组, 并使变量之间的共线性达到最小, 同时能大大减少建模所用变量的个数, 提高建模的速度和效率。

本文应用 UVE-SPA 提取的有效波长作为最小二乘支持向量机 (least square-support vector machine, LS-SVM)^[10] 的输入, 建立 UVE-SPA 变量选择模型, 应用于白虾品种鉴别的研究。并将此结果与分别采用全波段波长及仅应用 UVE 得到的有效波长作为 LS-SVM 模型输入的结果进行比较。

1 材料与方法

1.1 仪器设备

实验使用美国 ASD (Analytical Spectral Device, Boulder, USA) 公司的 Handheld FieldSpec 光谱仪, 其光谱范围为 325 ~ 1075nm, 探头视场角为 15°; 光谱扫描次数设为 30 次, 光源采用 14.5 V 卤素灯; 试验采用漫反射模式; 分析软件采用 Unscrambler V9.7 (CAMO AS, Oslo, Norway) 和 Matlab V7.6 (The Math Works, Natick, USA)。

1.2 样本制备及样本集划分

实验选用我国白虾属中的脊尾白虾、秀丽白虾和东方白虾进行鉴别分析。试验中将同种类的虾均匀平铺在光谱仪采集视角范围内, 样本覆盖背景, 进行光谱采集。脊尾白虾和秀丽白虾各采集 50 个样本, 东方白虾采集了 100 个样本。采用 Kennard-Stone 算法^[11]选取 150 个样本作为建模集, 剩余 50 个样本作为预测集。建模集样本用于模型建立, 预测集样本用于对模型的预测性能进行检验。

1.3 无信息变量消除算法原理及算法实现

在 PLS 模型中, 光谱矩阵 X 和浓度矩阵 Y 存在如下的关系:

$$Y = Xb + e \quad (1)$$

其中, b 是系数向量, e 是误差向量。UVE 就是把一定变量数目的随机变量矩阵加入光谱矩阵中, 然后通过交叉验证建立 PLS 模型, 得到系数矩阵 B , 分析系数向量 b 的平均值和标准偏差的商 C 的稳定性, 即:

$$C_i = \frac{\text{mean}(b_i)}{S(b_i)} \quad (2)$$

其中, $\text{mean}(b)$ 表示系数向量 b 的平均值, $S(b)$ 表示系数向量 b 的标准偏差, i 表示光谱矩阵中的第 i 列向量。根据 C_i 的绝对值大小确定是否把第 i 列变量用于最后的 PLS 模型中。具体算法见参考文献 [7]。

1.4 连续投影算法实现

SPA 算法的简要介绍如下^[12]:

记 $x_{k(0)}$ 为初始迭代向量, N 为需要提取的变量个数, 光谱矩阵为 J 列。

(1) 迭代开始前, 任选光谱矩阵的 1 列 j 把建模集的第 j 列赋值给 x_j , 记为 $x_{k(0)}$;

(2) 把未选入的列向量位置的集合记为 $s, s = \{j, 1, \dots, J, j \notin \{k(0), \dots, k(n-1)\}\}$;

(3) 分别计算 x_j 对剩余列向量的投影: $Px_j = x_j - (x_j^T x_{k(n-1)}) x_{k(n-1)} (x_{k(n-1)}^T x_{k(n-1)})^{-1}, j \in s$;

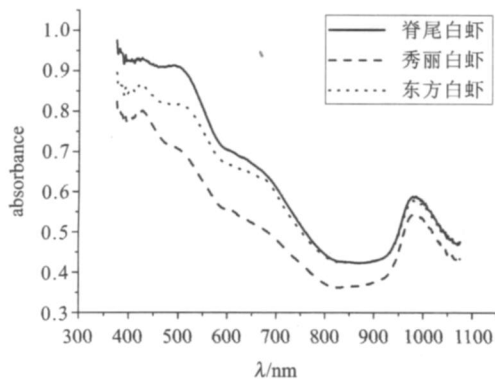


图 1 三种典型白虾的原始可见近红外吸收光谱

Fig 1 Original absorbance spectra of exopalaemon from three typical species

(4) 记 $k(n) = \arg(\max_j Px_j)$, $j \in s$;

(5) 令 $x_j = Px_j$, $j \in s$;

(6) $n = n + 1$, 如果 $n < N$, 回到 b 循环计算。

最后, 提取出的变量为 $\{x_{k(n)} = 0, \dots, N - 1\}$. 对应于每一个 $k(0)$ 和 N , 循环一次后进行多元线性回归分析 (MLR), 得到验证集的预测标准偏差 (RMSEV), 最小的 RMSEV 值对应的 $k(0)$ 和 N 就是最优值。

2 试验结果与分析

2.1 不同白虾种的可见近红外光谱图

由于测量到的光谱在 325 ~ 375 nm 范围内存在较大的噪声, 因此选用 376 ~ 1075 nm 波长范围内共计 700 个变量进行分析. 不同种类白虾的典型可见近红外吸收光谱图如图 1 所示, 图中横坐标为波长, 纵坐标为吸光度. 从图 1 可以看出, 不同白虾种的光谱曲线的趋势非常相似, 在 990 nm 附近有明显的水的吸收峰. 然而只从光谱特征上, 难以区分不同种类的白虾. 因此需要用化学计量学建模方法对光谱数据进行处理.

2.2 应用 UVE 选取有效波长

采用平滑、SNV、求导等对全波长光谱进行预处理

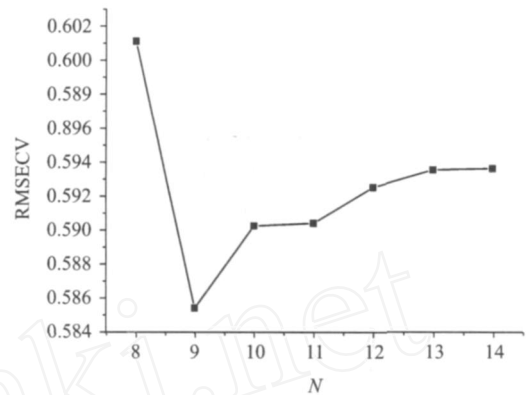


图 2 不同主成分数的 UVE 得到的 RMSECV 分布图

Fig 2 RMSECV plot of different numbers of principal components by UVE, where N is number of principal components

理, 然后作为 LS-SVM 的输入建立模型. 从表 1 可以看出, 采用原始光谱全波长作为输入变量时的结果最好, 预测集的正确率达到 90%. 采用原始光谱数据的结果优于平滑、求导等处理. 这可能是因为在预处理过程中, 引入或产生了新的噪声, 或者是处理后的有用信息量下降, 或者两者并存, 使得所建模型的预测性能下降. 虽然采用全波长作为输入变量已经达到了较好的分类结果, 但是 LS-SVM 模型的输入变量有 700 个, 部分变量可能仍包含无用的或者不相关的信息, 并且上百个输入变量造成模型冗余、复杂^[13]. 为此采用 UVE 对未经预处理的全波长 700 个光谱变量进行选择.

UVE 中产生的随机变量个数同样设置为 700 个. UVE 中的最优主成分数是经内部交叉验证得到的. 从图 2 可以看到当主成分数为 9 时, UVE 选择变量的分类效果最好. 主成分为 9 时的 UVE 变量选择结果见图 3. 黑色竖线的左侧为 700 个变量的稳定性 C 分布曲线, 右侧为 UVE 中产生的 700 个随机变量的稳定性 C 分布曲线. 两条水平虚线表示变量选择的阈值上下限, 在虚线外的数值对应的变量被保留, 在虚线内的数值对应的变量不用于建模. 阈值

表 1 基于不同预处理和变量选择方法的 LS-SVM 模型的三种典型白虾分类结果

Table 1 Classification results of exopalaemon from three typical species by LS-SVM models by using different pretreatment and variable selection methods

pretreatment	variable selection methods	variables	calibration set ($n=150$)		prediction set ($n=50$)	
			correct rate (%)		correct rate (%)	
raw	none	700	99.33	90.00	99.33	90.00
SG + SNV	none	700	100.00	88.00	100.00	88.00
1 - Der	none	700	96.67	84.00	96.67	84.00
2 - Der	none	700	86.67	78.00	86.67	78.00
raw	UVE	291	100.00	88.00	100.00	88.00
raw	UVE-SPA	12	98.67	92.00	98.67	92.00

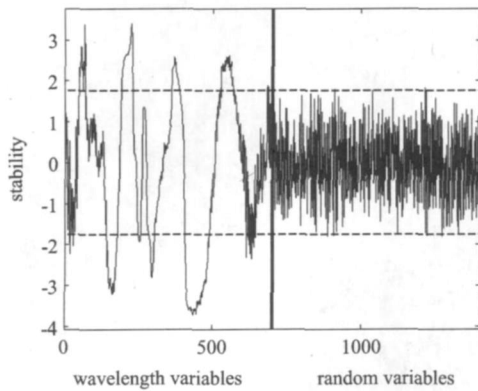


图3 主成分个数为9时的UVE稳定性分布曲线. 两条水平虚线表示变量选择的阈值上下限

Fig. 3 Stability distribution of UVE when nine components were considered. Two horizontal dot lines represent the threshold boundaries

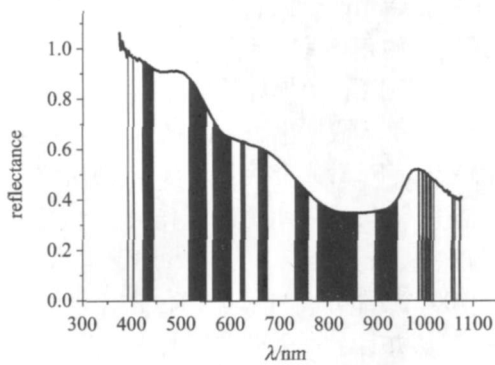


图4 UVE选择的291个波长分布图. 竖线为选中的波长
Fig. 4 Plot of 291 selected wavelengths by UVE. Columns represent selected wavelength

的选择标准为随机变量稳定性最大值的 99%. 经过 UVE的变量选择,最后得到了 291个波长,其分布情况见图 4,其中竖线表示选中的波长. 可以看到, UVE选择的波长在 423 ~ 443nm 的紫光部分、517 ~ 551nm 和 566 ~ 603nm 的绿光部分、659 ~ 676nm 和 734 ~ 760nm 的红光部分以及 780 ~ 860nm 和 1000nm附近的近红外范围都有一定的分布. 说明无论是可见光还是近红外光谱都包含关于白虾种分类的有用信息. 将得到的这 291个波长作为 LS-SVM 的输入变量,预测结果见表 1. 可以看到虽然模型的分正确率和 700个变量作为输入时相比差不多,但是输入变量数却从 700个减少到 291个,模型得到了优化.

2.3 应用 UVE-SPA 选取有效波长

经过 UVE选择过的变量个数仍过于庞大,因此采用 SPA对经过 UVE选择后的 291个变量做进一步选择,得到共线性最小的有效波长. 对以上 291个

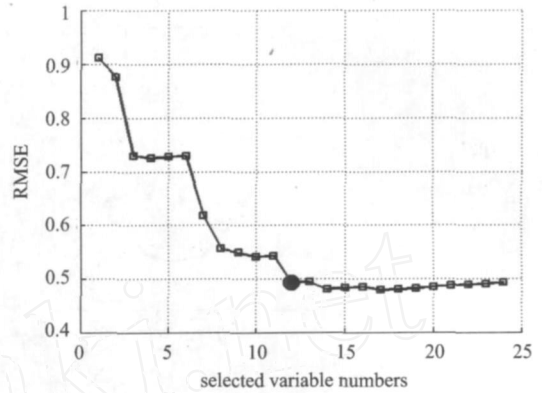


图5 SPA选择的不同变量数的RMSE分布
Fig. 5 RMSE plot of number of selected variables by SPA

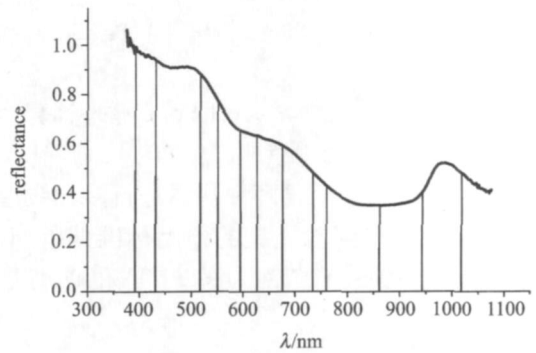


图6 UVE-SPA选择的12个波长(竖线表示)
Fig. 6 Twelve wavelengths selected by UVE-SPA. Columns represent selected wavelengths

变量进行建模,经过 SPA进一步选择的不同波长数的 RMSE分布如图 5所示,图中黑色实心圆点表示所选到的波长数. 由图 5可以看出,选取 3个波长时, RMSE有一个下降过程,然后到选择 8个波长时, RMSE又有一个下降过程,到提取出 12个有效波长时 RMSE到达一个低点,以后 RMSE趋于平坦,因此确定从 291个波长中选择出其中的 12个波长. 由 UVE-SPA得到的这 12个波长的分布情况如图 6所示,其中竖线表示选中的波长. 通过 UVE-SPA提取的 12个波长分别为 392, 431, 517, 551, 595, 627, 676, 734, 760, 861, 943和 1018nm. 将这 12个波长作为输入变量建立白虾品种的 UVE-SPA-LS-SVM模型,分类结果如表 1所示. 从表 1中可以看出,通过 UVE-SPA得到的 12个波长所建立的 LS-SVM模型的建模集分类正确率略低于采用未经预处理的全波段 700个波长以及仅仅采用 UVE选择的 291个波长建立的 LS-SVM模型,但是预测集的分类正确率却有一定的提高. 虽然提高的幅度不大,但是 UVE-SPA-LS-SVM模型仅采用 12个有效波长

建模.说明采用 UVE-SPA 所提取的有效波长能充分代表原始光谱的有效信息,可以作为波长提取的有效手段.选择的 12 个波长均匀地分布在可见和近红外光谱范围内,说明对于白虾种的分类不能仅仅考虑某个范围的波长.结果表明,采用 UVE-SPA 对可见近红外光谱进行有效波长选择后建立的 LS-SVM 模型对白虾种的分类是可行的,并且获得了满意的准确度.

3 结论

采用 UVE-SPA 对可见近红外光谱进行有效波长选择,选择后的波长作为输入变量建立 LS-SVM 模型对白虾种进行分类.通过 UVE-SPA 从全波长 700 个变量中选择得到了 12 个最能够反应光谱信息的波长.预测集样本的分类准确率达到 92.00%,优于采用全波段 700 个波长以及仅仅采用 UVE 选择的 291 个波长作为输入变量的 LS-SVM 模型.结果表明,应用可见近红外光谱技术对白虾种进行分类是可行的,采用 UVE-SPA-LS-SVM 组合模型能获得满意的分类结果,为白虾种的鉴别提供了方法依据.

REFERENCES

- [1] YAN Yan-Lu, ZHAO Long-Lian, HAN Dong-Hai, *et al* *The Foundation and Application of Near Infrared Spectroscopy Analysis*[M]. Beijing: China Light Industry Press (严衍祿,赵龙莲,韩东海,等.近红外光谱分析基础与应用.北京:中国轻工业出版社), 2005.
- [2] JWU Di, FENG Lei, ZHANG Chuan-Qing, *et al* Early detection of gray mold (*Cinerea*) on eggplant leaves based on Vis/Near infrared spectra [J]. *J. Infrared Millim. Waves* (吴迪,冯雷,张传清,等.基于可见近红外光谱技术的茄子叶片灰霉病早期检测研究.红外与毫米波学报) 2007, **26** (4): 269—273.
- [3] Min M, Lee W S. Determination of significant wavelengths and prediction of nitrogen content for citrus [J]. *Trans ASAE*, 2005, **48** (2): 455—461.
- [4] Wu D, He Y, Feng S. Short-wave near-infrared spectroscopy analysis of major compounds in milk powder and wavelength assignment [J]. *Anal. Chim. Acta*, 2008, **610** (2): 232—242.
- [5] Kalivas J H, Roberts N, Sutter J M. Global optimization by simulated annealing with wavelength selection for ultraviolet-visible spectrophotometry [J]. *Anal. Chem.*, 1989, **61** (18): 2024—2030.
- [6] Jouanrinbaud D, Massart D L, Leardi R, *et al* Genetic algorithms as a tool for wavelength selection in multivariate calibration [J]. *Anal. Chem.*, 1995, **67** (23): 4295—4301.
- [7] Centner V, Massart D L, De Noord O E, *et al* Elimination of uninformative variables for multivariate calibration [J]. *Anal. Chem.*, 1996, **68** (21): 3851—3858.
- [8] Araújo M C U, Saldanha T C B, Galvão R K H, *et al* The successive projections algorithm for variable selection in spectroscopic multicomponent analysis [J]. *Chemom. Intell. Lab. Syst.*, 2001, **57** (2): 65—73.
- [9] Galvão R K H, Araújo M C U, Frago W D, *et al* A variable elimination method to improve the parsimony of MLR models using the successive projections algorithm [J]. *Chemom. Intell. Lab. Syst.*, 2008, **92** (1): 83—91.
- [10] JWU Di, HE Yong, FENG Shui-Juan, *et al* Application of infrared spectra technique based on LS-SVM to the non-destructive measurement of fat content in milk powder [J]. *J. Infrared Millim. Waves* (吴迪,何勇,冯水娟,等.基于 LS-SVM 的红外光谱技术在奶粉脂肪含量无损检测中的应用.红外与毫米波学报), 2008, **27** (3): 180—184.
- [11] Macho S, Iusa R, Callao M P, *et al* Monitoring ethylene content in heterophasic copolymers by near-infrared spectroscopy standardisation of the calibration model [J]. *Anal. Chim. Acta*, 2001, **445** (2): 213—220.
- [12] CHEN Bin, MENG Xiang-long, WANG Hao. Application of successive projections algorithm in optimizing near infrared spectroscopic calibration model [J]. *Journal of Instrumental Analysis* (陈斌,孟祥龙,王豪.连续投影算法在近红外光谱校正模型优化中的应用.分析测试学报), 2007, **26** (1): 66—69.
- [13] Chauchard F, Cogdill R, Roussel S, *et al* Application of LS-SVM to non-linear phenomena in NIR spectroscopy: development of a robust and portable sensor for acidity prediction in grapes [J]. *Chemom. Intell. Lab. Syst.* 2004, **71** (2): 141—150.