

## 多任务最小二乘支持向量回归机及其在近红外光谱 分析技术中的应用研究

徐硕<sup>1</sup>, 乔晓东<sup>1</sup>, 朱礼军<sup>1</sup>, 安欣<sup>2</sup>, 张录达<sup>3\*</sup>

1. 中国科学技术信息研究所信息技术支持中心, 北京 100038
2. 对外经济贸易大学国际经济与贸易学院, 北京 100029
3. 中国农业大学理学院, 北京 100193

**摘要** 在近红外光谱定量分析中,许多模型分开考虑各种样品成分含量,失去了样品成分间潜在的联系。针对该问题,文章将建模分析每种样品成分含量的问题看作一个任务,将同时建模分析所有样品成分含量的问题转换为多任务学习问题。在 LS-SVR 的基础上,提出了多任务 LS-SVR(MTLS-SVR),并给出一种有效的大规模问题求解算法。最后,以高粱样品数据集为实验材料,建立了三种样品成分(蛋白质, 赖氨酸及淀粉)的同时定量分析模型。三种样品成分的预测值与实际值的平均相对误差分别为 1.52%, 3.04% 和 1.01%, 相关系数分别为 0.993 1, 0.894 0 和 0.940 6, 经分析比较,发现 MTLS-SVR 模型优于 PLS, LS-SVR 以及多因变量 LS-SVR(MLS-SVR),从而验证了 MTLS-SVR 模型的可行性和有效性。

**关键词** 近红外光谱; 化学计量学; 多任务 LS-SVR

**中图分类号:** O657.3    **文献标识码:** A    **DOI:** 10.3964/j.issn.1000-0593(2011)05-00-04

### 引言

近红外光谱技术<sup>[1]</sup>具有操作简单、分析速度快以及测定一次光谱可同时获得样品多种成分含量的独特优点,使其在作物品质分析上得到了广泛应用。目前近红外光谱定量分析采用的化学计量学建模方法比较多,有些已取得了不错的效果,比如偏最小二乘(PLS)<sup>[2]</sup>、支持向量回归机(SVR)<sup>[3]</sup>等。与其他方法相比,SVR 具有出色的学习及推广能力,但它需要求解一个二次规划问题,非常耗时。而最小二乘 SVR(LS-SVR)<sup>[4]</sup>用等式约束代替不等式约束,只需求解一个线性方程组,大大加快了求解速度,受到了人们越来越多地重视。

近红外光谱包含了样品中所有成分的光谱信息,但目前许多模型分开考虑样品成分含量,失去了样品成分间潜在的联系。近年来,机器学习领域中多任务学习(MTL)研究逐渐成熟。如果将建模分析每种样品成分含量的问题看作一个任务,则可将样品成分含量同时建模分析的问题转换为多任务学习问题,从真正意义上实现了样品多成分含量的同时建模分析。目前大部分 MTL 模型<sup>[5,6]</sup>是基于 Bayesian 的,但 Ev-

geniou 等<sup>[7]</sup>借助 Bayesian MTL 的思想,设计了一种正则化 MTL 模型,不过该模型仍需求解一个二次规划问题。本文在 Evgeniou 等<sup>[7]</sup>工作的基础上,提出了多任务 LS-SVR (MTLS-SVR)模型。类似于 LS-SVR,该模型也只需要求解一个线性方程组,而且本文也给出了一种有效的大规模问题求解算法。最后建模分析了高粱样品的三种成分(蛋白质, 赖氨酸及淀粉)含量,平均相对误差和相关系数指标均显著优于 PLS, LS-SVR 和 MLS-SVR。

### 1 最小二乘支持向量回归机(LS-SVR)

对于  $\forall n \in \mathbb{N}$ , 记  $\mathbf{N}_n = \{1, 2, \dots, n\}$ 。给定训练数据集  $\{(x_i, y_i)\}_{i=1}^N$ , 其中  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , 记  $\mathbf{y} = (y_1; y_2; \dots; y_N)$ 。则 LS-SVR<sup>[4]</sup>的原始问题可表述为

$$\begin{aligned} \min_{w \in \mathbb{R}^{n_h}, b \in \mathbb{R}, \xi \in \mathbb{R}^N} J(w, \xi) &= \frac{1}{2} w^T w + \gamma \frac{1}{2} \xi^T \xi & (1) \\ \text{s. t. } y_i &= w^T \varphi(x_i) + b + \xi_i, i \in \mathbf{N}_N & (2) \end{aligned}$$

其中  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^{n_h}$  为输入空间到某一高维(可能为无穷维)特征空间  $H$  的映射,  $n_h$  表示特征空间  $H$  的维度,  $\xi = (\xi_1; \xi_2; \dots; \xi_N)$ 。

收稿日期: 2010-10-24, 修订日期: 2011-03-09

基金项目: “十一五”国家科技支撑计划(2006BAH03B03), 中国科学技术信息研究所重点工作项目(2009KP01-3-2)和中央高校基本科研业务费专项资金(2009-2-05)资助

作者简介: 徐硕, 1979 年生, 中国科学技术信息研究所博士后

e-mail: xush@istic.ac.cn

\* 通讯联系人 e-mail: zhangld@cau.edu.cn

$\dots; \xi_N$ ) 为由松弛变量组成的向量。

通过 Lagrange 函数, 问题(1)~(2)的求解可转换为以下线性方程组的求解

$$\begin{bmatrix} 0 & \mathbf{e}^T \\ \mathbf{e} & \mathbf{K} + \gamma^{-1} \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{b} \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{y} \end{bmatrix} \quad (3)$$

其中  $\boldsymbol{\alpha} = (\alpha_1; \alpha_2; \dots; \alpha_N)$  为 Lagrange 乘子向量,  $\mathbf{e} = (1; 1; \dots; 1)$ ,  $\mathbf{I}$  为单位矩阵,  $K_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$  ( $i, j \in \mathbb{N}_N$ ),  $K(\cdot, \cdot)$  为满足 Mercer 定理<sup>[3]</sup> 的核函数。

记线性方程组(3)的解为  $\boldsymbol{\alpha}^* = (\alpha_1^*; \alpha_2^*; \dots; \alpha_N^*)$ , 则决策函数为

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^{*T} \varphi(\mathbf{x}) + b^* = \sum_{i=1}^N \alpha_i^* \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}) + b^* \\ &= \sum_{i=1}^N \alpha_i^* K(\mathbf{x}_i, \mathbf{x}) + b^* \end{aligned} \quad (4)$$

## 2 多任务 LS-SVR(MTLS-SVR)

假设有  $M (\geq 1)$  个任务, 对于任务  $m \in \mathbb{N}_M$ , 给定  $N_m$  个训练数据  $(\mathbf{x}_{m,i}, \mathbf{y}_{m,i})_{i=1}^{N_m}$ , 其中  $\mathbf{x}_{m,i} \in \mathbb{R}^d$ ,  $\mathbf{y}_{m,i} \in \mathbb{R}$ , 这样就有  $N = \sum_m N_m$  个训练数据。为方便, 记  $\mathbf{y} = (\mathbf{y}_1; \mathbf{y}_2; \dots; \mathbf{y}_M)$ , 其中  $\mathbf{y}_m = (\mathbf{y}_{m,1}; \mathbf{y}_{m,2}; \dots; \mathbf{y}_{m,N_m})$ 。对于近红外光谱技术, 所有样品成分对应的光谱数据是一样的, 即对所有任务, 输入向量  $\mathbf{x}_{m,i} \in \mathbb{R}^d$  ( $i \in \mathbb{N}_{N_m}$ ) 是相同的, 但本文提出的方法不限于此, 即它具有更广的适用性。

### 2.1 原始问题及推导过程

对于任务  $m \in \mathbb{N}_M$ , 记需要建模的回归函数为  $f_m(\mathbf{x}) = \mathbf{w}_m^T \varphi(\mathbf{x}) + b_m$ 。为引入样品成分间潜在的联系, 假设所有权重向量  $\mathbf{w}_m$  在某个均值向量  $\mathbf{w}_0$  周围波动, 波动的幅度用向量  $\mathbf{v}_m$  来描述, 即  $\mathbf{w}_m = \mathbf{w}_0 + \mathbf{v}_m$ 。如果任务间联系比较强,  $\mathbf{v}_m \rightarrow 0$ , 否则  $\mathbf{v}_m \rightarrow 0$ 。类似于正则化 MTL<sup>[7]</sup>, MTLS-SVR 的原始问题可表述为

$$\min_{\mathbf{w}_0 \in \mathbb{R}^{n_h}, \mathbf{v}_m \in \mathbb{R}^{n_h}, b \in \mathbb{R}^M, \boldsymbol{\xi} \in \mathbb{R}^N} J(\mathbf{w}_0, \mathbf{v}_m, \boldsymbol{\xi}) = \frac{1}{2} \mathbf{w}_0^T \mathbf{w}_0 + \frac{1}{2} \frac{\lambda}{M} \sum_{m=1}^M \mathbf{v}_m^T \mathbf{v}_m + \gamma \frac{1}{2} \sum_{m=1}^M \boldsymbol{\xi}_m^T \boldsymbol{\xi}_m \quad (5)$$

$$\text{s. t. } \mathbf{y}_{m,i} = (\mathbf{w}_0 + \mathbf{v}_m)^T \varphi(\mathbf{x}_{m,i}) + b_m + \xi_{m,i}, m \in \mathbb{N}_M, i \in \mathbb{N}_{N_m} \quad (6)$$

其中  $\mathbf{b} = (b_1; b_2; \dots; b_M)$  为偏置向量,  $\boldsymbol{\xi} = (\xi_1; \xi_2; \dots; \xi_M)$ ,  $\boldsymbol{\xi}_m = (\xi_{m,1}; \xi_{m,2}; \dots; \xi_{m,N_m})$  ( $m \in \mathbb{N}_M$ ) 为由松弛变量组成的向量。

问题(5)~(6)的 Lagrange 函数为

$$L(\mathbf{w}_0, \mathbf{v}_m, \mathbf{b}, \boldsymbol{\xi}, \boldsymbol{\alpha}) = J(\mathbf{w}_0, \mathbf{v}_m, \boldsymbol{\xi}) - \sum_{m=1}^M \sum_{i=1}^{N_m} \alpha_{m,i} \times ((\mathbf{w}_0 + \mathbf{v}_m)^T \varphi(\mathbf{x}_{m,i}) + b_m + \xi_{m,i} - \mathbf{y}_{m,i}) \quad (7)$$

其中  $\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1; \boldsymbol{\alpha}_2; \dots; \boldsymbol{\alpha}_M)$ ,  $\boldsymbol{\alpha}_m = (\alpha_{m,1}; \alpha_{m,2}; \dots; \alpha_{m,N_m})$  ( $m \in \mathbb{N}_M$ ) 表示 Lagrange 乘子向量, 则对应的 KKT 条件如下

$$\frac{\partial L}{\partial \mathbf{w}_0} = 0 \Rightarrow \mathbf{w}_0 = \sum_{m=1}^M \sum_{i=1}^{N_m} \alpha_{m,i} \varphi(\mathbf{x}_{m,i})$$

$$\frac{\partial L}{\partial \mathbf{v}_m} = 0 \Rightarrow \mathbf{v}_m = \frac{M}{\lambda} \sum_{i=1}^{N_m} \alpha_{m,i} \varphi(\mathbf{x}_{m,i})$$

$$\frac{\partial L}{\partial b_m} = 0 \Rightarrow \sum_{i=1}^{N_m} \alpha_{m,i} = 0, \quad m \in \mathbb{N}_M$$

$$\frac{\partial L}{\partial \xi_{m,i}} = 0 \Rightarrow \alpha_{m,i} = \gamma \xi_{m,i}, \quad m \in \mathbb{N}_M, i \in \mathbb{N}_{N_m}$$

$$\frac{\partial L}{\partial \alpha_{m,i}} = 0 \Rightarrow (\mathbf{w}_0 + \mathbf{v}_m)^T \varphi(\mathbf{x}_{m,i}) + b_m +$$

$$\xi_{m,i} - \mathbf{y}_{m,i} = 0, \quad m \in \mathbb{N}_M, i \in \mathbb{N}_{N_m} \quad (8)$$

类似于 LS-SVR<sup>[4]</sup>, 式(8)可表示为线性方程组

$$\begin{bmatrix} 0 & \mathbf{A}^T \\ \mathbf{A} & \mathbf{Q} + \gamma^{-1} \mathbf{I}_N + (M/\lambda) \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{b} \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{y} \end{bmatrix} \quad (9)$$

其中  $\mathbf{I}_N$  为  $N$  阶单位矩阵;  $\mathbf{A} = \text{blockdiag}\{\mathbf{e}_{N_1}, \mathbf{e}_{N_2}, \dots, \mathbf{e}_{N_M}\}$ ,  $\mathbf{e}_{N_m}$  ( $m \in \mathbb{N}_M$ ) 为包含  $N_m$  个元素的全 1 列向量;  $\mathbf{B} = \text{blockdiag}\{\boldsymbol{\Omega}_1, \boldsymbol{\Omega}_2, \dots, \boldsymbol{\Omega}_M\}$ ,  $\boldsymbol{\Omega}_m = \mathbf{Z}_m^T \mathbf{Z}_m$ ,  $\mathbf{Z}_m = (\varphi(\mathbf{x}_{m,1}), \varphi(\mathbf{x}_{m,2}), \dots, \varphi(\mathbf{x}_{m,N_m}))$  ( $m \in \mathbb{N}_M$ );  $\mathbf{Q} = \mathbf{Z}^T \mathbf{Z}$ ,  $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_M)$ 。

记线性方程组(9)的解为  $\boldsymbol{\alpha}^* = (\alpha_1^*; \alpha_2^*; \dots; \alpha_M^*)$ ,  $\boldsymbol{\alpha}_m = (\alpha_{m,1}^*; \alpha_{m,2}^*; \dots; \alpha_{m,N_m}^*)$  ( $m \in \mathbb{N}_M$ ),  $\mathbf{b}^* = (b_1^*; b_2^*; \dots; b_M^*)$ , 则任务  $m \in \mathbb{N}_M$  的决策函数为

$$\begin{aligned} f_m(\mathbf{x}) &= (\mathbf{w}_0^* + \mathbf{v}_m^*)^T \varphi(\mathbf{x}) + b_m^* = \sum_{m=1}^M \sum_{i=1}^{N_m'} \alpha_{m',i}^* \varphi(\mathbf{x}_{m',i})^T \varphi(\mathbf{x}) + \\ &\quad \frac{M}{\lambda} \sum_{i=1}^{N_m} \alpha_{m,i}^* \varphi(\mathbf{x}_{m,i})^T \varphi(\mathbf{x}) + b_m^* = \sum_{m=1}^M \sum_{i=1}^{N_m'} \alpha_{m',i}^* K(\mathbf{x}_{m',i}, \mathbf{x}) \\ &\quad + \frac{M}{\lambda} \sum_{i=1}^{N_m} \alpha_{m,i}^* K(\mathbf{x}_{m,i}, \mathbf{x}) + b_m^* \end{aligned} \quad (10)$$

### 2.2 大规模问题求解算法

线性方程组(9)包含  $N+M$  个方程, 而且涉及到的矩阵通常为稠密型的, 目前有许多求解线性方程组的成熟方法<sup>[8-10]</sup>。研究表明<sup>[11]</sup>, HS(hestenes-stiefel)共轭梯度法在求解大规模线性方程组具有一定优势, 但该方法要求  $\mathbf{Ax}=\mathbf{b}$  中的矩阵  $\mathbf{A}$  必须是正定对称阵, 但是式(9)中的矩阵  $\mathbf{A}$  是对称阵却不是正定阵。为利用该方法, 类似于 Suykens 等<sup>[4]</sup>, 可等价变换式(9)。为方便, 记式(9)为

$$\begin{bmatrix} 0 & \mathbf{A}^T \\ \mathbf{A} & \mathbf{H} \end{bmatrix} \begin{bmatrix} \boldsymbol{\zeta}_1 \\ \boldsymbol{\zeta}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \end{bmatrix} \quad (11)$$

其中  $\mathbf{H} = \boldsymbol{\Omega} + \gamma^{-1} \mathbf{I}_N + (M/\lambda) \mathbf{B}$ ,  $\boldsymbol{\zeta}_1 = \mathbf{b}$ ,  $\boldsymbol{\zeta}_2 = \boldsymbol{\alpha}$ ,  $\mathbf{d}_1 = 0$ ,  $\mathbf{d}_2 = \mathbf{y}$ 。因为矩阵  $\mathbf{H}$  是正定对称阵, 故式(11)可等价变换为

$$\begin{bmatrix} \mathbf{S} & 0 \\ 0 & \mathbf{H} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\zeta}_1 \\ \mathbf{H}^{-1} \mathbf{A} \boldsymbol{\zeta}_1 + \boldsymbol{\zeta}_2 \end{bmatrix} = \begin{bmatrix} -\mathbf{d}_1 + \mathbf{A}^T \mathbf{H}^{-1} \mathbf{d}_2 \\ \mathbf{d}_2 \end{bmatrix} \quad (12)$$

其中  $\mathbf{S} = \mathbf{A}^T \mathbf{H}^{-1} \mathbf{A}$ , 而且容易证明矩阵  $\mathbf{S}$  是可逆的。从而可得算法如下:

(1) 利用 HS 共轭梯度法求解  $\mathbf{H}\boldsymbol{\eta}=\mathbf{A}$  和  $\mathbf{H}\mathbf{v}=\mathbf{d}_2$ , 并记最优解分别为  $\boldsymbol{\eta}^*$  和  $\mathbf{v}^*$ ;

(2) 计算  $\mathbf{S}=\mathbf{A}^T \boldsymbol{\eta}^*$ ;

(3) 线性方程组(9)的解为:  $\mathbf{b}=\boldsymbol{\zeta}_1=\mathbf{S}^{-1} \boldsymbol{\eta}^{*\top} \mathbf{d}_2$ ,  $\boldsymbol{\alpha}=\boldsymbol{\zeta}_2=\mathbf{v}^*-\boldsymbol{\eta}^* \boldsymbol{\zeta}_1$ 。

## 3 实验方法及实验数据

128 个高粱样品由中国农业科学院品种资源所提

供, 样品用带有 1.0 mm 空径筛网的粉碎机磨碎, 其在 19 个波长点处的近红外漫反射光谱由 LA450 近红外光谱仪测定。该仪器由 BRAN+LUEBBE(德国)公司生产, 具有从 1 445 ~ 2 348 nm 之间的 19 个分隔的滤光片。并对每个样品分别测定了蛋白质、赖氨酸以及淀粉三种成分的含量, 其中蛋白质含量采用国家标准 GB290521982 半微量凯氏定氮法测定; 赖氨酸含量的测定采用国家标准 GB480121984 谷物籽粒赖氨酸测定(染料法); 淀粉含量采用国家标准 BG500621985 谷物籽粒淀粉测定法测定。这样, 对应于三种样品成分含量, 需要同时建模三个任务。

## 4 实验结果及分析

本文从 128 个高粱样品中随机选取 96 个样品组成训练集, 余下的 32 个样品构成测试集。为说明 MTLS-SVR 模型的优势, 以 LS-SVR<sup>[4]</sup>、多因变量 LS-SVR(MLS-SVR)<sup>[12]</sup>和 PLS<sup>[2]</sup>三种模型做参照, 平均相对误差  $\delta$  和相关系数  $R$  作为评价指标。核函数采用径向基(RBF)核:  $K(x, z) = \exp(-\rho \|x - z\|^2)$ ,  $\rho > 0$ 。

### 4.1 参数选择

参数( $\gamma$ ,  $\rho$ ,  $\lambda$ )对 LS-SVR, MLS-SVR 和 MTLS-SVR 模型的性能影响比较大。为了选择一组合适的参数, 本文采用两步格网搜索选参法<sup>[13,14]</sup>: 第一步粗格网搜索, 参数  $\gamma$ ,  $\rho$  和  $\lambda$  按对数增长的方式分别从集合  $\{2^{-5}, 2^{-3}, \dots, 2^{15}\}$ ,  $\{2^{-15}, 2^{-13}, \dots, 2^3\}$  及  $\{2^{-10}, 2^{-8}, \dots, 2^{10}\}$  中取值, 根据总体平均相对误差最小的原则采用留一法确定合适的参数  $\gamma^*$ ,  $\rho^*$  和  $\lambda^*$ 。如果某个参数处于搜索空间的边界, 将相应空间按相同的乘积步长( $2^{\pm 2}$ )增加一个元素, 直到没有进一步的

性能改进为止; 第二步细格网搜索, 首先对  $c \in \{\gamma, \rho, \lambda\}$  构造新的搜索空间:  $\{2^{-1.75} \times c^*, 2^{-1.5} \times c^*, \dots, 2^{1.75} \times c^*\}$ , 然后重复上一步, 将该步确定的值作为最优参数值, 见表 1。对于 MLS-SVR 还需确定一个最优权重向量, 本文采用了安欣等<sup>[12]</sup>给出的方法。对于 PLS 模型, 本文将输入矩阵的行秩与列秩的最小值作为主成分的个数。

**Table 1 Optimal parameters for LS-SVR, MLS-SVR and MTLS-SVR on the broomcorn data set**

|     | LS-SVR             |                     | MLS-SVR            |                     | MTLS-SVR        |                    |                 |
|-----|--------------------|---------------------|--------------------|---------------------|-----------------|--------------------|-----------------|
|     | $\gamma$           | $\rho$              | $\gamma$           | $\rho$              | $\gamma$        | $\rho$             | $\lambda$       |
| 蛋白质 | 2 <sup>27.5</sup>  | 2 <sup>-16.25</sup> | 2 <sup>26.25</sup> | 2 <sup>-16.25</sup> |                 |                    |                 |
| 赖氨酸 | 27.75              | 2 <sup>-1</sup>     | 2 <sup>17.5</sup>  | 2 <sup>-9.5</sup>   | 2 <sup>20</sup> | 2 <sup>-9.75</sup> | 2 <sup>-1</sup> |
| 淀粉  | 2 <sup>17.25</sup> | 2 <sup>-7.25</sup>  | 2 <sup>16.5</sup>  | 2 <sup>-7.25</sup>  |                 |                    |                 |

### 4.2 结果分析

根据表 1 确定的最优参数, 分别利用 PLS, LS-SVR, MLS-SVR 和 MTLS-SVR 建立定量分析模型, 确定高粱样品各成分含量的预测值并计算相应的评价指标, 见表 2。容易看出, MTLS-SVR 的效果最好, MLS-SVR 与 LS-SVR 的性能相当, PLS 的性能最差。本文认为原因有三个方面: (1) MTLS-SVR 同时考虑了各个样品成分的特性( $v_m$ )和共性( $w_0$ ), 而 LS-SVR 只建模了各个样品成分的特性; (2) MLS-SVR 尽管为各个样品成分设置了不同的权重, 但对偶问题仍需分开求解, 因而也只建模了样品成分的特性, 而且由于样品成分权重的引入, 使得参数优化变得更加困难; (3) 本文采用的 PLS 可看作原始空间中的一种线性模型<sup>[15]</sup>, 而近红外光谱数据与样品成分含量间有时存在非线性关系。

**Table 2 Comparisons of the predicted results with PLS, LS-SVR, MLS-SVR and MTLS-SVR quantitative analysis models**

|      | $\delta/\%$ |        |         |          | $R$    |        |         |          |
|------|-------------|--------|---------|----------|--------|--------|---------|----------|
|      | PLS         | LS-SVR | MLS-SVR | MTLS-SVR | PLS    | LS-SVR | MLS-SVR | MTLS-SVR |
| 蛋白质  | 1.77        | 1.52   | 1.52    | 1.52     | 0.9917 | 0.9922 | 0.9923  | 0.9931   |
| 赖氨酸  | 5.75        | 4.98   | 4.91    | 3.04     | 0.7317 | 0.7627 | 0.7725  | 0.8940   |
| 淀粉   | 1.36        | 1.11   | 1.11    | 1.01     | 0.8842 | 0.9223 | 0.9223  | 0.9406   |
| Avg. | 2.96        | 2.54   | 2.51    | 1.85     | 0.8692 | 0.8924 | 0.8957  | 0.9426   |

## 5 结论及讨论

受正则化 MTL 的启发, 本文在 LS-SVR 的基础上设计了一种新的 MTL 模型—MTLS-SVR。将同时建模分析多种样品成分含量的问题转换为 MTL 的问题, 从而可充分利用

各样品成分含量间潜在的联系, 从真正意义上实现了样品多成分含量的同时建模分析。针对样品成分较多的情形, 给出了一种大规模问题求解算法。以高粱样品数据集为实验材料, 通过与 PLS、LS-SVR 以及 MLS-SVR 的对比分析, 验证了 MTLS-SVR 的有效性和高效性, 同时为发展化学计量学提供了一种新的定量分析建模方法。

## References

- [1] YAN Yan-lu, ZHAO Long-lian, HAN Dong-hai, et al(严衍禄, 赵龙莲, 韩东海, 等). Foundation of Near Infrared Spectral Analysis and its Applications(近红外光谱分析基础与应用). Beijing: China Light Industry Press(北京: 中国轻工业出版社), 2005.
- [2] Abdi H. Partial Least Squares (PLS) Regression. Encyclopedia for Research Methods for the Social Sciences, Lewis-Beck M, Bryman A,

- Futing T, eds. Sage, Thousand Oaks, CA, pp. 792.
- [3] Vapnik V N. The Nature of Statistical Learning Theory, 2nd Edition. New York: Springer Verlag, 1999.
  - [4] Suykens J A K, Gestel T V, Brabanter J D, et al. Least Squares Support Vector Machines. Singapore: World Scientific Pub. Co., 2002.
  - [5] Bakker B, Heskes T. Journal of Machine Learning Research, 2003, 4(May): 83.
  - [6] Heskes T. Empirical Bayes for Learning to Learn. Proceedings of the 17th International Conference on Machine Learning (ICML), San Francisco, CA, USA, 2000. 367.
  - [7] Evgeniou T, Pontil M. Regularized Multi-Task Learning. Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), Seattle, WA, USA, 2004. 109.
  - [8] Golub G H, Van Loan C F. Matrix Computations, 3rd Edition. Baltimore and London: Johns Hopkins University Press, 1996.
  - [9] Press W H, Teukolsky S A, Vetterling W T, et al. Numerical Recipes in C: The Art of Scientific Computing. New York: Cambridge University Press, 1992.
  - [10] Saad Y. Iterative Methods for Sparse Linear Systems, 2nd Edition. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, USA, 2003.
  - [11] Hamers B, Suykens J A K, Moor B D. A Comparison of Iterative Methods for Least Squares Vector Machine Classifiers. Internal Report 01-110, 2001. ESAT-SISTA, K. U. Leuven, Leuven, Belgium.
  - [12] AN Xin, XU Shuo, ZHANG Lu-da, et al(安 欣, 徐 硕, 张录达, 等). Spectroscopy and Spectral Analysis (光谱学与光谱分析), 2009, 29(1): 127.
  - [13] Hsu C-W, Chang C-C, Lin C-J. A Practical Guide to Support Vector Classification. Available [online]: <http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>.
  - [14] Xu S, Ma F J, Tao L. Learn from the Information Contained in the False Splice Sites as well as in the True Splice Sites using SVM. Proceedings of the International Conference on Intelligent Systems and Knowledge Engineering (ISKE), Chengdu, China, 2007. 1360.
  - [15] Rosipal R, Trejo L J. Journal of Machine Learning Research, 2001, 2(Dec): 97.

## Multi-Task Least-Squares Support Vector Regression Machines with Applications in NIR Spectral Analysis

XU Shuo<sup>1</sup>, QIAO Xiao-dong<sup>1</sup>, ZHU Li-jun<sup>1</sup>, AN Xin<sup>2</sup>, ZHANG Lu-da<sup>3\*</sup>

1. Information Technology Supporting Center, Institute of Scientific and Information of China, Beijing 100038, China
2. School of International Trade and Economics, University of International Business and Economics, Beijing 100029, China
3. College of Science, China Agricultural University, Beijing 100193, China

**Abstract** In near infrared spectral quantitative analysis, many models consider separately each composition when modeling sample compositions' content, this disregarding the underlying relatedness among sample compositions. To address this problem, the paper views modeling each sample composition's content as a task, thus one can transform the problem that models simultaneously all sample compositions' content to a multi-task learning problem. On the basis of the LS-SVR, a multi-task LS-SVR (MTLS-SVR) model is proposed. Furthermore, an efficient large-scale algorithm is given. The broomcorn samples are taken as experimental material, and corresponding quantitative analysis models are constructed for three sample compositions' content (protein, lysine and starch) with LS-SVR, PLS, multiple dependent variables LS-SVR (MLS-SVR) and MTLS-SVR. For the MTLS-SVR model, the average relative errors between actual values and predicted ones for the three sample compositions' content are 1.52%, 3.04% and 1.01%, respectively, and the correlation coefficients are 0.993 1, 0.894 0 and 0.940 6, respectively. Experimental results show MTLS-SVR model outperforms significantly the three others, which verifies the feasibility and efficiency of the MTLS-SVR model.

**Keywords** Near infrared spectrum; Chemometrics; Multi-task LS-SVR

(Received Oct. 24, 2010; accepted Mar. 9, 2011)

\* Corresponding author