

1-2-2019

## Online Modeling with Semi-Supervised Locally Weighted Partial Least Squares Based on Expectation Maximization Algorithm

Weili Xiong

*1.School of Internet of Things Engineering, Institute of Automation, Jiangnan University, Wuxi 214122, China; ;2.Key Laboratory of Advanced Process Control for Light Industry (Ministry of Education), Jiangnan University, Wuxi 214122, China;*

Mingchen Xue

*1.School of Internet of Things Engineering, Institute of Automation, Jiangnan University, Wuxi 214122, China; ;*

Yanjun Li

*1.School of Internet of Things Engineering, Institute of Automation, Jiangnan University, Wuxi 214122, China; ;*

Follow this and additional works at: <https://dc-china-simulation.researchcommons.org/journal>



Part of the [Artificial Intelligence and Robotics Commons](#), [Computer Engineering Commons](#), [Numerical Analysis and Scientific Computing Commons](#), [Operations Research, Systems Engineering and Industrial Engineering Commons](#), and the [Systems Science Commons](#)

---

This Paper is brought to you for free and open access by Journal of System Simulation. It has been accepted for inclusion in Journal of System Simulation by an authorized editor of Journal of System Simulation.

---

## Online Modeling with Semi-Supervised Locally Weighted Partial Least Squares Based on Expectation Maximization Algorithm

### Abstract

*Abstract: As only small proportion of labeled data can be obtained from chemical processes, an online soft sensing modeling method based on semi-supervised locally weighted partial least squares is proposed. The labeled and unlabeled historical data are accumulated to construct training database. The similarity between the newly measured data and the sample points in the database are calculated and used as the weight of each data point. The semi-supervised locally weighted partial least squares model is constructed, and the Expectation Maximization (EM) algorithm is employed to estimate the parameters of the model. Online prediction achieves simulation results of debutanizer distillation processes, which suggests that the proposed method has good prediction accuracy and stable generalization performance.*

### Keywords

semi-supervised, locally weighted partial least squares, EM algorithm, online modeling

### Recommended Citation

Xiong Weili, Xue Mingchen, Li Yanjun. Online Modeling with Semi-Supervised Locally Weighted Partial Least Squares Based on Expectation Maximization Algorithm[J]. Journal of System Simulation, 2018, 30(1): 8-10.

# 基于EM算法的半监督局部加权PLS在线建模方法

熊伟丽<sup>1,2</sup>, 薛明晨<sup>1</sup>, 李妍君<sup>1</sup>

(1.江南大学 物联网工程学院 自动化研究所, 江苏 无锡 214122; 2.江南大学 轻工过程先进控制教育部重点实验室, 江苏 无锡 214122)

**摘要:** 针对化工过程采样分析获得的有标签样本数量较少的问题, 提出一种基于半监督学习的局部加权偏最小二乘在线软测量建模方法。将过程收集到的有标签及无标签训练样本放入同一数据库中; 对于在线测得的新数据点, 计算其与数据库中各样本点之间的相似度, 将其作为各数据点的权重; 建立半监督局部加权偏最小二乘在线软测量模型, 并采用EM(Expectation Maximization)算法估计模型的参数, 得到模型的在线预测输出。通过对脱丁烷塔过程的仿真研究, 验证了所提方法具有良好的预测精度和泛化性能。

**关键词:** 半监督; 局部加权偏最小二乘; EM算法; 在线建模

中图分类号: TP391

文献标识码: A

文章编号: 1004-731X (2018) 01-0008-10

DOI: 10.16182/j.issn1004731x.joss.201801002

## Online Modeling with Semi-Supervised Locally Weighted Partial Least Squares Based on Expectation Maximization Algorithm

Xiong Weili<sup>1,2</sup>, Xue Mingchen<sup>1</sup>, Li Yanjun<sup>1</sup>

(1. School of Internet of Things Engineering, Institute of Automation, Jiangnan University, Wuxi 214122, China;

2. Key Laboratory of Advanced Process Control for Light Industry (Ministry of Education), Jiangnan University, Wuxi 214122, China)

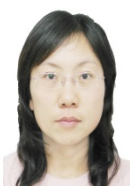
**Abstract:** As only small proportion of labeled data can be obtained from chemical processes, an online soft sensing modeling method based on semi-supervised locally weighted partial least squares is proposed. The labeled and unlabeled historical data are accumulated to construct training database. The similarity between the newly measured data and the sample points in the database are calculated and used as the weight of each data point. The semi-supervised locally weighted partial least squares model is constructed, and the Expectation Maximization (EM) algorithm is employed to estimate the parameters of the model. Online prediction achieves simulation results of debutanizer distillation processes, which suggests that the proposed method has good prediction accuracy and stable generalization performance.

**Keywords:** semi-supervised; locally weighted partial least squares; EM algorithm; online modeling

## 引言

在实际化工生产过程中, 由于操作条件复杂、工况变化快或测量技术的局限, 一些与产品质量密

切相关的过程变量的在线测量往往很困难, 软测量技术可以用来解决某些复杂工业生产过程中质量变量的在线预测, 所以建立高精度的软测量模型显得十分重要。目前常采用的软测量建模方法主要有基于机理模型的方法、基于数据驱动以及两者相结合的混合建模方法。基于机理模型的建模方法需要了解过程的详细信息, 运用物料守恒、能量守恒、化学反应动力学等原理, 对对象的操作过程和工艺



收稿日期: 2015-11-17 修回日期: 2016-01-30;  
基金项目: 国家自然科学基金(61773182), 江苏省“六大人才高峰”计划(2013-DZXX-043);  
作者简介: 熊伟丽(1978-), 女, 河南洛阳, 博士, 教授, 硕导, 研究方向为复杂工业过程建模及优化, 智能优化算法及应用。

<http://www.china-simulation.com>

原理进行分析, 建立机理模型, 对于复杂的化工过程而言, 这些是很难获得的, 建立相应的机理模型也变得非常困难, 由于通过 DCS 系统获得的过程数据越来越多, 基于数据驱动的建模方法得到越来越多的关注。

近几十年来, 基于数据驱动的软测量模型无论在工业领域或学术领域均受到了广泛的关注。常用的建模方法有多元线性回归 (Multiple Linear Regression, MLR)<sup>[1]</sup>、主元回归法 (Principal Component Regression, PCR)<sup>[2]</sup>、部分最小二乘 (Partial Least Squares, PLS)<sup>[3-4]</sup> 等, 这些均为线性建模方法, 对于非线性严重的过程, 常利用非线性建模方法如支持向量机 (Support Vector Machines, SVM)<sup>[5]</sup>、非线性 PLS<sup>[6]</sup>、局部加权回归 (Locally Weighted Regression, LWR)<sup>[7]</sup> 等。局部加权偏最小二乘 (Locally Weighted Partial Least Squares, LWPLS) 是基于即时 (Just-In-Time, JIT) 建模思想在传统 PLS 上的改进<sup>[8]</sup>, 根据数据库中不同样本点与新数据点之间的相似度, 建立局部模型时考虑它们的权重, 从而减小了相似度较低的样本对建模的影响, 提高了建模精度。另外, 由于局部加权偏最小二乘建模方法是对每一在线测得的数据均建立局部模型, 用局部线性模型来逼近非线性过程, 所以能够较好地描述非线性过程的动态特征。

但是, 传统的基于数据驱动的软测量建模方法往往需要既有输入又有输出的样本数据来进行建模, 即有标签样本, 将只有输入的样本数据, 即无标签样本丢弃, 过程数据的信息得不到充分利用, 建立的模型性能往往也得不到保证。而对应品质变量的输出样本一般很难测量, 这样会导致采集到的输入样本与输出样本数量不对应, 在获得大量输入样本数据的同时获得的输出样本数据仅占小部分。针对此类问题, 章军等<sup>[9]</sup>利用无标签样本中的支持向量, 训练出来的 SVM 回归模型具有更好的估计精度和泛化性能; 葛志强等<sup>[10]</sup>提出基于半监督学习思想的 SBPCR 建模方法, 充分利用有标签样本及无标签样本, 提高了模型的预测精度及鲁棒性;

吴继明等<sup>[11]</sup>提出一种基于 Boosting 思想的半监督学习算法 IMSB, 解决了多分类问题, 并在训练样本集中有标签样本比较少的环境下取得较好的效果。

结合实际化工过程, 为了解决有标签样本数量较少, 充分利用无标签样本包含的过程信息, 本文提出一种基于半监督学习的局部加权偏最小二乘在线软测量建模方法, 利用有标签样本的同时, 充分利用无标签样本所携带的信息, 建立半监督局部加权偏最小二乘模型, 并采用 EM 算法求取实时更新的半监督局部加权偏最小二乘模型的参数, 得到模型的在线预测输出。通过对脱丁烷塔过程的仿真研究, 结果表明本文所提方法可以实现对工业过程中关键变量的有效预测。

## 1 局部加权偏最小二乘

LWPLS 是一种基于 JIT 建模思想的建模方法, 在建立局部模型时, 考虑了新来数据与数据库中各样本点之间的相似度, 计算出它们之间的相似度值并将其作为各数据点的权重, 继而建立局部模型。因此, 相比 PLS 建模而言, LWPLS 能够更好描述化工过程中的非线性特征。假设存入数据库中的输入输出变量矩阵为  $\mathbf{X} \in R^{N \times M}$ ,  $\mathbf{Y} \in R^{N \times R}$ , 其中第  $i$  个输入输出样本分别为  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{iM}]^T$  及  $\mathbf{y}_i = [y_{i1}, y_{i2}, \dots, y_{iR}]^T$ ,  $i = 1, 2, \dots, N$ ,  $N$  为样本个数;  $M$  和  $R$  分别为输入、输出样本的维度。PLS 模型可以表示成式(1)形式:

$$\begin{cases} \mathbf{X} = \mathbf{TP}^T + \mathbf{E}_X \\ \mathbf{Y} = \mathbf{TC}^T + \mathbf{F}_Y \end{cases} \quad (1)$$

这里,  $\mathbf{P} \in R^{M \times K}$  为负荷矩阵,  $\mathbf{T} \in R^{N \times K}$  为主成分矩阵,  $K$  为选择的主成分的个数,  $\mathbf{C} \in R^{R \times K}$  为回归矩阵,  $\mathbf{E}_X \in R^{N \times M}$  及  $\mathbf{E}_Y \in R^{N \times R}$  分别对应输入输出变量矩阵的残差矩阵。

对于新的数据点  $\mathbf{x}_q = [x_{q1}, x_{q2}, \dots, x_{qM}]^T$ , 首先根据式(2)计算出它与数据库中各样本点之间的相似度:

$$\begin{cases} \omega_i = e^{-\frac{d_i}{\sigma_d \varphi}} \\ d_i = \sqrt{(\mathbf{x}_i - \mathbf{x}_q)^T (\mathbf{x}_i - \mathbf{x}_q)} \end{cases} \quad (2)$$

式中,  $\sigma_d$  是  $d_i$  的标准差,  $\varphi$  是位置参数, 一般取 0.1~1.5 之间<sup>[12]</sup>。将由(2)式计算得到的  $\omega_i$  表示成相似矩阵  $\boldsymbol{\Omega} \in R^{N \times N}$  的形式为:

$$\boldsymbol{\Omega} \in \text{diag}(\omega_1, \omega_2, \dots, \omega_N) \quad (3)$$

LWPLS 的算法步骤如下:

Step1: 确定主成分  $K$  的值, 并设置其在算法中的初值  $k=1$ 。 $K$  的值通常可通过主元分析法(PCA)来确定。

Step2: 根据公式(2)、(3)计算出相似矩阵  $\boldsymbol{\Omega}$ 。

Step3: 对输入输出矩阵及新来数据点进行预处理, 计算出  $\mathbf{X}_k, \mathbf{Y}_k$  及  $\mathbf{x}_{qk}$

$$\begin{cases} \mathbf{X}_k = \mathbf{X} - \mathbf{I}_N [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_M] \\ \mathbf{Y}_k = \mathbf{Y} - \mathbf{I}_N [\bar{y}_1, \bar{y}_2, \dots, \bar{y}_R] \\ \mathbf{x}_{qk} = \mathbf{x}_q - [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_M]^T \end{cases} \quad (4)$$

$$\bar{x}_m = \frac{\sum_{i=1}^N \omega_i x_{im}}{\sum_{i=1}^N \omega_i}, \quad \bar{y}_r = \frac{\sum_{i=1}^N \omega_i y_{ir}}{\sum_{i=1}^N \omega_i} \quad (5)$$

其中  $\mathbf{I}_N \in R^N$  为全 1 的列向量,  $m=1, 2, \dots, M$ ,  $r=1, 2, \dots, R$ 。

Step4: 推断出局部线性模型

$$\mathbf{X} \text{ 的主成分向量: } \mathbf{t}_k = \mathbf{X}_k \mathbf{w}_k \quad (6)$$

$$\mathbf{X} \text{ 的负荷向量: } \mathbf{p}_k = \frac{\mathbf{X}_k^T \boldsymbol{\Omega} \mathbf{t}_k}{\mathbf{t}_k^T \boldsymbol{\Omega} \mathbf{t}_k} \quad (7)$$

$$\text{模型回归系数向量: } \mathbf{c}_k = \frac{\mathbf{Y}_k^T \boldsymbol{\Omega} \mathbf{t}_k}{\mathbf{t}_k^T \boldsymbol{\Omega} \mathbf{t}_k} \quad (8)$$

$$\text{新来数据点的得分: } t_{qk} = \mathbf{x}_{qk}^T \mathbf{w}_k \quad (9)$$

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \mathbf{t}_k \mathbf{p}_k^T \quad (10)$$

$$\mathbf{Y}_{k+1} = \mathbf{Y}_k - \mathbf{t}_k \mathbf{c}_k^T \quad (11)$$

$$\mathbf{x}_{q(k+1)} = \mathbf{x}_{qk} - t_{qk} \mathbf{p}_k \quad (12)$$

如果  $k=K$ , 则转入下一步; 否则, 使  $k=k+1$  并转到 Step4。这里,  $\mathbf{w}_k$  是  $\mathbf{X}_k^T \boldsymbol{\Omega} \mathbf{Y}_k \mathbf{Y}_k^T \boldsymbol{\Omega} \mathbf{X}_k$  的最大特征值所对应的特征向量。

Step5: 通过上面步骤得到输入样本  $\mathbf{X}$  的得分矩阵

$\mathbf{T}=[t_1, t_2, \dots, t_k]$ , 负荷矩阵  $\mathbf{P}=[\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k]$ , 模型回归系数矩阵  $\mathbf{C}=[\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k]$  以及新来数据点的得分向量  $\mathbf{t}=[t_{q1}, t_{q2}, \dots, t_{qk}]$ , 可以计算新来数据点的输出估计值  $\hat{\mathbf{y}}_q$ :

$$\hat{\mathbf{y}}_q = \mathbf{t} \mathbf{C}^T = [t_{q1}, t_{q2}, \dots, t_{qk}] [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k]^T \quad (13)$$

## 2 基于局部加权偏最小二乘回归的半监督学习

### 2.1 半监督局部加权偏最小二乘模型

对实际工业过程中获得的无标签样本, 传统的 LWPLS 方法无法利用这些样本来进行建模, 常见的处理方法是将这些无标签样本丢弃, 这样可能会影响模型的预测精度和泛化性能。半监督学习可以在利用有标签样本的同时, 利用无标签样本去增强机器学习能力, 以提高模型的性能<sup>[13-14]</sup>。对于每一个新来样本数据, 基于半监督学习的 LWPLS 模型可以表示为<sup>[15]</sup>:

$$\begin{cases} \mathbf{x}_i = \mathbf{P} \mathbf{t}_i + \mathbf{e}_{xi} \\ \mathbf{y}_j = \mathbf{C} \mathbf{t}_j + \mathbf{f}_{yj} \end{cases} \quad (14)$$

其中  $i=1, 2, \dots, n$ ,  $j=1, 2, \dots, n_1$ ;  $n_1$  为数据库中有标签的样本的个数,  $n_2=n-n_1$  为数据库中无标签样本的个数。 $\mathbf{P} \in R^{m \times K}$ ,  $\mathbf{C} \in R^{r \times K}$  分别表示对应的负荷矩阵及回归矩阵, 其中  $m$  为输入变量的个数,  $r$  为输出变量的个数,  $\mathbf{t} \in R^{K \times 1}$  为主成分向量,  $\mathbf{e}_x \in R^{m \times 1}$ ,  $\mathbf{f}_y \in R^{r \times 1}$  分别为过程噪声的输入输出变量。在半监督 LWPLS 模型中, 一般假设主成分及过程噪声的概率密度函数均服从高斯分布, 即  $p(\mathbf{t}) = N(0, \mathbf{I})$ ,  $p(\mathbf{e}) = N(0, \sigma_x^2 / \omega_i \mathbf{I})$ ,  $p(\mathbf{f}) = N(0, \sigma_y^2 / \omega_j \mathbf{I})$ , 其中  $\mathbf{I}$  为单位矩阵,  $\sigma_x^2 / \omega_i$ ,  $\sigma_y^2 / \omega_j$  分别为输入变量和输出变量的噪声方差,  $\omega_i, \omega_j$  越小, 则表明包含的噪声信息越多, 对模型的贡献越小。

### 2.2 基于 EM 算法估计半监督 LWPLS 模型参数

#### 2.2.1 期望最大化算法

期望最大化 (Expectation-Maximization, EM)

算法是由 Dempster 提出的一种由已知数据估计未知参数的迭代方法, 旨在求解数据有缺失或隐含情况下参数的极大似然估计<sup>[16-18]</sup>, 适用于解决不完全数据建模问题。该算法由 E 步和 M 步组成, 在 E 步中根据隐含变量的后验分布计算条件期望, 即 Q 函数; 在 M 步中根据 E 步得到的条件期望求解参数的极大似然估计。E 步和 M 步间交替进行计算, 直到参数的估计值收敛。

### 2.2.2 EM 算法在半监督 LWPLS 模型参数估计中的应用

E 步: 选取主成分向量  $\{t_i\}_{i=1}^n$  作为隐含变量, 即  $C_{miss} = \{T\} = \{t_i, i=1, 2, \dots, n\}$ ; 可观测数据集为  $C_{obs} = \{X, Y\}$ ; 则完整数据集可定义为  $C = \{C_{obs}, C_{miss}\} = \{X, Y, T\}$ , 模型参数集为  $\Theta = \{P, C, \sigma_x^2, \sigma_y^2\}$ , 根据 EM 算法, 首先计算 Q 函数为:

$$\begin{aligned} Q(\Theta|\Theta^S) &= E_{C_{miss}|C_{obs}, \Theta^S} \ln P(C_{obs}, C_{miss}|\Theta) = \\ &E_{T|X, Y, \Theta^S} \ln P(X, Y, T|\Theta) = \\ &E_{T|X, Y, \Theta^S} \sum_{j=1}^{n_1} \ln P(x_j|t_j, \Theta) + E_{T|X, Y, \Theta^S} \sum_{j=1}^{n_2} \ln P(x_{j'}|t_{j'}, \Theta) + \\ &E_{T|X, Y, \Theta^S} \sum_{j=1}^{n_1} \ln P(y_j|t_j, \Theta) + E_{T|X, Y, \Theta^S} \sum_{j=1}^{n_2} \ln P(t_j|\Theta) + \\ &E_{T|X, Y, \Theta^S} \sum_{j=1}^{n_2} \ln P(t_{j'}|\Theta) = \\ &\int_T \sum_{j=1}^{n_1} \ln P(x_j|t_j, \Theta) P(T|X, Y, \Theta^S) dT + \\ &\int_T \sum_{j=1}^{n_2} \ln P(x_{j'}|t_{j'}, \Theta) P(T|X, Y, \Theta^S) dT + \\ &\int_T \sum_{j=1}^{n_1} \ln P(y_j|t_j, \Theta) P(T|X, Y, \Theta^S) dT + \\ &\int_T \sum_{j=1}^{n_2} \ln P(t_j|\Theta) P(T|X, Y, \Theta^S) dT + \\ &\int_T \sum_{j=1}^{n_2} \ln P(t_{j'}|\Theta) P(T|X, Y, \Theta^S) dT \end{aligned} \quad (15)$$

将式(15)最后一个等式等号右端第一项进一

步展开可得:

$$\begin{aligned} &\int_T \sum_{j=1}^{n_1} \ln P(x_j|t_j, \Theta) P(T|X, Y, \Theta^S) dT = \\ &\sum_{j=1}^{n_1} \int P(t_j|x_j, y_j, \Theta^S) \ln P(x_j|t_j, \Theta) dt_j = \\ &\sum_{j=1}^{n_1} \int P(t_j|x_j, y_j, \Theta^S) \left\{ \ln(2\pi) \frac{m}{2} \left| \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right|^{\frac{1}{2}} - \right. \\ &\left. \frac{1}{2} (x_j - P t_j)^T \left( \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right)^{-1} (x_j - P t_j) \right\} dt_j = \\ &\sum_{j=1}^{n_1} \left\{ -\frac{m}{2} \ln 2\pi - \frac{1}{2} \ln \left| \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right| \right\} - \\ &\sum_{j=1}^{n_1} \frac{\omega_j}{2\sigma_x^2} \int (x_j^T x_j - x_j^T P t_j - t_j^T P^T x_j + t_j^T P^T P t_j) \cdot \\ &P(t_j|x_j, y_j, \Theta^S) dt_j = \\ &\sum_{j=1}^{n_1} \left\{ -\frac{m}{2} \ln 2\pi - \frac{1}{2} \ln \left| \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right| \right\} - \\ &\sum_{j=1}^{n_1} \frac{\omega_j}{2\sigma_x^2} \left\{ x_j^T x_j - \int x_j^T P t_j P(t_j|x_j, y_j, \Theta^S) dt_j - \right. \\ &\left. \int t_j^T (x_j^T P)^T P(t_j|x_j, y_j, \Theta^S) dt_j + \right. \\ &\left. \int t_j^T P^T P t_j P(t_j|x_j, y_j, \Theta^S) dt_j \right\} = \\ &\sum_{j=1}^{n_1} \left\{ -\frac{m}{2} \ln 2\pi - \frac{1}{2} \ln \left| \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right| \right\} - \sum_{j=1}^{n_1} \frac{\omega_j}{2\sigma_x^2} \left\{ x_j^T x_j - \right. \\ &2x_j^T P E(t_j|x_j, y_j, \Theta^S) + \\ &\left. \text{trace} \left[ P^T P E(t_j t_j^T|x_j, y_j, \Theta^S) \right] \right\} = \\ &\sum_{j=1}^{n_1} \left\{ -\frac{m}{2} \ln 2\pi - \frac{1}{2} \ln \left| \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right| \right\} - \sum_{j=1}^{n_1} \frac{\omega_j}{2\sigma_x^2} \left\{ x_j^T x_j - \right. \\ &2x_j^T P E(t_j|x_j, y_j, \Theta^S) + \\ &\left. \text{trace} \left[ P^T P E(t_j t_j^T|x_j, y_j, \Theta^S) \right] \right\} \end{aligned} \quad (16)$$

同理, 将公式(15)最后一个等式等号右端第

二、三、四、五展开分别可得:

$$\int_T \sum_{j=1}^{n_2} \ln P(\mathbf{x}_j | \mathbf{t}_j, \boldsymbol{\theta}) P(T | X, Y, \boldsymbol{\theta}^S) dT = \sum_{j=1}^{n_2} \left\{ -\frac{m}{2} \ln 2\pi - \frac{1}{2} \ln \left| \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right| \right\} - \sum_{j=1}^{n_2} \frac{\omega_j}{2\sigma_x^2} \left\{ \mathbf{x}_j^T \mathbf{x}_j - 2 \mathbf{x}_j^T \mathbf{P} \mathbf{E}(\mathbf{t}_j | \mathbf{x}_j, \boldsymbol{\theta}^S) + \text{trace}[\mathbf{P}^T \mathbf{P} \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \boldsymbol{\theta}^S)] \right\} \quad (17)$$

$$\int_T \sum_{j=1}^{n_1} \ln P(\mathbf{y}_j | \mathbf{t}_j, \boldsymbol{\theta}) P(T | X, Y, \boldsymbol{\theta}^S) dT = \sum_{j=1}^{n_1} \left\{ -\frac{r}{2} \ln 2\pi - \frac{1}{2} \ln \left| \frac{\sigma_y^2}{\omega_j} \mathbf{I} \right| \right\} - \sum_{j=1}^{n_1} \frac{\omega_j}{2\sigma_y^2} \left\{ \mathbf{y}_j^T \mathbf{y}_j - 2 \mathbf{y}_j^T \mathbf{C} \mathbf{E}(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) + \text{trace}[\mathbf{C}^T \mathbf{C} \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S)] \right\} \quad (18)$$

$$\int_T \sum_{j=1}^{n_1} \ln P(\mathbf{t}_j | \boldsymbol{\theta}) P(T | X, Y, \boldsymbol{\theta}^S) dT = \sum_{j=1}^{n_1} \left\{ -\frac{K}{2} \ln 2\pi - \frac{1}{2} \text{trace}[\mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S)] \right\} \quad (19)$$

$$\int_T \sum_{j=1}^{n_2} \ln P(\mathbf{t}_j | \boldsymbol{\theta}) P(T | X, Y, \boldsymbol{\theta}^S) dT = \sum_{j=1}^{n_2} \left\{ -\frac{K}{2} \ln 2\pi - \frac{1}{2} \text{trace}[\mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S)] \right\} \quad (20)$$

M 步: 根据 E 步化简得到的 Q 函数, 分别对模型参数  $\boldsymbol{\theta}$  求偏导并置 0, 得到:

$$\frac{\partial Q(\boldsymbol{\theta} | \boldsymbol{\theta}^S)}{\partial \mathbf{P}} = 0 \Rightarrow \mathbf{P}^{\text{new}} = \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{x}_j \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) + \sum_{j=1}^{n_2} \omega_j \mathbf{x}_j \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \boldsymbol{\theta}^S) \right] \cdot \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) + \sum_{j=1}^{n_2} \omega_j \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \boldsymbol{\theta}^S) \right]^{-1} \quad (21)$$

$$\frac{\partial [Q(\boldsymbol{\theta} | \boldsymbol{\theta}^S)]}{\partial \mathbf{C}} = 0 \Rightarrow \mathbf{C}^{\text{new}} = \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{y}_j \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) \right] \cdot \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) \right]^{-1} \quad (22)$$

$$\frac{\partial [Q(\boldsymbol{\theta} | \boldsymbol{\theta}^S)]}{\partial \sigma_x^2} = 0 \Rightarrow (\sigma_x^2)^{\text{new}} = \frac{1}{m \left( \sum_{j=1}^{n_1} \omega_j + \sum_{j=1}^{n_2} \omega_j \right)} \left\{ \sum_{j=1}^{n_1} \omega_j \mathbf{x}_j^T \mathbf{x}_j + \sum_{j=1}^{n_2} \omega_j \mathbf{x}_j^T \mathbf{x}_j + \text{trace} \left\{ \mathbf{P}^{\text{newT}} \mathbf{P}^{\text{new}} \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) + \sum_{j=1}^{n_2} \omega_j \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \boldsymbol{\theta}^S) \right] \right\} - 2 \text{trace} \left\{ \mathbf{P}^{\text{newT}} \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{x}_j \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) + \sum_{j=1}^{n_2} \omega_j \mathbf{x}_j \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \boldsymbol{\theta}^S) \right] \right\} \right\} \quad (23)$$

$$\frac{\partial [Q(\boldsymbol{\theta} | \boldsymbol{\theta}^S)]}{\partial \sigma_y^2} = 0 \Rightarrow (\sigma_y^2)^{\text{new}} = \frac{1}{r \sum_{j=1}^{n_1} \omega_j} \left\{ \sum_{j=1}^{n_1} \omega_j \mathbf{y}_j^T \mathbf{y}_j + \text{trace} \left\{ \mathbf{C}^{\text{newT}} \mathbf{C}^{\text{new}} \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) \right] \right\} - 2 \text{trace} \left\{ \mathbf{C}^{\text{newT}} \left[ \sum_{j=1}^{n_1} \omega_j \mathbf{y}_j \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}^S) \right] \right\} \right\} \quad (24)$$

其中, 对于有标签样本以及无标签样本, 分别有<sup>[19]</sup>:

$$\mathbf{E}(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}) = \left( \frac{\omega_j}{\sigma_x^2} \mathbf{P}^T \mathbf{P} + \frac{\omega_j}{\sigma_y^2} \mathbf{C}^T \mathbf{C} + \mathbf{I} \right)^{-1} \cdot \left( \frac{\omega_j}{\sigma_x^2} \mathbf{P}^T \mathbf{x}_j + \frac{\omega_j}{\sigma_y^2} \mathbf{C}^T \mathbf{y}_j \right)$$

$$\mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}) = \left( \frac{\omega_j}{\sigma_x^2} \mathbf{P}^T \mathbf{P} + \frac{\omega_j}{\sigma_y^2} \mathbf{C}^T \mathbf{C} + \mathbf{I} \right)^{-1} + \mathbf{E}(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}) \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \mathbf{y}_j, \boldsymbol{\theta}) \quad (25)$$

$$\mathbf{E}(\mathbf{t}_j | \mathbf{x}_j, \boldsymbol{\theta}) = \left( \mathbf{P}^T \mathbf{P} + \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right)^{-1} \mathbf{P}^T \mathbf{x}_j$$

$$\mathbf{E}(\mathbf{t}_j \mathbf{t}_j^T | \mathbf{x}_j, \boldsymbol{\theta}) = \frac{\sigma_x^2}{\omega_j} \left( \mathbf{P}^T \mathbf{P} + \frac{\sigma_x^2}{\omega_j} \mathbf{I} \right)^{-1} + \mathbf{E}(\mathbf{t}_j | \mathbf{x}_j, \boldsymbol{\theta}) \mathbf{E}^T(\mathbf{t}_j | \mathbf{x}_j, \boldsymbol{\theta}) \quad (26)$$

将  $\theta$  代入 E 步重新计算期望值, 如此 E 步、M 步反复循环进行, 直至  $\|\theta - \theta^s\|$  收敛为止。

对于新来数据点  $x_q$ , 模型的输出估计为:

$$\hat{y}_q = C(P^T P + \sigma_x^2 I)^{-1} P^T x_q \quad (27)$$

### 3 基于 EM 算法的半监督局部加权偏最小二乘建模算法步骤

本文基于半监督 LWPLS 算法进行软测量建模, 该方法能充分融合无标签样本所携带的信息, 从而避免丢弃无标签样本导致训练所得模型预测精度不高等缺陷。对一个在线测得的新数据点  $x_q$  的建模流程如图 1 所示, 算法步骤描述如下:

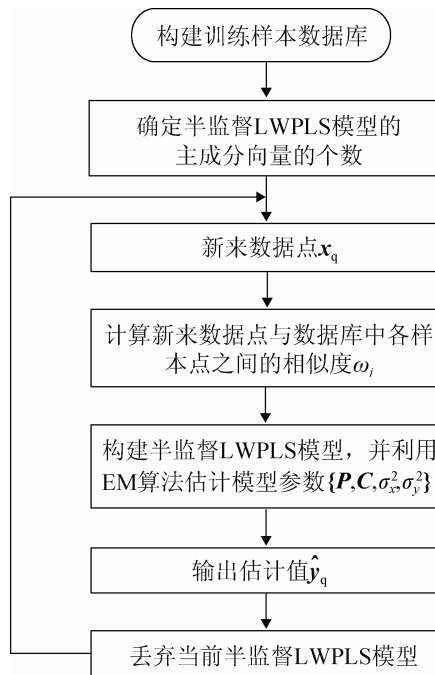


图 1 基于半监督 LWPLS 的在线软测量建模流程

Fig. 1 Online soft sensor modeling process based on semi-supervised LWPLS

Step 1: 将实际过程中采集到的有标签样本  $\{(x_j, y_j), x_j \in R^m, y_j \in R^r, j=1, 2, \dots, n_1\}$  及无标签样本  $\{x_{j'}, x_{j'} \in R^m, j'=1, 2, \dots, n_2\}$  放入数据库中;

Step 2: 利用 PCA 确定需要选取的半监督 LWPLS 模型主成分向量的个数  $K$  (即选取的半监督 LWPLS 模型投影的个数)。PCA 可以对高维数据进行降维, 通过累计方差贡献率准则得到模型的  $K$

个投影;

Step 3: 对于新数据点  $x_q$ , 根据式(2)计算其与数据库中各样本点之间的相似度  $\omega_i (i=1, 2, \dots, n)$ ;

Step 4: 根据式(14)构建半监督 LWPLS 模型, 并利用 EM 算法估计半监督 LWPLS 模型参数  $\theta = \{P, C, \sigma_x^2, \sigma_y^2\}$ , 即本文公式(15)~(26)部分;

Step 5: 通过公式(27)得到新建立的半监督模型在线预测输出  $\hat{y}_q$ ;

Step 6: 将均方根误差(RMSE)作为性能指标, 评价本文建模方法的预测性能, 其计算公式如式(28)所示:

$$RMSE = \sqrt{\frac{1}{n-1} \sum_{q=1}^n |y_q - \hat{y}_q|^2} \quad (28)$$

式中:  $y_q$  为真值;  $\hat{y}_q$  为模型预测值;  $n$  为测得的数据点的个数。

### 4 仿真实例

脱丁烷塔是石油炼制生产过程中脱硫和石脑油分离装置的重要组成部分<sup>[20]</sup>, 其中脱丁烷塔塔底的丁烷浓度对石油炼制影响较大且难以直接测得, 需要对其进行严格的监测与控制。为了提高脱丁烷塔的控制质量, 实时监测塔底丁烷浓度, 对丁烷浓度的实时估计显示十分重要。该工艺过程如图 2 所示。实际工业过程中选取 7 个测量变量作为辅助变量, 主导变量为脱丁烷塔塔底丁烷浓度, 选取的辅助变量如表 1 所示。

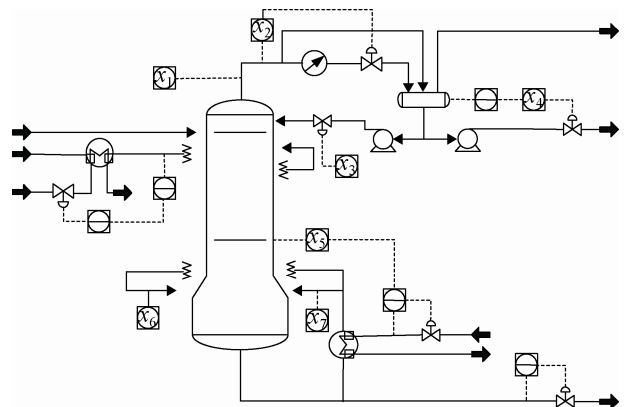


图 2 脱丁烷塔工艺流程图

Fig.2 Diagram of Debutanizer column process



表1 脱丁烷塔过程输入变量  
Tab.1 Input variables for debutanizer column process

Input variables	Description
$x_1$	Top temperature
$x_2$	Top pressure
$x_3$	Reflux flow
$x_4$	Flow to next process
$x_5$	6 <sup>th</sup> tray temperature
$x_6$	Bottom temperature
$x_7$	Bottom temperature

为了预测该过程丁烷浓度,通过脱丁烷塔过程平台,共收集2000组数据,其中1000组作为训练样本集,1000组作为测试样本集;同时,训练样本集又被分为有标签样本集及无标签样本集。为了验证本文所提基于半监督学习建模方法的有效性,利用本文方法对有标签训练样本集占总训练样本集的10%进行建模研究,在1000个测试样本集中每隔10个样本选一个作为有标签样本,构成100组有标签样本集,即有输入又有输出的样本数据,其余900组作为无标签样本集,即只有输入的样本数据,则有标签样本集可表示为:  $\{(x_j, y_j), x_j \in R^7, y_j \in R, j=1,2,\dots,100\}$ , 无标签样本集可表示为:  $\{x_{j'}, x_{j'} \in R^7, j'=1,2,\dots,900\}$ 。为显示有标签训练样本集占总训练样本集的比例,图3给出了总训练样本集中输出变量与样本个数之间的效果图。其中,绿点为有标签样本,红点为无标签样本,这里为突出显示有标签样本与无标签样本的效果,将无标签样本的输出暂时补上,在建模时,红点部分只有相对应的输入变量。

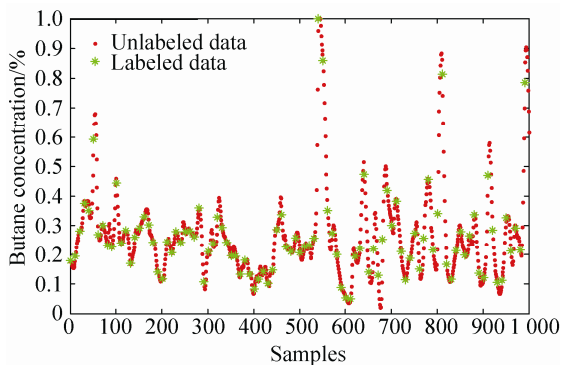


图3 有标签训练样本与无标签训练样本  
Fig.3 Labeled and unlabeled training samples

利用PCA确定需要选取的半监督LWPLS模型主成分向量的个数 $K$ ,通过累计方差贡献率准则确定主成分个数以保证这些提取的主成分能够解释过程数据90%以上的信息,根据该准则提取出5个主成分,即 $K=5$ ,如图4所示。

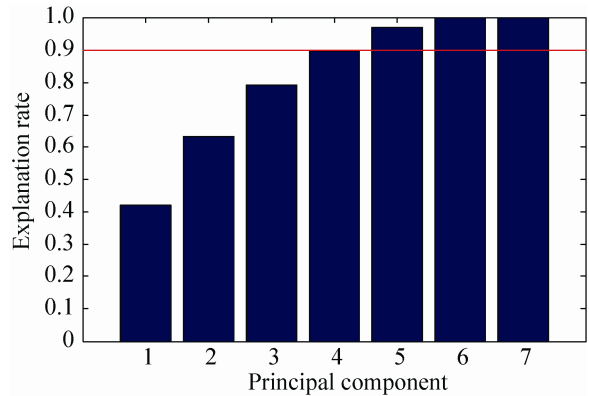


图4 PCA特征提取  
Fig.4 Feature extraction of PCA

基于EM算法的半监督LWPLS建模方法是一种在线建模方法,本仿真共有1000组测试样本,对于其中一组新来数据点 $x_q$ ,首先根据式(2)计算出其与数据库中有标签训练样本集及无标签训练样本集各样本点之间的相似度 $\omega_j(j=1,2,\dots,100)$ 及 $\omega_{j'}(j'=1,2,\dots,900)$ ;然后,将其带入公式(21)~(24)中不断迭代,公式中迭代算法参数的初值 $\theta = \{P, C, \sigma_x^2, \sigma_y^2\}$ 可以任意设定,当 $\|\theta - \theta^s\|$ 收敛时可求出半监督LWPLS模型参数 $\theta = \{P, C, \sigma_x^2, \sigma_y^2\}$ ;根据推导出的半监督LWPLS模型的参数 $\theta$ ,代入公式(27),即可得到模型的输出估计 $\hat{y}_q$ 。其中,半监督LWPLS模型中的位置参数 $\varphi$ 通过试验法来确定, $\varphi$ 一般取0.1~1.5之间<sup>[12]</sup>,通过不断调整位置参数,当 $\varphi$ 调整为0.2时建立的在线模型预测效果较佳。

为验证本文基于半监督学习建模方法的有效性,将本文建模方法与LWPLS建模方法进行比较,如图5所示。传统的LWPLS建模方法只使用有标签样本,将无标签样本丢弃,过程数据的信息得不到充分利用。由图5可见,基于EM算法的半监督LWPLS在线建模方法的预测结果与真值的吻合程度更高,预测效果更好。

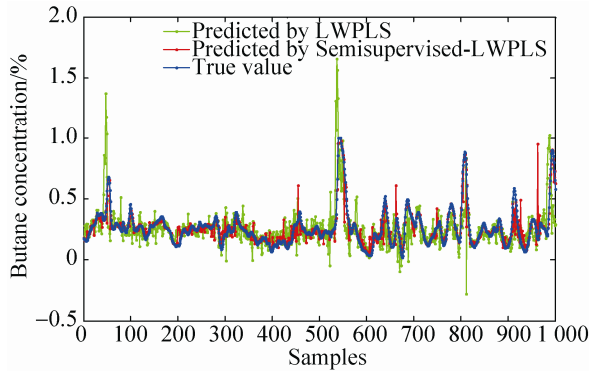


图 5 LWPLS 与半监督 LWPLS 模型对丁烷浓度的预测结果

Fig. 5 Prediction of butane concentration between LWPLS and semi-supervised LWPLS Model

同时,将本文建模方法与 PLS 建模方法(文献[2])及 SBPCR 建模方法(文献[10])进行比较来验证本文建模方法的预测精度。在有标签训练样本集占总训练样本集 10%情况下,3 种不同软测量建模方法的预测结果如图 6 所示。

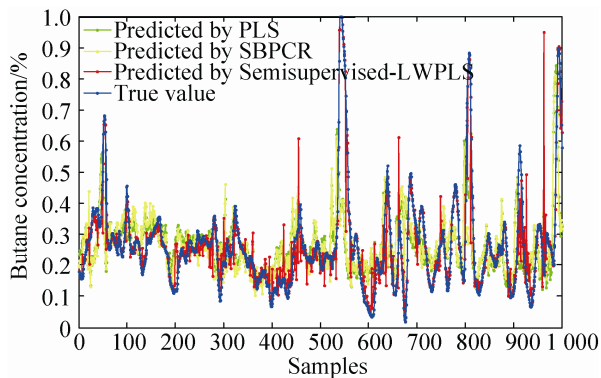


图 6 3 种建模方法对丁烷浓度的预测结果

Fig. 6 Prediction results of Butane concentration by three modeling methods

由图 6 可以看出,基于 EM 算法的半监督局部加权偏最小二乘在线建模方法的预测结果与丁烷浓度的真值基本一致,而采用 PLS 模型以及 SBPCR 模型的预测结果则相对较差。为了进一步说明本文方法的有效性,图 7 给出了有标签训练样本集占总训练样本集 10%情况下丁烷浓度真值与预测值的散点图,图中散点越接近黑色对角线说明预测值更接近真值,模型精度越高。由图 7 可以看出,本文建模方法得到的散点分布更靠近黑色对

角线,更加紧凑地分布在对角线两侧,而其它两种建模方法的散点分布则相对分散,说明本文建模方法得到的丁烷浓度预测值更接近真值,模型预测精度更高。

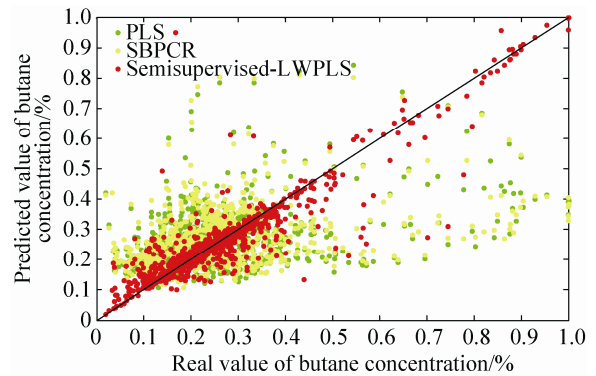


图 7 丁烷浓度真值与预测值散点分布

Fig. 7 Scatters of true and predicted values of Butane concentration

图 8 给出了 3 种建模方法的预测误差曲线。通过比较不难发现,基于 EM 算法的半监督 LWPLS 建模方法所得到的误差曲线在 0 附近波动幅度更小,这表明本文采用的建模方法预测误差较小。同时,仿真过程中还计算了有标签训练样本集占总训练样本集的 10%情况下模型对 1000 组样本进行预测的均方根误差,如表 2 所示。可以看出,采用基于 EM 算法的半监督局部加权偏最小二乘在线建模方法得到的均方根误差低于其他两种建模方法所得。所以,本文提出的建模方法具有更高的精度和泛化能力。

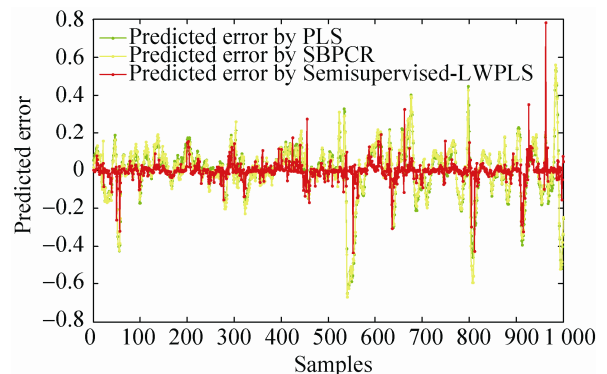


图 8 对丁烷浓度三种建模方法的预测误差

Fig. 8 Prediction errors of three modeling methods for Butane concentration

表2 对丁烷浓度预测均方根误差  
Tab. 2 Prediction root mean square error  
for Butane concentration

预测模型	PLS	SBPCR <sup>[10]</sup>	半监督 LWPLS
RMSE	0.1948	0.1810	0.0635

结合基于半监督学习的建模方法能够利用无标签样本中所携带信息的特性,利用本文方法、PLS建模方法及SBPCR建模方法分别对有标签训练样本集占总训练样本集的10%~50%进行建模结果分析,即有标签样本为100组至500组。对于5种不同比例的有标签训练样本集,3种不同软测量建模方法的预测结果如表3所示,进一步验证了本文所提建模方法的预测精度。

表3 不同有标签训练样本集比例下  
3种建模方法预测结果

Tab. 3 The prediction results of three modeling methods  
with different label ratios

模型\比例/%	10	20	30	40	50
PLS	0.194 8	0.186 2	0.176 5	0.157 9	0.156 8
SBPCR <sup>[10]</sup>	0.181 0	0.178 5	0.168 3	0.153 0	0.150 5
半监督 LWPLS	0.063 5	0.061 8	0.060 6	0.059 7	0.059 3

## 5 结论

对于实际化工过程中通过采样分析获得的有标签的样本数量较少,利用有标签样本建立的软测量模型不能充分反应过程特性的问题,本文提出一种基于半监督学习思想的在线建模方法,在利用有标签样本的同时融合无标签样本中所携带的信息。由于半监督LWPLS建模方法是对在线测得的数据不断的构建局部模型,用局部线性模型来逼近非线性过程,所以能够较好地描述非线性过程的动态特征。通过对脱丁烷塔过程的仿真研究,验证了本文所提建模方法的有效性,对于输出变量难以测得的非线性化工过程的动态建模研究具有一定的参考价值。

## 参考文献:

[1] 王惠文, 孟洁. 多元线性回归的预测建模方法[J]. 北京航空航天大学学报, 2007, 33(4): 500-504.

- Wang Huiwen, Meng Jie. Predictive modeling on multivariate linear regression[J]. Journal of Beijing University of Aeronautics and Astronautics, 2007, 33(4): 500-504.
- [2] 郑小霞, 钱锋. 基于PCA和最小二乘支持向量机的软测量建模[J]. 系统仿真学报, 2006, 18(3): 739-741.  
ZHENG Xiaoxia, QIAN Feng. Soft Sensor Modeling Based on PCA and Support Vector Machines[J]. Journal of System Simulation, 2006, 18(3): 739-741.
- [3] 王巍, 柴天佑, 赵立杰. 带有稳定学习的递归神经网络动态偏最小二乘建模[J]. 控制理论与应用, 2012, 29(3): 337-341.  
WANG Wei, CHAI Tian-you, ZHAO Li-jie. Dynamic partial least squares modeling with recurrent neural networks of stable learning[J]. Control Theory & Applications, 2012, 29(3): 337-341.
- [4] Shao Weiming, Tian Xueming. Adaptive Soft Sensor for Quality Prediction of Chemical Processes Based on Selective Ensemble of Local Partial Least Squares Models[J]. Chemical Engineering Research and Design (S0263-8762), 2015, 95(3): 113-132.
- [5] 王宏伟, 韩云涛, 彭继慎. 基于TSPSO支持向量机红外甲烷传感器动态补偿[J]. 传感技术学报, 2013, 26(9): 1193-1197.  
WANG Hongwei, HAN Yuntao, PENG Jishen. Dynamic Compensation of Infrared Methane Sensor Based on TSPSO- $\epsilon$ -SVM[J]. CHINESE JOURNAL OF SENSORS AND ACTUATORS, 2013, 26(9): 1193-1197.
- [6] 李春富, 郑松, 葛铭. 基于递推非线性部分最小二乘模型的间歇过程批到批优化[J]. 计算机与应用化学, 2011, 28(7): 939-942.  
Li Chunfu, Zheng Song, Ge Ming. Batch-to-batch optimization of batch processes based on recursive nonlinear partial least squares model[J]. Computers and Applied Chemistry, 2011, 28(7): 939-942.
- [7] Cleveland W S. Robust Locally Weighted Regression and Smoothing Scatterplots[J]. Journal of American Statistical Association (S0162-1459), 1979, 74(368): 829-836.
- [8] Kim S, Kano M, Nakagawa H, et al. Estimation of Active Pharmaceutical Ingredients Content Using Locally Weighted Partial Least Squares and Statistical Wavelength Selection[J]. International Journal of Pharmaceutics (S0378-5173), 2011, 421(2): 269-274.
- [9] 章军, 杨慧中. 基于SVM的苯酚浓度半监督软测量方法[J]. 计算机与应用化学, 2013, 30(20): 1453-1456.  
Zhang Jun, Yang Huizhong. Semi-supervised soft sensor

