Contents lists available at ScienceDirect



Chinese Journal of Chemical Engineering

journal homepage: www.elsevier.com/locate/CJCHE



# Process Monitor Adaptive Local Outlier Probability for Dynamic Process Monitoring<sup>☆</sup>



Key Laboratory of Advanced Control and Optimization for Chemical Processes of Ministry of Education, East China University of Science and Technology, Shanghai 200237, China

#### ARTICLE INFO

Article history: Received 24 December 2013 Received in revised form 28 January 2014 Accepted 7 February 2014 Available online 19 June 2014

Yuxin Ma, Hongbo Shi\*, Mengling Wang

Keywords: Time-varying Complex data distribution Local outlier probability Multi-mode Fault detection

## ABSTRACT

Complex industrial processes often have multiple operating modes and present time-varying behavior. The data in one mode may follow specific Gaussian or non-Gaussian distributions. In this paper, a numerically efficient moving window local outlier probability algorithm is proposed. Its key feature is the capability to handle complex data distributions and incursive operating condition changes including slow dynamic variations and instant mode shifts. First, a two-step adaption approach is introduced and some designed updating rules are applied to keep the monitoring model up-to-date. Then, a semi-supervised monitoring strategy is developed with an updating switch rule to deal with mode changes. Based on local probability models, the algorithm has a superior ability in detecting faulty conditions and fast adapting to slow variations and new operating modes. Finally, the utility of the proposed method is demonstrated with a numerical example and a non-isothermal continuous stirred tank reactor.

© 2014 Chemical Industry and Engineering Society of China, and Chemical Industry Press. All rights reserved.

#### 1. Introduction

In industrial processes, operating conditions are usually affected by some slow variations denoted as time-varying characteristics, caused by some dynamic behavior such as seasonal fluctuation, catalyst deactivation, equipment aging, sensor or process drifting, preventive maintenance and cleaning [1]. Generally, effects of the time-varying behavior on the mean and covariance of variables cannot be neglected, so there may be many false alarms if conventional multivariate statistical process monitoring (MSPM) methods are applied directly [2]. In order to maintain process efficiency for a long period of time, numerous adaptive methods have been developed. Recursive MSPM methods and methods based on the moving window strategy are two alternative widely used approaches [3,4].

Multimodality is another important feature of industrial processes due to changes of market demands, alternations of feedstock or variations of manufacturing strategy. The difference between the characteristics of nearby operating conditions is always significant, so intensive studies have been carried out with either multiple local models or a single global model [5,6]. While it is more practical to accommodate the time-varying behavior and multimode features together. The developed methods can be divided into two categories. One is the adaptive clustering methods. Teppola et al. [7] applied adaptive fuzzy C-means algorithms on the score values of principle component analysis (PCA) to monitor a wastewater treatment plant. Liu [8] used an adaptive Takagi-Sugeno fuzzy model on PCA subspace to model a large scale nonlinear system containing many operating regions. Since PCA is used as a preprocessing tool, monitoring results of these two methods more or less depend on and be restricted by the capability of PCA. Petković et al. [9] designed an on-line adaptive clustering method utilizing a generalized information potential. Although previously unseen functioning modes can be included by introducing an adaptive expert system, the method suffers from a non negligible detection delay. The other category is adaptive statistical methods. Improved recursive algorithms based on recursive PCA or the signed digraph were proposed with some if-then rules to distinguish process condition changes from disturbances [10-12]. Ge and Song [13] introduced the just-in-timelearning strategy to the modeling procedure of local least squares support vector regression and the residuals between the real output and the predicted one was analyzed by a two-step information extraction strategy. Xie and Shi [14] and Yu [15] developed two different dynamic fashions of Gaussian mixture model (GMM) separately based on the moving window strategy and a particle filter resampling method.

The problem of complex data distributions in time-varying and multimode processes has scarcely been addressed. Although the moving window strategy has been proven to be effective, it still encounters some limitations when incorporated with statistical methods such as

1004-9541/© 2014 Chemical Industry and Engineering Society of China, and Chemical Industry Press. All rights reserved.

<sup>&</sup>lt;sup>A</sup> Supported by the National Natural Science Foundation of China (61374140), Shanghai Postdoctoral Sustentation Fund (12R21412600), the Fundamental Research Funds for the Central Universities (WH1214039), and Shanghai Pujiang Program (12PJ1402200).

Corresponding author.

E-mail address: hbshi@ecust.edu.cn (H. Shi).

PCA, partial least squares (PLS) or GMM. Since the variables of an industrial process may satisfy specified Gaussian or non-Gaussian distributions, and high order statistics are usually helpful to reveal more information from the data [16–18], adaptive monitoring algorithms should be developed, which can explore both Gaussianity and non-Gaussianity of process data. Local outlier probability (LoOP) is an unsupervised data mining technique proposed for outlier detection [19]. It combines the idea of local, density-based outlier scoring with a probabilistic, statistically-oriented approach, and assigns the probability of being an outlier to all data records. Since a normalization procedure is included, LoOP is independent of any specific data distribution. Therefore, a combination of LoOP and moving window strategy should be potential to tackle these problems.

The main contribution of this paper is to propose a numerically efficient moving window LoOP algorithm for monitoring industrial processes with complex data distributions, time-varying property and multiple operating modes. Some designed rules are introduced and incorporated with a two-step adaption approach to ensure that the monitoring model can be updated at a high speed. To cope with the multimode features, a semi-supervised monitoring strategy is employed, and an update termination rule is developed to prevent the monitoring model contaminated by faults or disturbances. Since the method is based on local probabilistic models, the accuracy of model is higher and it will be much easier to detect faulty conditions.

### 2. Adaptive Process Monitoring Based on Moving Window Loop

For low computation burden and practical applications, it is fast and reasonable to only update the information of those samples whose neighbors have changed due to the insertion and discard of samples. The key problems addressed in this section are how to find the affected samples and how to update their information.

#### 2.1. Offline initialization

To make an initialization and calculate the LoOP value for each sample  $\mathbf{x}_j$  (j = 1, 2, ..., L) with dimension D in the initial window  $\mathbf{W}_1$ , its k nearest neighbors are found as follows, with its neighborhood set in  $\mathbf{W}_1$  can be recognized as  $\mathbf{knn}_1(\mathbf{x}_i)$ .

$$d\left(\boldsymbol{x}_{j}, \boldsymbol{x}_{p}\right) = \sqrt{\sum_{n=1}^{D} \left| \boldsymbol{x}_{jn} - \boldsymbol{x}_{pn} \right|^{2}} \qquad \left( p \neq j \text{ and } \boldsymbol{x}_{j}, \boldsymbol{x}_{p} \in \boldsymbol{W}_{1} \right)$$
(1)

Assuming that samples in **knn**<sub>1</sub>( $x_j$ ) are centered around  $x_j$ , then we can define probabilistic set distance as:

$$pdist_1(\boldsymbol{x}_j) = \lambda \cdot \sqrt{\sum_{\boldsymbol{x}_p \in \mathbf{knn}_1(\boldsymbol{x}_j)} d(\boldsymbol{x}_j, \boldsymbol{x}_p)^2 / k}$$
(2)

where  $\lambda$  is a weighted factor usually taken as 2. For estimating the density around  $x_j$ , the probabilistic local outlier factor (PLOF) is defined as follows with function E(.) used to compute the expectation of PLOF in the current window.

$$PLOF_1(\boldsymbol{x}_j) = pdist_1(\boldsymbol{x}_j) / \left( E_{\boldsymbol{x}_p \in \boldsymbol{knn}_1(\boldsymbol{x}_j)} \left[ pdist_1(\boldsymbol{x}_p) \right] \right) - 1$$
(3)

To achieve normalization, the aggregate value  $nPLOF_1$  which can be considered as a standard deviation of PLOF values is obtained:

$$nPLOF_1 = \lambda \cdot \sqrt{E[(PLOF_1)^2]}$$
(4)

Finally, by applying the Gaussian error function, the local outlier probability indicating the probability that a sample is an outlier can be calculated as:

$$LoOP_1(\boldsymbol{x}_j) = \max\left\{0, \operatorname{erf}\left(PLOF_1(\boldsymbol{x}_j) / \left(\sqrt{2} \cdot nPLOF_1\right)\right)\right\}$$
(5)

where erf(.) is the Gaussian error function applied to obtain a probabilistic value.

#### 2.2. Online updating and process monitoring

By applying the moving window strategy, a two-step adaption procedure is introduced to update the monitoring model. Some more details of the adaption procedure for a window size L are as follows.

#### Step 1: discard

The effect of eliminating the oldest sample  $x_i$  from the previous window  $W_i$  on the mean and variance can be evaluated as follows.

$$\widetilde{\boldsymbol{\mu}} = (L\boldsymbol{\mu}_i - \boldsymbol{x}_i)/(L - 1) \tag{6}$$

$$\Delta \widetilde{\boldsymbol{\mu}} = \boldsymbol{\mu}_i - \widetilde{\boldsymbol{\mu}} \tag{7}$$

$$\widetilde{\boldsymbol{\sigma}}(m)^{2} = \frac{1}{L-2} \left( (L-1) \cdot \left( \left[ \boldsymbol{\sigma}_{i}(m) \right]^{2} - \left[ \Delta \widetilde{\boldsymbol{\mu}}(m) \right]^{2} \right) - \left[ \boldsymbol{x}_{i}(m) - \boldsymbol{\mu}_{i}(m) \right]^{2} \right) \quad (8)$$
$$(m = 1, 2..., D)$$

$$\widetilde{\boldsymbol{\Sigma}} = \operatorname{diag}[\widetilde{\boldsymbol{\sigma}}(1), \widetilde{\boldsymbol{\sigma}}(2), \cdots, \widetilde{\boldsymbol{\sigma}}(D)]$$
(9)

where diag(.) is the function used to calculate the diagonal matrix. Eq. (6) describes the updating of the variable mean while Eqs. (7)–(9) describe the updating of the variable variance.

After moving all the information about  $\mathbf{x}_i$  from the current monitoring model, a set  $\mathbf{S}_{i-1}^0$  (i > 1) is constructed to store the samples, in which  $\mathbf{x}_i$  is one of their k nearest neighbors.

$$\boldsymbol{S}_{i-1}^{0} = \boldsymbol{S}_{i-1}^{0} \cup \left\{ \boldsymbol{x}_{j} \right\}, \text{ if } i \neq j, \ \boldsymbol{x}_{j} \in \boldsymbol{W}_{i} \text{ and } \boldsymbol{x}_{i} \in \mathbf{knn}_{i} \left( \boldsymbol{x}_{j} \right)$$
(10)

where **knn**<sub>*i*</sub>(**x**<sub>*j*</sub>) represents the neighborhood set of sample **x**<sub>*j*</sub> in window **W**<sub>*i*</sub>. Obviously, if **x**<sub>*j*</sub>  $\in$  **S**<sup>0</sup><sub>*i*-1</sub>, due to the deletion of **x**<sub>*i*</sub>, the neighborhood set **knn**<sub>*i*</sub>(**x**<sub>*j*</sub>) will change.

Step 2: insertion

When a new sample  $\mathbf{x}_{i+L}$  is judged normal and added into the data matrix, the updated mean vector and variance in  $\mathbf{W}_{i+1}$  are computed as follows.

$$\boldsymbol{\mu}_{i+1} = \left[ (L-1)\widetilde{\boldsymbol{\mu}} + \boldsymbol{x}_{i+L} \right] / L \tag{11}$$

$$\Delta \boldsymbol{\mu}_{i+1} = \boldsymbol{\mu}_{i+1} - \widetilde{\boldsymbol{\mu}} \tag{12}$$

$$\boldsymbol{\sigma}_{i+1}(m)^{2} = \frac{1}{L-1} \left( (L-2) \cdot \left[ \tilde{\boldsymbol{\sigma}}(m) \right]^{2} + (L-1) \cdot \left[ \Delta \boldsymbol{\mu}_{i+1}(m) \right]^{2} + \left[ \boldsymbol{x}_{i+L}(m) - \boldsymbol{\mu}_{i+1}(m) \right]^{2} \right) \quad (m = 1, 2..., D)$$
(13)

$$\boldsymbol{\Sigma}_{i+1} = \operatorname{diag}[\boldsymbol{\sigma}_{i+1}(1), \boldsymbol{\sigma}_{i+1}(2), \cdots, \boldsymbol{\sigma}_{i+1}(D)]$$
(14)

Eq. (11) describes the updating of the mean vector while Eqs. (12)–(14) describe the updating of the variance. However, only for those with new sample  $\mathbf{x}_{i+L}$  among their *k* nearest

Download English Version:

# https://daneshyari.com/en/article/166012

Download Persian Version:

https://daneshyari.com/article/166012

Daneshyari.com