

Long-term dependency slow feature analysis for dynamic process monitoring

Xinrui Gao*, Yuri A.W. Shardt*

* *Technical University of Ilmenau, Ilmenau, Thuringia, Germany, 98684*
(e-mail: {xinrui.gao,yuri.shardt}@tu-ilmenau.de).

Abstract: Industrial processes are large scale, highly complex systems. The complex flow of mass and energy, as well as the compensation effects of closed-loop control systems, cause significance cross-correlation and autocorrelation between process variables. To operate the process systems stably and efficiently, it is crucial to uncover the inherent characteristics of both variance structure and dynamic relationship. Long-term dependency slow feature analysis (LTSFA) is proposed in this paper to overcome the Markov assumption of the original slow feature analysis to understand the long-term dynamics of processes, based on which a monitoring procedure is designed. A simulation example and the Tennessee Eastman process benchmark are studied to show the performance of LTSFA. The proposed method can better extract the system dynamics and monitor the process variations using fewer slow features.

Keywords: Process monitoring, slow feature analysis, fault diagnosis, latent variable model

1. INTRODUCTION

Modern industrial processes are large scale, highly complex systems with many units and equipment. A large number of measurements are set to monitor process operation, and much process data is automatically collected. To operate the process systems stably and efficiently, developing process models based on process data to describe the process characteristics is crucial. The process data has two distinctive characteristics. The first is the serious collinearity caused by mass and energy balances and the compensation effects of the closed-loop control systems. Secondly, the process system always shows significant time-dependent characteristics. For instance, a transition process would occur while state switching due to the long settling time. Similarly, oscillations would occur after a disturbance introduced because of recycle streams, heat integrations, and other connections of materials, energy, and information (Shardt et al., 2012).

Multivariate statistical process monitoring (MSPM) approaches have been extensively studied to deal with the above two problems (Qin, 2012). The most widely used method is principal component analysis (PCA) (Gajjar, Kulahci, & Palazoglu, 2018; Wise, Ricker, Veltkamp, & Kowalski, 1990), which extracts uncorrelated lower dimensional latent variables to concisely describe the main variance structure of the process observation space. Independent component analysis (ICA) recovers the independent latent variables from complex process data by leveraging high order moment information (Kano, Tanaka, Hasebe, Hashimoto, & Ohno, 2003). For the description of time-dependent characteristics, dynamic PCA (DPCA) performs the standard PCA on the augmented data with time lags to extract process dynamics (Ku, Storer, & Georgakis, 1995). However, it fails to decouple the dynamics and static variance information. Slow feature analysis (SFA) is used to

concurrently monitor the process by describing the stationary distribution and dynamic behaviours separately (Shang et al., 2015). However, assuming the standard Markov property, namely that the one-step time dependency is sufficient, limits the performance of the standard SFA and has to be improved in the same way as DPCA in practice.

Thus, this paper proposes a new long-term dependency slow feature analysis (LTSFA) method that can extract longer time dependencies for the design of a process monitoring method. This approach is validated using a simulated example and the Tennessee Eastman process.

2. LINEAR SLOW FEATURE ANALYSIS

Given an m -dimensional ergodic observation signal $\mathbf{x}(t) = [x_1(t) \ \cdots \ x_m(t)]$, SFA seeks to find a latent signal $\mathbf{s}(t) = [s_1(t) \ \cdots \ s_k(t)]^T$ with the slowest variation, also named slow features (SFs), to describe its time varying characteristics (Wiskott & Sejnowski, 2002), that is,

$$\begin{aligned} & \min_{g_j(\bullet)} \langle \dot{s}_j^2(t) \rangle \\ \text{s.t.} \quad & 1) \langle s_j(t) \rangle = 0, \\ & 2) \langle s_j^2(t) \rangle = 1, \\ & 3) \forall i \neq j, \langle s_i(t)s_j(t) \rangle = 0 \end{aligned} \tag{1}$$

where $g_j(\bullet)$ is the input-output function that needs to be found, $\langle \bullet \rangle$ is the temporal averaging operator, and $\dot{s}(t)$ is the first-order derivative or difference with respect to time and can be approximated for discrete time series as

$$\dot{s}(t) \approx s(t) - s(t-1) \tag{2}$$

Objective (1) uses the average squared temporal difference to define the signal slowness. Constraints 1) and 2) simplify the problem without loss of generality, and constraint

2) avoids the trivial solution $s_j = \text{constant}$. Constraint 3) guarantees that the SFs are statistically uncorrelated so that they carry different information. It also implicitly implies that the extracted SFs are sorted by their slowness. In the linear case, the input-output function set $\mathbf{g} = [g_j(\cdot)]_{j=1}^k$ is a linear transformation

$$\mathbf{s}(t) = \mathbf{g}(\mathbf{x}(t)) = \mathbf{P}\mathbf{x}(t) \quad (3)$$

where $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k]^\top$ is the parameter matrix.

For a mean centred signal $\mathbf{x}(t)$, due to

$$\begin{aligned} \langle \hat{s}_j^2(t) \rangle &= \langle (s_j(t))^2 + (s_j(t-1))^2 - 2(s_j(t)s_j(t-1)) \rangle \\ &= 2 - 2\langle s_j(t)s_j(t-1) \rangle \end{aligned} \quad (4)$$

the objective function of SFA is equivalent to the following formulation with the same constraints, that is,

$$\max_{\mathbf{p}} \langle s_j(t)s_j(t-1) \rangle \quad (5)$$

Using Lagrange multipliers, this optimization problem can be translated into the generalized eigenvalue decomposition (GED) problem, that is

$$\begin{aligned} C^x(1)\mathbf{P} &= C^x(0)\mathbf{P}\mathbf{\Lambda} \\ \mathbf{\Lambda} &= \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_k\}, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \end{aligned} \quad (6)$$

where λ_j is the generalized eigenvalue of matrix pair $(C^x(1), C^x(0))$, and $C^x(\tau)$ is the symmetric version of the τ -step autocovariance matrix of $\mathbf{x}(t)$ defined as

$$C^x(\tau) \equiv \frac{1}{2} (\langle \mathbf{x}(t)\mathbf{x}^\top(t-\tau) \rangle + \langle \mathbf{x}(t-\tau)\mathbf{x}^\top(t) \rangle) \quad (7)$$

3. LONG-TERM DEPENDENCY SFA

3.1. Optimization problem for long-term dependency SFA

From Equation (4), it can be concluded that minimizing the average squared temporal difference of SF is equivalent to maximizing its one-step lagged autocorrelation. However, in practical systems, one-step lagged autocorrelation is far from sufficient to describe the dynamics, especially manufacturing processes with typical long-term dependency. Hence, it is necessary to improve SFA to possess long-term dependency modelling ability.

From a predictability perspective, the long-term dependency can be represented as

$$\hat{s}(t) = \beta_1 s(t-1) + \beta_2 s(t-2) + \dots + \beta_l s(t-l) \quad (8)$$

where l is the maximum time delay. Thus, the objective of the long-term dependency SFA can be formulated as

$$\begin{aligned} \langle s(t)\hat{s}(t) \rangle &= \frac{1}{2} \langle s(t)\hat{s}(t) + \hat{s}(t)s(t) \rangle \\ &= \sum_{\tau=1}^l \beta_\tau C^s(\tau) = \mathbf{p}^\top \left(\sum_{\tau=1}^l \beta_\tau C^x(\tau) \right) \mathbf{p} \end{aligned} \quad (9)$$

where $C^s(\tau)$ and $C^x(\tau)$ are defined as in Equation (7). Let $\mathbf{X} = [\mathbf{x}(1) \ \dots \ \mathbf{x}(N+l)]$ be the total sample collection of $\mathbf{x}(t)$. Construct the following matrices

$$\mathbf{X}_i = [\mathbf{x}(i) \ \dots \ \mathbf{x}(i+N-1)] \text{ for } i=1, \dots, l+1 \quad (10)$$

$$\mathbf{A}_\tau = [\mathbf{X}_{i+1} \ \mathbf{X}_{i+1-\tau}], \ \mathbf{B}_\tau = [\mathbf{X}_{i+1-\tau} \ \mathbf{X}_{i+1}]$$

where $\tau \in \{1, 2, \dots, l\}$ is the time delay. Then, the τ -step covariance matrix can be written as

$$\begin{aligned} C^x(\tau) &= \frac{1}{2} (\langle \mathbf{x}(t)\mathbf{x}^\top(t-\tau) \rangle + \langle \mathbf{x}(t-\tau)\mathbf{x}^\top(t) \rangle) \\ &= \frac{1}{2N} \sum_{t=l+1}^{l+N} (\mathbf{x}(t)\mathbf{x}^\top(t-\tau) + \mathbf{x}(t-\tau)\mathbf{x}^\top(t)) = \frac{\mathbf{A}_\tau \mathbf{B}_\tau^\top}{2N} \end{aligned} \quad (11)$$

Substituting Equation (11) into Equation (9), gives

$$\langle s(t)\hat{s}(t) \rangle = \frac{1}{2N} \mathbf{p}^\top \left(\sum_{\tau=1}^l \beta_\tau \mathbf{A}_\tau \mathbf{B}_\tau^\top \right) \mathbf{p} = \frac{1}{2N} \mathbf{p}^\top \mathbf{A} \mathbf{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} \quad (12)$$

where $\mathbf{A} = [\mathbf{A}_1 \ \mathbf{A}_2 \ \dots \ \mathbf{A}_l]$, $\mathbf{B} = [\mathbf{B}_1 \ \mathbf{B}_2 \ \dots \ \mathbf{B}_l]$,

$$\mathbf{\Lambda}_\beta = \text{diag}\{\beta_1 \mathbf{I}, \beta_2 \mathbf{I}, \dots, \beta_l \mathbf{I}\}.$$

Then, the optimization objective of LTSFA can be rewritten as

$$\begin{aligned} \max_{\mathbf{p}, \beta} \quad & \mathbf{p}^\top \mathbf{A} \mathbf{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} \\ \text{s.t.} \quad & 1) \ \mathbf{p}^\top \mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top \mathbf{p} = 1 \\ & 2) \ \|\beta\| = 1 \end{aligned} \quad (13)$$

where constraint 1) guarantees that the extracted SFs have unit variances to avoid the trivial solution. It is clear that LTSFA will reduce to SFA if the time delay is set to one.

3.2. Algorithm for long-term dependency SFA

3.2.1. Analysis of the optimization formulation

Introducing Lagrange multipliers λ_p and λ_β , Equation (13) can be converted into

$$J = \mathbf{p}^\top \mathbf{A} \mathbf{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} + \lambda_p (1 - \mathbf{p}^\top \mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top \mathbf{p}) + \lambda_\beta (1 - \beta^\top \beta) \quad (14)$$

Taking derivatives with respect to \mathbf{p} and β , and setting them to zero, gives

$$\frac{\delta J}{\delta \mathbf{p}} = 2\mathbf{A} \mathbf{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} - 2\lambda_p \mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top \mathbf{p} = 0 \quad (15)$$

$$\frac{\delta J}{\delta \beta} = (\mathbf{I} \otimes \mathbf{p})^\top ([\mathbf{\Gamma}_1 \ \dots \ \mathbf{\Gamma}_l]^\top \otimes \mathbf{B}) \mathbf{A}^\top \mathbf{p} - 2\lambda_\beta \beta = 0 \quad (16)$$

where \otimes is the Kronecker product, and $\mathbf{\Gamma}_i$ is a selection matrix:

$$\mathbf{\Gamma}_i = \text{diag}\{\delta_1 \mathbf{I}, \dots, \delta_j \mathbf{I}, \dots, \delta_l \mathbf{I}\}, \ \delta_j = \begin{cases} 1, & j=i \\ 0, & j \neq i \end{cases} \quad (17)$$

Premultiplying Equations (15) and (16) by \mathbf{p}^\top and β^\top respectively, the relationship $J^* = \lambda_p = 2\lambda_\beta$ can be obtained from

$$2\lambda_p = 2\lambda_p \mathbf{p}^\top \mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top \mathbf{p} = 2\mathbf{p}^\top \mathbf{A} \mathbf{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} = 2J^* \quad (18)$$

$$2\lambda_\beta = 2\lambda_\beta \boldsymbol{\beta}^\top \boldsymbol{\beta} = \boldsymbol{\beta}^\top (\mathbf{I} \otimes \mathbf{p})^\top \left([\boldsymbol{\Gamma}_1 \ \cdots \ \boldsymbol{\Gamma}_l]^\top \otimes \mathbf{B} \right) \mathbf{A}^\top \mathbf{p} \quad (19)$$

$$= \mathbf{p}^\top (\mathbf{A} \boldsymbol{\Lambda}_\beta \mathbf{B}^\top)^\top \mathbf{p} = \mathbf{p}^\top \mathbf{A} \boldsymbol{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} = J^*$$

From Equations (15) and (18), the optimal value of Equation (13) J^* is the largest generalized eigenvalue of the matrix pair

$$\left(\mathbf{A} \boldsymbol{\Lambda}_\beta \mathbf{B}^\top, \mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top \right) = \left(\sum_{\tau=1}^l \beta_\tau C^x(\tau), C^x(0) \right)$$

The optimal solution \mathbf{p}^* is the corresponding generalized eigenvector. However, \mathbf{p} and $\boldsymbol{\beta}$ are coupled together and thus cannot be solved using analytical methods.

3.2.2. Iterative solution for one slow feature

Let $s(t) = \mathbf{x}^\top(t) \mathbf{p}$ be the SFs at time instant t , then we can form the following data matrices,

$$\mathbf{s}_i = [s(i) \ \cdots \ s(i+N-1)]^\top = \mathbf{X}_i^\top \mathbf{p} \text{ for } i=1, \dots, l+1 \quad (20)$$

$$\widehat{\mathbf{S}} = [\mathbf{s}_1 \ \mathbf{s}_{l-1} \ \cdots \ \mathbf{s}_l], \mathbf{Z} = [\mathbf{X}_1 \ \mathbf{X}_{l-1} \ \cdots \ \mathbf{X}_l]$$

where \mathbf{X}_i is defined as in Equation (10). Then, Equation (15) can be rewritten as

$$\begin{aligned} \lambda_p \mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top \mathbf{p} &= \mathbf{A} \boldsymbol{\Lambda}_\beta \mathbf{B}^\top \mathbf{p} = \sum_{\tau=1}^l (\beta_\tau \mathbf{A}_\tau \mathbf{B}_\tau^\top \mathbf{p}) \\ &= \mathbf{X}_{l+1} \sum_{\tau=1}^l (\beta_\tau \mathbf{X}_{l+1-\tau}^\top \mathbf{p}) + \sum_{\tau=1}^l (\beta_\tau \mathbf{X}_{l+1-\tau}) \mathbf{X}_{l+1}^\top \mathbf{p} \\ &= \mathbf{X}_{l+1} [\mathbf{s}_1 \ \cdots \ \mathbf{s}_l] \boldsymbol{\beta} + [\mathbf{X}_1 \ \cdots \ \mathbf{X}_l] (\boldsymbol{\beta} \otimes \mathbf{I}) \mathbf{s}_{l+1} \\ &= \mathbf{X}_{l+1} \widehat{\mathbf{S}} \boldsymbol{\beta} + \mathbf{Z} (\boldsymbol{\beta} \otimes \mathbf{I}) \mathbf{s}_{l+1} \end{aligned} \quad (21)$$

Thus, we have

$$\lambda_p \mathbf{p} = (\mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top)^\dagger (\mathbf{X}_{l+1} \widehat{\mathbf{S}} \boldsymbol{\beta} + \mathbf{Z} (\boldsymbol{\beta} \otimes \mathbf{I}) \mathbf{s}_{l+1}) \quad (22)$$

where $(\cdot)^\dagger$ is the Moore-Penrose inverse. Equation (16) can be rewritten as

$$\begin{aligned} 2\lambda_\beta \boldsymbol{\beta} &= (\mathbf{I} \otimes \mathbf{p})^\top \left([\boldsymbol{\Gamma}_1 \ \cdots \ \boldsymbol{\Gamma}_l]^\top \otimes \mathbf{B} \right) \mathbf{A}^\top \mathbf{p} \\ &= (\mathbf{I} \otimes \mathbf{p})^\top \text{diag}\{\mathbf{B}_1, \dots, \mathbf{B}_l\} \mathbf{A}^\top \mathbf{p} = 2\widehat{\mathbf{S}}^\top \mathbf{s}_{l+1} \end{aligned} \quad (23)$$

Then, the iterative solution is:

- (1) Scale \mathbf{x} to zero mean and unit variance.
- (2) Initialize \mathbf{p} with a random unit vector.
- (3) Iteratively solve \mathbf{p} and $\boldsymbol{\beta}$ until the objective function J converges to the optimal value, that is,

$$\mathbf{s} = \mathbf{X}^\top \mathbf{p}, \quad \mathbf{s} := \mathbf{s} / \|\mathbf{s}\|$$

$$\boldsymbol{\beta} = \widehat{\mathbf{S}}^\top \mathbf{s}_{l+1}, \quad \boldsymbol{\beta} := \boldsymbol{\beta} / \|\boldsymbol{\beta}\|$$

$$\mathbf{p} = (\mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top)^\dagger (\mathbf{X}_{l+1} \widehat{\mathbf{S}} \boldsymbol{\beta} + \mathbf{Z} (\boldsymbol{\beta} \otimes \mathbf{I}) \mathbf{s}_{l+1})$$

$$J = \mathbf{s}^\top \sum_{\tau=1}^l \beta_\tau \mathbf{s}_{l+1-\tau} = \mathbf{s}_{l+1}^\top \widehat{\mathbf{S}} \boldsymbol{\beta}$$

If the predefined stopping condition is fulfilled, the obtained \mathbf{p} and $\boldsymbol{\beta}$ are the optimal solution corresponding to the SF.

3.2.3. Multiple slow feature extraction through matrix deflation

After the first SF is obtained, the next one can be extracted by applying Step (3) in Section 3.2.2 to the residual data excluding the previous SF information. The residual data can be obtained by matrix deflation

$$\mathbf{X} \equiv \mathbf{X} - \mathbf{q} \mathbf{s}^\top \quad (24)$$

where the loading vector $\mathbf{q} = \mathbf{X} \mathbf{s} / \mathbf{s}^\top \mathbf{s}$ is the solution of

$$\min \|\mathbf{X} - \mathbf{q} \mathbf{s}^\top\|_F \quad (25)$$

Iteratively performing the above procedure can obtain all the k SFs in the descending order of slowness or predictability. The different SFs are orthogonal to each other, which is the same as in the standard SFA.

3.2.4. Model developed for historical slow feature extraction and prediction

After all the SFs have been extracted, the optimal projection vectors, loading vectors, and regression weights are collected as

$$\mathbf{P} = [\mathbf{p}_1 \ \cdots \ \mathbf{p}_k]^\top, \mathbf{Q} = [\mathbf{q}_1 \ \cdots \ \mathbf{q}_k]^\top, \boldsymbol{\Phi} = [\boldsymbol{\beta}_1 \ \cdots \ \boldsymbol{\beta}_k]^\top.$$

Let $\mathbf{X}^{(i)}$ be the residual data matrix after the i -th deflation, $\mathbf{X}^{(i)} = \mathbf{X}$ and $\mathbf{s}^{(i)} = (\mathbf{X}^{(i)})^\top \mathbf{p}_i$ be the i -th SF, the observation space can be divided using Equation (24) into the SF subspace and residual subspace

$$\mathbf{X} = \mathbf{q}_1 (\mathbf{s}^{(1)})^\top + \cdots + \mathbf{q}_k (\mathbf{s}^{(k)})^\top + \mathbf{X}^{(k+1)} = \mathbf{Q}^\top \mathbf{S}^\top + \mathbf{X}^{(k+1)} \quad (26)$$

where $\mathbf{S} = [\mathbf{s}^{(1)} \ \cdots \ \mathbf{s}^{(k)}]$ is the k SFs of the historical data \mathbf{X} , and is also a set of the orthogonal basis of the SF subspace. Since the residuals $\mathbf{X}^{(k+1)}$ belong to the nullspace of the projection matrix \mathbf{P} , we have

$$\mathbf{X}^\top \mathbf{P}^\top - \mathbf{S} \mathbf{Q} \mathbf{P}^\top = (\mathbf{X}^{(k+1)})^\top \mathbf{P}^\top = \mathbf{0} \quad (27)$$

Hence, the historical SF model is

$$\mathbf{S} = \mathbf{X}^\top \mathbf{P}^\top (\mathbf{Q} \mathbf{P}^\top)^{-1} \quad (28)$$

For an observation $\mathbf{x}(t)$, its projection onto SF subspace is

$$\mathbf{s}(t) = (\mathbf{P} \mathbf{Q}^\top)^{-1} \mathbf{P} \mathbf{x}(t) \quad (29)$$

Given the historical data $\bar{\mathbf{X}} = [\mathbf{x}(t) \ \cdots \ \mathbf{x}(t-l+1)]$, the one-step ahead prediction model can be obtained as

$$\begin{aligned} \hat{\mathbf{s}}(t+1) &= [\hat{\mathbf{s}}^{(1)}(t+1) \ \cdots \ \hat{\mathbf{s}}^{(k)}(t+1)]^\top \\ &= [\boldsymbol{\beta}_1^\top \bar{\mathbf{s}}^{(1)} \ \cdots \ \boldsymbol{\beta}_k^\top \bar{\mathbf{s}}^{(k)}]^\top = \boldsymbol{\Lambda}_{\bar{\mathbf{S}}} \text{vec}(\boldsymbol{\Phi}^\top) \end{aligned} \quad (30)$$

where $\bar{\mathbf{s}}^{(i)} = [s^{(i)}(t) \ \cdots \ s^{(i)}(t-l+1)]^\top$, $\text{vec}(\cdot)$ is the column vectorization operator,

$\boldsymbol{\Lambda}_{\bar{\mathbf{S}}} = \text{diag}\left\{(\bar{\mathbf{s}}^{(1)})^\top \ \cdots \ (\bar{\mathbf{s}}^{(k)})^\top\right\}$ is a blocked diagonal

matrix deduced from $\bar{\mathbf{S}}$, and $\bar{\mathbf{S}}$ is the slow feature matrix of $\bar{\mathbf{X}}$, that is,

$$\bar{\mathbf{S}} = \begin{bmatrix} \bar{s}^{(1)} & \dots & \bar{s}^{(k)} \end{bmatrix}^\top = \left(\bar{\mathbf{X}}^\top \mathbf{P}^\top (\mathbf{Q}\mathbf{P}^\top)^{-1} \right)^\top \quad (31)$$

Furthermore, one-step ahead observation can be predicted by

$$\hat{\mathbf{x}}(t+1) = \mathbf{Q}^\top \hat{\mathbf{s}}(t+1) \quad (32)$$

In summary, the detailed procedure of long-term dependency SFA algorithm is given in Table 1.

Table 1: The algorithm for LTSFA

Input:	$\mathbf{X} = [\mathbf{x}(1) \ \dots \ \mathbf{x}(N+l)]$,
	$\mathbf{x}(t) = [x_1(t) \ \dots \ x_m(t)]^\top$	
1.	Scale \mathbf{X} to zero mean and unit variance	
2.	Initialize \mathbf{p} with a random unit vector	
3.	Iterate the following steps until convergence to obtain the optimal parameter pair $(\mathbf{p}^*, \boldsymbol{\beta}^*)$:	
1)	$\mathbf{s} = \mathbf{X}^\top \mathbf{p}$, $\mathbf{s} := \mathbf{s} / \ \mathbf{s}\ $	
2)	Form the data matrices: $\mathbf{X}_i = [\mathbf{x}(i) \ \dots \ \mathbf{x}(i+N-1)]$ for $i=1, \dots, l+1$, $\mathbf{s}_i = [s(i) \ \dots \ s(i+N-1)]^\top$ for $i=1, \dots, l+1$, $\hat{\mathbf{S}} = [\mathbf{s}_1 \ \dots \ \mathbf{s}_1]$, $\mathbf{Z} = [\mathbf{X}_1 \ \dots \ \mathbf{X}_1]$	
3)	$\boldsymbol{\beta} = \hat{\mathbf{S}}^\top \mathbf{s}_{l+1}$, $\boldsymbol{\beta} := \boldsymbol{\beta} / \ \boldsymbol{\beta}\ $	
4)	$\mathbf{p} = (\mathbf{X}_{l+1} \mathbf{X}_{l+1}^\top)^\dagger (\mathbf{X}_{l+1} \hat{\mathbf{S}} \boldsymbol{\beta} + \mathbf{Z}(\boldsymbol{\beta} \otimes \mathbf{I}) \mathbf{s}_{l+1})$	
5)	$J = \mathbf{s}_{l+1}^\top \sum_{\tau=1}^l \beta_\tau \mathbf{s}_{l+1-\tau} = \mathbf{s}_{l+1}^\top \hat{\mathbf{S}} \boldsymbol{\beta}$	
4.	Collect \mathbf{p}^* , $\boldsymbol{\beta}^*$ into $\mathbf{P} = [\mathbf{p}_1 \ \dots \ \mathbf{p}_k]^\top$, $\boldsymbol{\Phi} = [\boldsymbol{\beta}_1 \ \dots \ \boldsymbol{\beta}_k]^\top$	
5.	Data matrix deflation. 1) $\mathbf{X} := \mathbf{X} - \mathbf{q}\mathbf{s}^\top$, where $\mathbf{q} = \mathbf{X}\mathbf{s} / \mathbf{s}^\top \mathbf{s}$ 2) Collect \mathbf{q} into $\mathbf{Q} = [\mathbf{q}_1 \ \dots \ \mathbf{q}_k]^\top$ Return to Step 3 to extract the next SF until all k SFs have been obtained.	
6.	SF extraction of historical observation \mathbf{X} $\mathbf{S} = \mathbf{X}^\top \mathbf{P}^\top (\mathbf{Q}\mathbf{P}^\top)^{-1}$, $\mathbf{s}(t) = (\mathbf{P}\mathbf{Q}^\top)^{-1} \mathbf{P}\mathbf{x}(t)$	
7.	One-step prediction given historical data $\bar{\mathbf{X}} = [\mathbf{x}(t) \ \dots \ \mathbf{x}(t-l+1)]$ 1) $\hat{\mathbf{s}}(t+1) = \boldsymbol{\Lambda}_{\bar{\mathbf{S}}} \text{vec}(\boldsymbol{\Phi}^\top)$, where $\boldsymbol{\Lambda}_{\bar{\mathbf{S}}} = \text{diag} \left\{ (\bar{s}^{(1)})^\top \ \dots \ (\bar{s}^{(k)})^\top \right\}$, $\bar{\mathbf{S}} = [\bar{s}^{(1)} \ \dots \ \bar{s}^{(k)}]^\top = \left(\bar{\mathbf{X}}^\top \mathbf{P}^\top (\mathbf{Q}\mathbf{P}^\top)^{-1} \right)^\top$ 2) $\hat{\mathbf{x}}(t+1) = \mathbf{Q}^\top \hat{\mathbf{s}}(t+1)$	

3.2.5. Determining the hyperparameters

The model construction of LTSFA needs to determine two hyperparameters: the number of SF k and dynamic order l . In machine learning, the data set is always split into three subsets: training set, validation set, and test set. The model performance of each pair of hyperparameters is validated by the validation set after model training using the training set to select the optimal hyperparameters with the best performance.

For LTSFA, after all k SFs have been extracted, the residual $\mathbf{X}^{(k+1)}$ is expected to be white, and thus the autocorrelation and cross-correlation for nonzero lags will be approximately zero. The 95% confidence interval (CI) of the correlations (Shardt, 2015) will be used to verify model performance. The optimal hyperparameters correspond to the least violation for the 95% CI.

4. LTSFA-BASED PROCESS MONITORING

4.1. Observation-space partition

The residual after performing LTSFA on the observation is essentially white. However, it contains a static structure, which can be further explored using standard PCA. Hence, the observation space can be decoupled into a slow feature subspace, a static principal component subspace, and a residual subspace, that is,

$$\mathbf{X} = \mathbf{Q}^\top \mathbf{S}^\top + \mathbf{X}^{(k+1)} = \mathbf{Q}^\top \mathbf{S}^\top + \mathbf{W}\mathbf{T}^\top + \mathbf{E} \quad (33)$$

where $\mathbf{W} = [\mathbf{w}_1 \ \dots \ \mathbf{w}_r]$ is the loading matrix and $\mathbf{T} = (\mathbf{X}^{(k+1)})^\top \mathbf{W}$ is the score matrix of the principal components. The number of static principal components is determined by the method of the cumulative percentage of variance (CPV) (Valle, Li, & Qin, 1999).

4.2. Process monitoring indices design

After partitioning the space, we can monitor the variations in the subspaces to detect abnormal situations using Hotelling's T^2 and the squared prediction error (SPE), that is,

$$T_d^2 = \mathbf{s}^\top \mathbf{s} = (\mathbf{P}\mathbf{Q}^\top)^{-1} \mathbf{P}\mathbf{x}^\top \mathbf{x}\mathbf{P}^\top (\mathbf{Q}\mathbf{P}^\top)^{-1} \quad (34)$$

$$T_s^2 = \mathbf{t}^\top \boldsymbol{\Omega}^{-1} \mathbf{t} = \left(\mathbf{x}^{(k+1)} \right)^\top \mathbf{W}\boldsymbol{\Omega}^{-1} \mathbf{W}^\top \mathbf{x}^{(k+1)} \quad (35)$$

$$SPE = \mathbf{e}^\top \mathbf{e} = \|(\mathbf{I} - \mathbf{W}\mathbf{W}^\top) \mathbf{x}^{(k+1)}\|^2 \quad (36)$$

where $\boldsymbol{\Omega} = \frac{1}{N} \mathbf{T}^\top \mathbf{T} = \text{diag} \{ \lambda_1, \dots, \lambda_r \}$ and λ_i is the

eigenvalue of $\mathbf{X}^{(k+1)} (\mathbf{X}^{(k+1)})^\top$. The above three indices actually monitor the stationary variations in the subspaces. The dynamic behaviour anomaly can be detected by using the innovation as

$$T_v^2 = \mathbf{v}^\top \mathbf{P}_v \boldsymbol{\Lambda}_v^{-1} \mathbf{P}_v^\top \mathbf{v} \quad (37)$$

where $\mathbf{v}(t) = \mathbf{s}(t) - \hat{\mathbf{s}}(t)$ is the innovation and $\langle \mathbf{v}(t) \mathbf{v}^\top(t) \rangle = \mathbf{P}_v^\top \boldsymbol{\Lambda}_v \mathbf{P}_v$ is the singular value decomposition of the innovation covariance.

Under an assumption of a Gaussian distribution, the control limits for the indices can be determined as follow given a significance level α (Chiang, Russell, & Braatz, 2000; Qin, 2003)

$$T_d^2 \leq T_d^2(\alpha) = \chi_k^2(\alpha) \quad (38)$$

$$T_s^2 \leq T_s^2(\alpha) = \frac{r(N-1)(N+1)}{N(N-r)} F_{r,N-r}(\alpha) \quad (39)$$

$$SPE \leq SPE(\alpha) = g\chi_h^2(\alpha) \quad (40)$$

$$T_v^2 \leq \frac{k(N-1)(N+1)}{N(N-k)} F_{k,N-k}(\alpha) \quad (41)$$

where χ_k^2 and χ_h^2 are χ^2 -distribution with k and h degrees of freedom, $F_{r,N-r}$ and $F_{k,N-k}$ are F -distributions with r and $N-r$, k and $N-k$ degrees of freedom, $g = \theta_2/\theta_1$, $h = \theta_1^2/\theta_2$, and

$\theta_j = \sum_{i=r+1}^m \lambda_i^j$ for $j=1,2$. If the indices of an observation exceed the control limits, then it is considered that a fault has occurred.

5. CASE STUDIES

5.1. Simulation example

A second-order autoregressive model is designed

$$\begin{cases} \mathbf{s}_t = \mathbf{A}_1 \mathbf{s}_{t-1} + \mathbf{A}_2 \mathbf{s}_{t-2} + \mathbf{v}_t \\ \mathbf{x}_t = \mathbf{B} \mathbf{s}_t + \mathbf{e}_t \end{cases} \quad (42)$$

$$\mathbf{A}_1 = \begin{bmatrix} 0.7904 & 0 & 0 \\ 0 & 0.5026 & 0 \\ 0 & 0 & 0.7904 \end{bmatrix}$$

$$\mathbf{A}_2 = \begin{bmatrix} 0.1004 & 0 & 0 \\ 0 & 0.3026 & 0 \\ 0 & 0 & -0.2104 \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} -0.8207 & 0.5970 & 0.0253 \\ -0.8164 & -0.4752 & -0.6769 \\ -0.9511 & -0.5484 & 0.2549 \\ -0.3652 & -0.5081 & 0.4817 \\ 0.6245 & 0.7796 & 0.7808 \end{bmatrix}$$

where $\mathbf{e}_t \in \mathbb{R}^5 \sim N(0, 0.5^2 \mathbf{I})$ and $\mathbf{v}_t \in \mathbb{R}^3 \sim N(0, 0.5^2 \mathbf{I})$ are independent and identically distributed random processes.

Three thousand data samples are generated and split into training, validation, and test set with 2000, 500, and 500 data samples respectively. Using the method described in Section 3.2.5, the hyperparameters are determined as $l=2$ and $k=3$. The autocorrelation and cross-correlations of the original data $\mathbf{x}(t)$ and residual $\tilde{\mathbf{x}}(t)$ are shown in Figure 1 and Figure 2. It can be seen that the dynamics of the data, which is indicated by both autocorrelation and cross-correlation with nonzero delays, is filtered by LTSFA. Figure 3 is the residual correlations of the standard SFA. Assuming a maximum time delay of 20 samples, the percentage of times a point lies

outside the 95% CI was 3.2% for LTSFA SFs, but 40.4% for the standard SFA. The correlations of the latent variables $\mathbf{s}(t)$ and the innovation $\mathbf{v}(t)$ in Figure 4 show that LTSFA can also uncover the dynamics of the latent states, while the standard SFA does not give an explicit expression of latent state dependency.

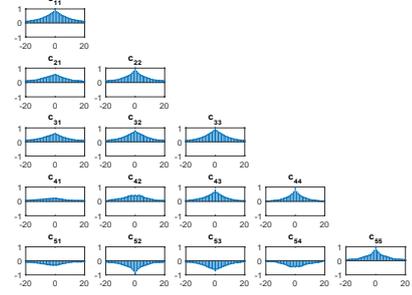


Figure 1: Autocorrelation and cross-correlation of $\mathbf{x}(t)$

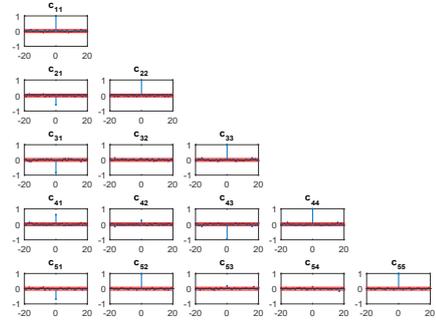


Figure 2: Autocorrelation and cross-correlation for the LTSFA residual

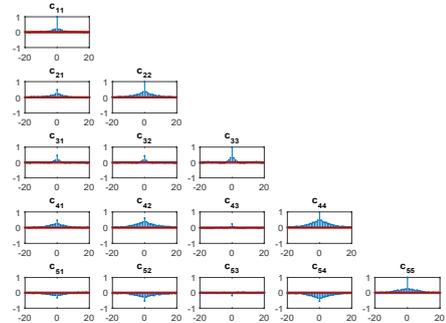


Figure 3: Autocorrelation and cross-correlation for the standard SFA residual

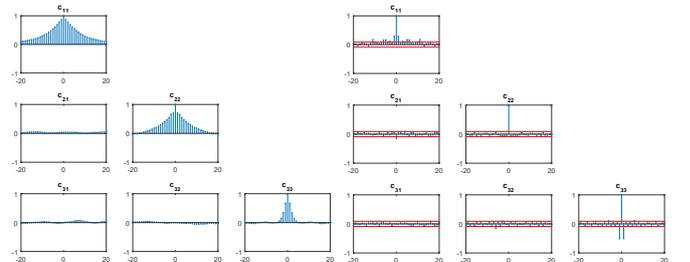


Figure 4: Correlations for left $\mathbf{s}(t)$ and right innovation $\mathbf{v}(t)$

5.2. Tennessee Eastman process benchmark

The Tennessee Eastman process (TEP) is a widely used benchmark in process control (Downs & Vogel, 1993). The revised version is adopted in this study (Lyman, Georgakis, & engineering, 1995). It contains 12 manipulated variables (XMV (1-12)) and 41 measurement variables which consist of 22 process variables (XMEAS (1-22)) and 19 quality variables (XMEAS (23-41)). In this paper, we use 11 manipulated variables (XMV (1-11)) and all 22 process variables. Five hundred normal samples sampled at 3-min intervals are used to train the models. The hyperparameters of LTSFA are determined as $l = 3$ and $k = 19$. The dynamic SFA provided in reference (Shang et al., 2015) is built to compare the monitoring results. The optimal hyperparameters of dynamic SFA are selected as $d = 2$ and $q = 0.55$. The monitoring results of LTSFA can give more detailed information of the process variations. For instance, Figure 5 shows that the variation of IDV(4) mainly occurs in the SF subspace and the residual subspace, while the static principal component subspace is less affected. Figure 6 shows that the anomaly caused by the high-frequency oscillation of IDV(14) can also be monitored in the SF subspace immediately. Furthermore, LTSFA only uses 19 SFs, while the dynamic SFA 46. Thus, the dynamic SFA may not be able to perform dimension reduction, especially in the case of large delays.

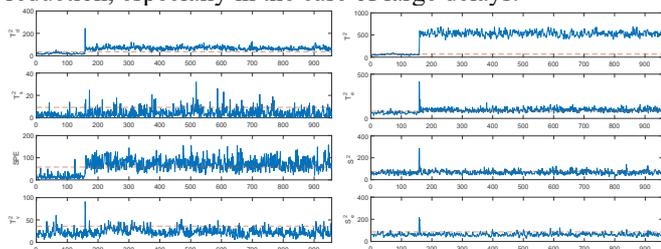


Figure 5: Monitoring results for IDV(4) using left LTSFA and right the dynamic SFA

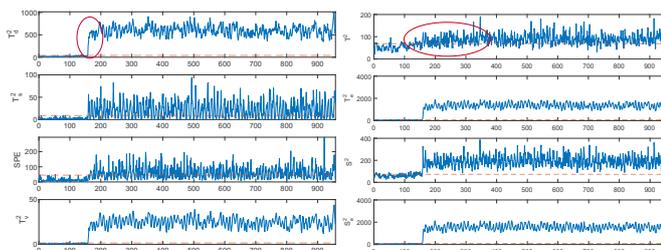


Figure 6: Monitoring results for IDV(14) using left LTSFA and right the dynamic SFA

6. CONCLUSIONS

This paper presented a new model called LTSFA to overcome the Markov assumption of the original SFA and improve the long-term dynamics modelling ability. The objective of LTSFA is established based on a high-order regressive model and an iterative algorithm is developed to solve the objective. In addition to the better dynamic modelling capability, LTSFA also gives an explicit expression of the long-term temporal dependency of latent states. A process monitoring strategy is developed using this approach. It is tested on a simulated example and the Tennessee Eastman

process. It is shown that the proposed method can better extract the system dynamics and give more detailed process variations using fewer slow features than the standard SFA and the dynamic SFA methods. Further work will consider extending the results to nonlinear systems.

REFERENCES

- Chiang, L. H., Russell, E. L., & Braatz, R. D. (2000). *Fault detection and diagnosis in industrial systems*: Springer Science & Business Media.
- Downs, J. J., & Vogel, E. F. (1993). A plant-wide industrial process control problem. *Computers chemical engineering*, 17(3), 245-255.
- Gajjar, S., Kulahci, M., & Palazoglu, A. (2018). Real-time fault detection and diagnosis using sparse principal component analysis. *Journal of Process Control*, 67, 112-128.
- Kano, M., Tanaka, S., Hasebe, S., Hashimoto, I., & Ohno, H. (2003). Monitoring independent components for fault detection. *AIChE Journal*, 49(4), 969-976.
- Ku, W., Storer, R. H., & Georgakis, C. (1995). Disturbance detection and isolation by dynamic principal component analysis. *Chemometrics intelligent laboratory systems*, 30(1), 179-196.
- Lyman, P. R., Georgakis, C. J. C., & engineering, c. (1995). Plant-wide control of the Tennessee Eastman problem. 19(3), 321-331.
- Qin, S. J. (2003). Statistical process monitoring: basics and beyond. *Journal of Chemometrics: A Journal of the Chemometrics Society*, 17(8-9), 480-502.
- Qin, S. J. (2012). Survey on data-driven industrial process monitoring and diagnosis. *Annual reviews in control*, 36(2), 220-234.
- Shang, C., Yang, F., Gao, X., Huang, X., Suykens, J. A., & Huang, D. (2015). Concurrent monitoring of operating condition deviations and process dynamics anomalies with slow feature analysis. *AIChE Journal*, 61(11), 3666-3682.
- Shardt, Y. A.W. (2015). *Statistics for Chemical and Process Engineers*: Springer.
- Shardt, Y. A.W., Zhao, Y., Qi, F., Lee, K., Yu, X., Huang, B., & Shah, S. (2012). Determining the state of a process control system: Current trends and future challenges. *The Canadian Journal of Chemical Engineering*, 90(2), 217-245.
- Valle, S., Li, W., & Qin, S. J. (1999). Selection of the number of principal components: the variance of the reconstruction error criterion with a comparison to other methods. *Industrial Engineering Chemistry Research*, 38(11), 4389-4401.
- Wise, B. M., Ricker, N., Veltkamp, D., & Kowalski, B. R. (1990). A theoretical basis for the use of principal component models for monitoring multivariate processes. *Process control quality*, 1(1), 41-51.
- Wiskott, L., & Sejnowski, T. J. (2002). Slow feature analysis: Unsupervised learning of invariances. *Neural computation*, 14(4), 715-770.