# Calibration, Prediction and Process Monitoring Model Based on Factor Analysis for Incomplete Process Data

Dong Soon KIM[1], Chang Kyoo YOO[1], Young Il KIM[2], Jae Hak JUNG[3] and In-Beum LEE[1]

[1]*School of Environmental Engineering/Department of Chemical Engineering, Pohang University of Science and Technology, San 31 Hyoja Dong, Pohang 790-784, Korea*

[2]*Energy Research Department, Research Institute of Industrial Science & Technology (RIST), San 32 Hyoja Dong, Pohang 790-330, Korea*

[3]*School of Chemical Engineering Technology, Yeungnam University, Daedong, Gyeongsan, 712-749, Korea*

Unspecific missing of values in real chemical and biological industries have been found. Regardless of the incompleteness of the measured sample, a monitoring system should be designed to tackle the missing data problem and be applied to on-line systems immediately. A calibration method of a factor analysis (FA) model for incomplete data sets is proposed. And a prediction method based on the calibrated model is suggested in order to estimate missing values in incomplete calibration sets and incomplete test sets. An expectation and maximization (EM) algorithm is used to calibrate the model and expectation of conditional density is used to predict the model result. The proposed method is compared with the well-known iterative singular values decomposition (iSVD) method, i.e. a principal component analysis (PCA) based method; and a simple data set is tested as an illustrative example. The proposed method gives better estimation results for the missing values than the well-known PCA based method. There are several advantages of the proposed method over the PCA based projection methods: (1) data pretreatment is not an essential step since the FA model is scale invariant whereas the PCA model is not, (2) since the proposed method utilizes probability information of all variables directly, to apply it as a statistical process monitoring technique is preferable to others, and (3) the single model can be extended to a mixture of such models by the virtue of the EM algorithm.

## Introduction

The real-world multivariate data set can have missing values at random by some reasons, e.g. troubles of sensors for industrial data, incomplete answers of respondents for statistical survey data, etc. The on-line measured sample can also have missing values in general since different sensing periods among sensors may exist; some measurements not yet received in the sample are regarded as missing values. Some techniques can be used to handle such types of data set or data, like a calibration technique of a multivariate model from an incomplete calibration set, and a prediction method for the missing values in every sample. Suppose that there is an incomplete calibration set. Calibrating a multivariate model from the set, should we discard a whole sample just for small portion of missing in it? Again suppose a partial sample from

missing data would be received successively. When diagnosing the sample, should we wait all measurements' scores to be gathered? This paper primarily aims at answering these two questions.

A principal component analysis (PCA) based regression methods to estimate the missing values are well known (Grung and Mamme, 1998; Walczak and Massart, 2001). However, in spite of their algorithmic simplicity, they have a critical defect that the PCA model is scale variant inherently, that is, differently scaled data sets result in totally different models. It is the natural consequence of PCA which is designed to focus on variances of variables. When working with PCA, a pretreatment step on variables to have equal importance is essential.

The factor analysis (FA) model by an EM algorithm is scale invariant, where its parameters are the factor loading matrix, mean vector and noise variance matrix. The FA pursues learning the most probable model parameters that can best explain correlations between the measurements and the factors under Gaussian probability density assumption on the

measurement space. When making the FA model, we may not worry about a scale problem of the measurements. Mathematically, PCA is a special case of FA in which set noise variances infinitesimal (Tipping and Bishop, 1997a, 1997b; Kim and Lee, 2003).

After setting several notations in Section 1, we will summarize the PCA based calibration and prediction method from incomplete data in Section 2. The FA model based probabilistic approach will be proposed in Section 3. Comparisons of the two methods will be presented in Section 4 through the case study of a simple example. Finally, the conclusions will be addressed.

## 1. Incomplete Data

Suppose that there is a calibration data set $\mathbf{X} = \{\boldsymbol{x}_n\}_{n \in N} \in \mathbb{R}^{P \times N}$ with $\boldsymbol{x}_n = [\boldsymbol{x}^o_n; \boldsymbol{x}^m_n] \in \mathbb{R}^P$. Here, subscript '$n$' represents a calibration sample index, and superscript 'o' and 'm' symbolize the observed and the missing data, respectively; the dimension of the observed vector of the $n$-th sample, $\dim(\boldsymbol{x}^o_n)$, can have one of $\{1, 2, \cdots, \dim(\boldsymbol{x})\}$ according to the sample index $n$ but $\dim(\boldsymbol{x}^o_n) + \dim(\boldsymbol{x}^m_n) = \dim(\boldsymbol{x})$ $\forall$ $n$. Notation $n \in N$ to express $n = 1, 2, \cdots, N$ just for its conciseness is used in this paper. This unusual notation should be distinguished from an 'element' symbol, e.g. $\boldsymbol{x} \in \{\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{z}\}$ or $\boldsymbol{x} \in \mathbb{R}^P$. In parenthesis '[·]', a semicolon is the symbol to denote the next row while comma represents the next column, e.g. $[1; 2] \in \mathbb{R}^{2 \times 1}$ and $[1, 2] \in \mathbb{R}^{1 \times 2}$. Let us denote the observed part in $\mathbf{X}$ as $\mathbf{X}^o = \{\boldsymbol{x}^o_n\}_{n \in N}$, and the missing part in $\mathbf{X}$ as $\mathbf{X}^m = \{\boldsymbol{x}^m_n\}_{n \in N}$. Hence $\mathbf{X}^o \cup \mathbf{X}^m = \mathbf{X}$ and $\mathbf{X}^o \cap \mathbf{X}^m = \{\phi\}$.

For a newly measured sample, $\boldsymbol{x}_n$, it can also have a missing value. Let us denote $\boldsymbol{x}_n = [\boldsymbol{x}^i_n; \boldsymbol{x}^t_n]$, where subscript '$n$' represents a newly measured sample index, and superscript 'i' and 't' symbolize the input and the target, respectively. Here, to distinguish $\boldsymbol{x}_n$ from $\boldsymbol{x}_n$ clearly, we gave the names 'target' and 'input' for the new sample; the target and the input in the test sample are equivalent to the missing and the observed data in the calibration sample, respectively. **Figure 1** shows an example of missing values in the multivariate data, where the hollowed squares represent the observed element and the grey squares indicate the missing element in multivariate data. For instance, in this figure, a calibration data set $\mathbf{X} = \{\boldsymbol{x}_n\}_{n \in 10} \in \mathbb{R}^{7 \times 10}$ in which $\boldsymbol{x}_1 = [\boldsymbol{x}^o_1 \in \mathbb{R}^6; \boldsymbol{x}^m_1 \in \mathbb{R}^1]$, and a test data vector $\boldsymbol{x}_n = [\boldsymbol{x}^i_n \in \mathbb{R}^5; \boldsymbol{x}^t_n \in \mathbb{R}^2]$. Then our tasks are clear: (1) to calibrate the FA model only from $\mathbf{X}^o$, and (2) to predict $\boldsymbol{x}^t_n$ from $\boldsymbol{x}^i_n$ using the calibrated FA model.

## 2. PCA from Incomplete Data

First of all, the well-known iterative PCA based regression method is addressed.
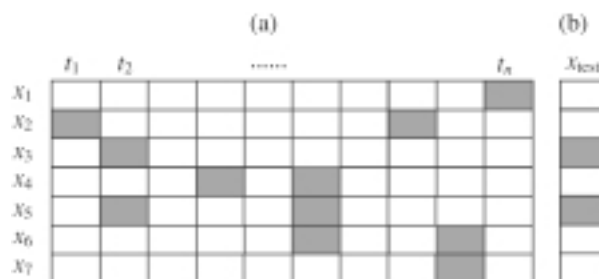


**Fig. 1** Examples of missing values at random: (a) a calibration data set, and (b) a test data vector

### 2.1 Calibration

The PCA model of $\mathbf{X}$ reduced to $L$ (=rank($\mathbf{Z}$) $\leq$ rank($\mathbf{X}$) = $P$) dimensional PC space is given by

$$\mathbf{X} = \mathbf{A} \cdot \mathbf{Z} + \mathbf{E} \tag{1}$$

where $\mathbf{A} = \{\boldsymbol{a}_l\}_{l \in L}$ with $\boldsymbol{a}_i^T \cdot \boldsymbol{a}_j \in \{0, 1\}$ for $\{i \neq j, i = j\}$, $\mathbf{Z} = \{\underline{\boldsymbol{z}}_l\}_{l \in L}$ and $\mathbf{E} = \mathbf{X} - \sum_{l \in L} \boldsymbol{a}_l \cdot \underline{\boldsymbol{z}}_l$ have been called the loading matrix, the score matrix, and the residual matrix, respectively. All vectors used in this article are column vectors unless specified; we will use a double under bar to denote a row vector, e.g. $\underline{\boldsymbol{z}}_l$ is the $l$-th row of an $L \times N$ matrix $\mathbf{Z}$ while $\boldsymbol{z}_n$ is the $n$-th column of $\mathbf{Z}$. If $\mathbf{X}^m = \{\phi\}$, i.e. there is no missing in $\mathbf{X}$, a least-square sense optimal pair of $\boldsymbol{a}_l$ and $\underline{\boldsymbol{z}}_l$ is easily extracted by a nonlinear iterative partial least square (NIPALS) algorithm successively or a singular value decomposition (SVD) technique at once. In fact NIPALS is a method to solve the SVD problem of $\mathbf{X}$, e.g. $[\mathbf{A}, \mathbf{S}, \mathbf{V}] = \text{SVDs}(\mathbf{X}, L)$ then $\mathbf{Z} = \mathbf{S} \cdot \mathbf{V}^T$ and $\mathbf{E} = \mathbf{X} - \mathbf{A} \cdot \mathbf{S} \cdot \mathbf{V}^T$. However, in the case of $\mathbf{X}^m \neq \{\phi\}$, the iterative SVD method (Walczak and Massart, 2001) can be applicable. The iSVD method enables us to estimate $\mathbf{X}^m$ as well as both $\mathbf{A}$ and $\mathbf{Z}$, e.g.
(1) initial guess of $\underline{\mathbf{X}} = \mathbf{X}^o \cup \underline{\mathbf{X}}^m_{\text{initial}}$
(2) until $\underline{\mathbf{X}}^m$ converges, iterate the followings:

$$\begin{aligned} &[\mathbf{A}, \mathbf{S}, \mathbf{V}] = \text{SVDs}(\underline{\mathbf{X}}, L), \\ &\mathbf{X}_{\text{PCA}} = \mathbf{A} \cdot \mathbf{S} \cdot \mathbf{V}^T, \\ &\underline{\mathbf{X}}^m = \mathbf{X}^m_{\text{PCA}}, \underline{\mathbf{X}} = \mathbf{X}^o \cup \underline{\mathbf{X}}^m \end{aligned} \tag{2}$$

where an 'under bar' denotes the estimate. We may set $L = P$. Then an orthogonal square matrix $\mathbf{A}$, which is to be a new set of basis vectors, transforms the original correlated variable $\boldsymbol{x}$ to a new uncorrelated variable $\boldsymbol{z}$ without any information loss, i.e. $\mathbf{E} = \{\phi\}$. This is the multiple linear regression (MLR) approach. But in general, we set $L < P$ to have robust estimate results. It is the principal component regression (PCR) approach. In essence MLR is a special case of PCR:

$$\mathbf{X}^m_{\text{PCR}} \leftarrow \text{SVDs}(\underline{\mathbf{X}}, L) \text{ and } \mathbf{X}^m_{\text{MLR}} \leftarrow \text{SVDs}(\underline{\mathbf{X}}, P) \tag{3}$$

where $L \in \{1, 2, \ldots, \leq P\}$. Notice that slight modification of NIPALS algorithm can generate almost all kinds of multivariate projection models, e.g. PLS, continuum regression (CR), cyclic subspace regression (CSR), and so on (Kalivas, 1999; Geladi, 2002). The iterative algorithm can be applied to the intermediate models without serious structural change.

## 2.2 Prediction

To analyze a new incomplete sample, $x_n = [x^i_n; x^t_n]$, the PCA model can be used, $x = A \cdot z + e$, as a combined two sub-models: an input model, $x^i = A^i \cdot z + e^i$, and a target model, $x^t = A^t \cdot z + e^t$. The combined model form can be expresses as follows:

$$[x^i; x^t] = [A^i; A^t] \cdot z + [e^i; e^t] \tag{4}$$

where $A^i$ and $A^t$ are submatrices of $A$ corresponding to $\dim(x^i)$ and $\dim(x^t)$, respectively. Suppose $x = [x_1; x_2; x_3]$, $z = [z_1; z_2]$, $e = [e_1; e_2; e_3]$ and $A = [a_{11}, a_{12}; a_{21}, a_{22}; a_{31}, a_{32}]$. If $x^i = [x_1; x_3]$, $x^t = [x_2]$, then $A^i = [a_{11}, a_{12}; a_{31}, a_{32}]$ and $A^t = [a_{21}, a_{22}]$; $e^i = [e_1; e_3]$, $e^t = [e_2]$. If $\dim(x^i_n) \neq 0$, the corresponding target vector, $x^t_n$, is predicted via two steps: first least-square sense optimal PC scores using the input model, then to predict the target scores using the target model:

$$\underline{z}_n = A^{i+}_n \cdot x^i_n \tag{5}$$

$$\underline{x}^t_n = A^t_n \cdot \underline{z}_n \tag{6}$$

where left pseudo-inverse of $A^i_n$ is defined by $A^{i+}_n \equiv (A^{iT}_n \cdot A^i_n)^{-1} \cdot A^{iT}_n \neq A^{iT}_n$. Notice that since $\dim(x^i_j) \neq \dim(x^i_k)$ for $j \neq k$ was permitted, the number of rows of $A^i_n$ and $A^t_n$ are varied according to the sample index $n$; however, $\text{row}(A^i_n) + \text{row}(A^t_n) = \dim(x)$, and $\text{column}(A^i_n) = \text{column}(A^t_n) = \dim(z)$ are preserved $\forall n$.

## 2.3 PCA model of scale variant

The PCA model is scale variant. To ease elucidation, let us define a scaling matrix $B \equiv \text{diag}(\{b_p\}_{p \in P})$ that changes the scale of $X$ to $B \cdot X$ while it does not affect the correlations among variables. However, the PCA model of $X$ is totally different from that of $B \cdot X$. It means that $X \to B \cdot X$ does not imply neither $A \cdot Z \to B \cdot A \cdot Z$ nor $E \to B \cdot E$. The PCA model is also shift variant. Let us introduce a mean-shift vector $m \in \mathbb{R}^P$ that moves $X$ to $X' = \{x_n + m\}_{n \in N}$. But the movement does not imply neither $A \cdot Z \to A \cdot Z'$ nor $E \to E'$, where $Z' = \{z_n + \alpha \cdot m\}_{n \in N}$ and $E' = \{e_n + (1 - \alpha) \cdot m\}_{n \in N}$, for any $0 \leq \alpha \leq 1$.

Therefore some pretreatments on a calibration set to have equal importance among variables are a prerequisite for the PCA model through mean-centering ($x \to x^{mc}$), auto-scaling ($x \to x^{auto}$), or whitening ($x \to x^{white}$), where $x^{mc} = x - \text{mean}(x)$, $x_p^{auto} = \text{var}(x_p)^{-0.5} \cdot (x_p - \text{mean}(x_p)) \; \forall \; p$, $x^{white} = \text{cov}(x)^{-0.5} \cdot (x - \text{mean}(x))$. The pretreatments are also needed to analyze the test sample since the sample should be transformed to have identical statistics of the calibration set. Without any missing in both $x_n$ and $X$, the pretreatments are simple. But if there are any missing values in $X$, they are cumbersome because we should estimate the means and the (co)variances of the variables only from $X^o$.

## 3. Factor Analysis (FA) from Incomplete Data

The factor analysis (FA) model (Anderson, 1984) has the form of

$$x = A \cdot z + \mu + e \tag{7}$$

where $A \in \mathbb{R}^{P \times L}$ denotes the loading matrix, $\mu \in \mathbb{R}^P$ symbolizes the mean vector of the measurement variable $x \in \mathbb{R}^P$; $z \in \mathbb{R}^L$ represents the factor variable, and $e \in \mathbb{R}^P$ signifies the noise variable. Learning task in the model is obtaining the likelihood-sense optimal parameters set $\Theta \equiv \{A, \mu, \Lambda\}$ from a (in)complete calibration data set $X = X^o \cup X^m$. To do this, we assume that probability density function (PDF) of the factor is $\mathcal{N}(z: 0, I)$, where a random variable $z$ can be completely characterized by a Gaussian parametric function with its mean $\mu$ and covariance $\Sigma$. And the noise is $\mathcal{N}(e: 0, \Lambda = \text{diag}(\{\lambda_p\}_{p \in P}))$, where each noise element, $e_p$, can have a different variance, $\lambda_p$. Furthermore, factors and noises are supposed to be statistically independent of each other, i.e. $p(e \mid z) = p(e)$. It means that there is no information in $z$ about $e$.

## 3.1 Calibration

We consider Eq. (7) as a combined form of two sub-models: an observed model and a missing model,

$$[x^o; x^m] = [A^o; A^m] \cdot z + [\mu^o; \mu^m] + [e^o; e^m] \tag{8}$$

It is obvious that Eq. (7) is a special case of Eq. (8) when $x^m = \{\phi\}$. Since the Gaussian PDF is closed to linear operations, and both $z$ and $e$ were assumed to Gaussians; the followings are clear: $\mathcal{N}(x^m: \mu^m, \Sigma^{mm})$, $\mathcal{N}(x^o: \mu^o, \Sigma^{oo})$, $\mathcal{N}(x^m \mid z: \mu^{m|z}, \Sigma^{m|z})$, and $\mathcal{N}(z|x^o: \mu^{z|o}, \Sigma^{z|o})$. Here, the mean and the covariance in each PDF are easily derived as follows:

$$\Sigma^{mm} = A^m \cdot A^{mT} + \Lambda^m \tag{9}$$

$$\Sigma^{oo} = A^o \cdot A^{oT} + \Lambda^o \tag{10}$$

$$\mu^{m|z} = A^m \cdot z + \mu^m \tag{11}$$

$$\Sigma^{m|z} = \Lambda^m \tag{12}$$

$$\mu^{z|o} = A^{oT} \cdot (\Sigma^{oo})^{-1} \cdot (x^o - \mu^o) \tag{13}$$

$$\Sigma^{z|o} = I - A^{oT} \cdot (\Sigma^{oo})^{-1} \cdot A^o \tag{14}$$

Notice that if $\dim(z) < \dim(x^o)$, say, $A^o$ is a thin matrix, inversion of $\Sigma^{oo}$ is often intractable. For Eqs.

(13) and (14), let us denote $\mathcal{N}([x^o; z]: [\mu^o; 0], [\Sigma^{oo}, \Sigma^{oz}; \Sigma^{zo}, \Sigma^{zz}])$ then $\mathcal{N}(z|x^o: \Sigma^{zo}\cdot\Sigma^{oo-1}\cdot(x^o - \mu^o), \Sigma^{zz} - \Sigma^{zo}\cdot\Sigma^{oo-1}\cdot\Sigma^{oz})$; where $\Sigma^{oo} = A^o\cdot A^{oT} + \Lambda^o$, $\Sigma^{oz} = A^o$, $\Sigma^{zo} = A^{oT}$, and $\Sigma^{zz} = I$. However, there is an efficient way to the type of inversion using the matrix inversion lemma,

$$(A^o\cdot A^{oT} + \Lambda^o)^{-1}$$
$$= \Lambda^{o-1} - \Lambda^{o-1}\cdot A^o\cdot(I + A^{oT}\cdot\Lambda^{o-1}\cdot A^o)^{-1}\cdot A^{oT}\cdot\Lambda^{o-1} \quad (15)$$

The lemma is essential to calculate both Eqs. (13) and (14), where a more detailed proof is shown in Appendix B. Without the lemma, numerical problems may be inevitable when the FA model is calibrated.

3.1.1 MLE for a complete data set    The maximum likelihood estimator (MLE) is a well-known technique to calibrate the FA model. Let's define an augmented data set $D \equiv \{d_n\}_{n\in N} \in \mathbb{R}^{(P+L)\times N}$ with $d_n = [x^o_n; x^m_n; z_n] \in \mathbb{R}^{(P+L)}$, where each data vector in $D$ consists of an observed vector, a missing vector, and a factor vectors. Notice that the factor vector is regarded as a permanent missing vector with the fixed dimension $L$ while the missing vector can have arbitrary dimension from zero to $(\dim(x) - 1)$ according to its sample index. Given both $x^m$ and $z$ in $d$, that is, $d$ is a complete vector, probability density of $d$, $p(d)$, and log likelihood of $d$, $l(d)$, are given by

$$p(d: \Theta) = p(z)\cdot p(x \mid z: \Theta)$$
$$= \mathcal{N}(z: 0, I)\cdot\mathcal{N}(x \mid z: A\cdot z + \mu, \Lambda) \quad (16)$$

$$l(\Theta: d) \equiv \log p(d: \Theta)$$
$$= 0.5\cdot\log \det (\Lambda^{-1})$$
$$- 0.5\cdot(x - A'\cdot z')^T\cdot\Lambda^{-1}\cdot(x - A'\cdot z') + \delta \quad (17)$$

where $A' \equiv [A, \mu] \in \mathbb{R}^{P\times(L+1)}$, $z' \equiv [z; 1] \in \mathbb{R}^{L+1}$, and $\delta$ represents $\Theta$ independent terms in the log likelihood function. The Gaussian PDF of $x$ is $\mathcal{N}(x: \mu, \Sigma) = (2\pi)^{-0.5\cdot\dim(x)}\cdot\det(\Sigma^{-1})^{0.5}\cdot\exp(-0.5\cdot(x-\mu)^T\cdot\Sigma^{-1}\cdot(x-\mu))$. Let us suppose an independent and identically distributed (*iid*) condition on $D$. Mathematically, *iid* on $D$ implies that $p(d_i \mid D \setminus d_i) = p(d_i) = p(d_j)$, where the first equality indicates an independent condition and the second identity represents an identically distributed condition. Given both $x^m_n$ and $z_n$ $\forall$ $n$ under *iid* on $D$, the log likelihood of $D$, $\mathcal{L}(D)$, is simply the sum of $l(\Theta: d_n)$ $\forall$ $n$.

$$\mathcal{L}(\Theta: D)$$
$$\equiv \sum_{n\in N} l(\Theta: x_n)$$
$$= \sum_{n\in N}[0.5\cdot\log \det(\Lambda^{-1})$$
$$- 0.5\cdot(x_n - A'\cdot z'_n)^T\cdot\Lambda^{-1}\cdot(x_n - A'\cdot z'_n)] + \delta \quad (18)$$

Therefore the likelihood-sense optimal $\Theta$ that maximizes $\mathcal{L}(\Theta: D)$ is found by

$$(\partial/\partial A')[\mathcal{L}(\Theta: D)] = 0$$
$$\Rightarrow A' = (\sum_{n\in N} x_n\cdot z_n'^T)\cdot(\sum_{n\in N} z_n'\cdot z_n'^T)^{-1} \quad (19)$$

$$(\partial/\partial \Lambda^{-1})[\mathcal{L}(\Theta: D)] = 0$$
$$\Rightarrow \Lambda = N^{-1}\cdot\text{diag} \sum_{n\in N}(x_n\cdot x_n^T - A'\cdot z_n'\cdot x_n^T) \quad (20)$$

where $(\partial/\partial A)\cdot(x - A\cdot z)^T\cdot\Lambda^{-1}\cdot(x - A\cdot z) = -2\cdot\Lambda^{-1}\cdot(x - A\cdot z)\cdot z^T$ for symmetric $\Lambda$, $(\partial/\partial\Lambda^{-1})\cdot[\log \det(\Lambda^{-1})] = \Lambda^T$, and $(\partial/\partial\Lambda^{-1})\cdot[(x - A\cdot z)^T\cdot\Lambda^{-1}\cdot(x - A\cdot z)] = (x - A\cdot z)\cdot(x - A\cdot z)^T$ and $\text{diag}(\cdot)$ denotes set off-diagonal elements in the parenthesis to zeros. Note that $A'$ in Eq. (20) should be the resultant of Eq. (19). If a mean-centered $X$ is used to these equations, Eq. (19) indicates the least-square solution, $A = X\cdot Z^+$, and Eq. (20) implies slightly biased, $N\cdot(N - 1)^{-1}$, the sample error variance matrix whose diagonals are the variances of the estimation errors, i.e. $\Lambda = N^{-1}\cdot\text{diag}(X\cdot(I - Z^+\cdot Z)\cdot X^T)$.

3.1.2 EM for incomplete data    Only given $X^o$ in $D = X^o \cup X^m \cup Z$, direct maximization of $\mathcal{L}(\Theta: D)$ is impossible through analytic ways. An expected complete data log likelihood function can be assumed, as it has been called an energy function, $q(\cdot)$ for $l(\cdot)$ first.

$$q(\Theta|\Theta_t)$$
$$\equiv \mathcal{E}_{z,xm|xo:\Theta t}[l(d: \Theta)]$$
$$= \mathcal{E}_{z,xm|xo:\Theta t}[0.5\cdot\log \det(\Lambda^{-1})$$
$$- 0.5\cdot(x - A'\cdot z')^T\cdot\Lambda^{-1}\cdot(x - A'\cdot z')] + \delta \quad (21)$$

where $\mathcal{E}_{z,xm|xo:\Theta t}[l(d: \Theta)] \equiv \iint l(x^o, x^m, z: \Theta)\cdot p(z, x^m|x^o: \Theta_t)dz dx^m$, and $\Theta_t$ denotes an estimate of $\Theta$ at iteration time $t$. If $x^m$ were given, then $q(\Theta|\Theta_t) = \mathcal{E}_{z|x:\Theta t}[l(d: \Theta)]$ which results in the standard FA. Since *iid* on $D$ was assumed, an energy function $\Theta(\cdot)$ for $\mathcal{L}(\cdot)$ is just the summation of $q_n(\Theta|\Theta_t)$ $\forall$ $n$.

$$Q(\Theta|\Theta_t)$$
$$\equiv \sum_{n\in N} q_n(\Theta|\Theta_t)$$
$$= \sum_{n\in N}\mathcal{E}_{zn,xmn|xon:\Theta t}[0.5\cdot\log \det(\Lambda^{-1})$$
$$- 0.5\cdot(x_n - A'\cdot z_n')^T\cdot\Lambda^{-1}\cdot(x_n - A'\cdot z_n')]$$
$$+ \delta \quad (22)$$

Instead of intractable $\mathcal{L}(\cdot)$, maximize Eq. (22) via two steps: an expectation step (E-step) and a maximization step (M-step). The E-step evaluates $Q$ given from previous M-step resultant $\Theta$, and the maximization step (M-step) finds maximizing $\Theta$ of the E-step resultant $Q$. The iterative algorithm altering the E-step and the M-step has been named an EM algorithm, and the algorithm can never decrease $\mathcal{L}(\cdot)$ as iteration proceeds (Dempster *et al.*, 1977), in other words, $\Theta_t$ converges to likelihood-sense (local) optimal $\Theta$ as $t$ increases. There might be local optimal since there is no way to find a global optimum solution for a general nonlinear functions. The local region of the solution candidates of the function can be extended via a genetic algorithm, a Boltzmann machine, and so on.

**E-Step:** To evaluate $Q(\Theta|\Theta_t)$ let us approximate $p(x^m|z, x^o) \approx p(x^m|\underline{z})$ where $\underline{z} \equiv \mathcal{E}[z|x^o]$. It implies that the observed vector is only used to infer the factor

scores which are responsible for all generated signals from the concerned process. In generative latent variable models, e.g. FA, probabilistic PCA, hidden markov model, mixture modeling, measurement scores, no matter how many they are, are treated just as shadows of latent variables' combinations. With the approximation, the conditional joint density over two unknowns, given a known is factorized to

$$p(z, x^{\mathrm{m}}|x^{\mathrm{o}}) \approx p(z|x^{\mathrm{o}}) \cdot p(x^{\mathrm{m}}|z) \qquad (23)$$

Since we knows $p(z|x^{\mathrm{o}}) = \mathcal{N}(z|x^{\mathrm{o}}: \boldsymbol{\mu}^{\mathrm{z|o}}, \boldsymbol{\Sigma}^{\mathrm{z|o}})$ and $p(x^{\mathrm{m}}|z) = \mathcal{N}(x^{\mathrm{m}}|z: \boldsymbol{\mu}^{\mathrm{m|z}}, \boldsymbol{\Sigma}^{\mathrm{m|z}})$, Eq. (23) is evaluated using Eqs. (11)–(14). Furthermore, using this factorization, the expected elements which implicitly existed in the $Q$ function are easily derived as follows:

$$\underline{z} \equiv \mathcal{E}_{z,xm|xo}[z] \approx \boldsymbol{\mu}^{\mathrm{z|o}} \qquad (24)$$

$$\underline{zz} \equiv \mathcal{E}_{z,xm|xo}[z \cdot z^{\mathrm{T}}] \approx \boldsymbol{\Sigma}^{\mathrm{z|o}} + \boldsymbol{\mu}^{\mathrm{z|o}} \cdot \boldsymbol{\mu}^{\mathrm{z|oT}} \qquad (25)$$

$$\underline{x}^{\mathrm{m}} \equiv \mathcal{E}_{z,xm|xo}[x^{\mathrm{m}}] \approx \boldsymbol{\mu}^{\mathrm{m|z}} \qquad (26)$$

$$\underline{xx}^{\mathrm{m}} \equiv \mathcal{E}_{z,xm|xo}[x^{\mathrm{m}} \cdot x^{\mathrm{mT}}] \approx \boldsymbol{\Sigma}^{\mathrm{m|z}} + \boldsymbol{\mu}^{\mathrm{m|z}} \cdot \boldsymbol{\mu}^{\mathrm{m|zT}} \qquad (27)$$

$$\underline{x}^{\mathrm{m}}\underline{z} \equiv \mathcal{E}_{z,xm|xo}[x^{\mathrm{m}} \cdot z^{\mathrm{T}}] \approx \boldsymbol{\mu}^{\mathrm{m|z}} \cdot \boldsymbol{\mu}^{\mathrm{z|oT}} \qquad (28)$$

$$\underline{z}\underline{x}^{\mathrm{m}} \equiv \mathcal{E}_{z,xm|xo}[z \cdot x^{\mathrm{mT}}] \approx \boldsymbol{\mu}^{\mathrm{z|o}} \cdot \boldsymbol{\mu}^{\mathrm{m|zT}} \qquad (29)$$

They are also calculated by Eqs. (11)–(14). Augmented forms of the expectation elements, which explicitly exist in the $Q$ function are evaluated as follows using Eqs. (24)–(29):

$$\underline{z}' \equiv \mathcal{E}_{z,xm|xo}[z'] = [\underline{z}; 1] \qquad (30)$$

$$\underline{zz}' \equiv \mathcal{E}_{z,xm|xo}[z' \cdot z'^{\mathrm{T}}] = [\underline{zz}, \underline{z}; \underline{z}^{\mathrm{T}}, 1] \qquad (31)$$

$$\underline{x} \equiv \mathcal{E}_{z,xm|xo}[x] = [x^{\mathrm{o}}; \underline{x}^{\mathrm{m}}] \qquad (32)$$

$$\underline{xx} \equiv \mathcal{E}_{z,xm|xo}[x \cdot x^{\mathrm{T}}] = [x^{\mathrm{o}} \cdot x^{\mathrm{oT}}, x^{\mathrm{o}} \cdot \underline{x}^{\mathrm{mT}}; \underline{x}^{\mathrm{m}} \cdot x^{\mathrm{oT}}, \underline{xx}^{\mathrm{m}}] \quad (33)$$

$$\underline{xz}' \equiv \mathcal{E}_{z,xm|xo}[x \cdot z'^{\mathrm{T}}] = \underline{x} \cdot \underline{z}'^{\mathrm{T}} \qquad (34)$$

$$\underline{z}'\underline{x} \equiv \mathcal{E}_{z,xm|xo}[z' \cdot x^{\mathrm{T}}] = \underline{z}' \cdot \underline{x}^{\mathrm{T}} \qquad (35)$$

Notice that the most probable inferences on all unobserved states, i.e. factors and missings, and their correlations are evaluated and used in the EM algorithm of the FA model calibration while the least square of the PCA model simply adapts the orthogonal projection of $x$ to the latent space as the best inference of the states.

**M-step:** Similar to Eqs. (19) and (20), likelihood-sense optimal $\Theta$ that maximizes $Q(\Theta|\Theta_t)$ is obtained by

$$(\partial/\partial\mathbf{A}')[Q(\Theta|\Theta_t)] = \mathbf{0}$$
$$\Rightarrow \mathbf{A}' = (\textstyle\sum_{n \in N} \underline{xz}'_n^{\mathrm{T}}) \cdot (\textstyle\sum_{n \in N} \underline{zz}'_n)^{-1} \qquad (36)$$

$$(\partial/\partial\boldsymbol{\Lambda}^{-1})[Q(\Theta|\Theta_t)] = \mathbf{0}$$
$$\Rightarrow \boldsymbol{\Lambda} = N^{-1} \cdot \mathrm{diag}\textstyle\sum_{n \in N}(\underline{xx}_n - \mathbf{A}' \cdot \underline{z}'\underline{x}_n) \qquad (37)$$

As before, $\mathbf{A}'$ in Eq. (37) should be the resultant of Eq. (36). Notice that $\underline{xx}_n \neq \underline{x}_n \cdot (\underline{x}_n)^{\mathrm{T}}$ and $\underline{zz}'_n \neq \underline{z}'_n \cdot (\underline{z}'_n)^{\mathrm{T}}$ but $\underline{xz}'_n = \underline{x}_n \cdot (\underline{z}'_n)^{\mathrm{T}}$ and $\underline{z}'\underline{x}_n = \underline{z}'_n \cdot (\underline{x}_n)^{\mathrm{T}}$ as we expressed in Eqs. (30)–(35).

EM algorithm to calibrate the FA model from incomplete data set is summarized as follows: unobserved states and their correlations are the likelihood-sense optimal inferred by E-step using Eqs. (24)–(29) given $\Theta$. Then we can update $\Theta$ by the M-step using Eqs. (30) and (31), given the inference results which have the forms represented in Eqs. (24)–(26). A lot of iterations altering the E-step and the M-step converge to a likelihood-sense (local) optimal $\Theta$ definitely. Appendix A is the summary of these procedures; it may be helpful to realize the method in computer programs, e.g. Matlab.

3.1.3 Variance explanation ratio    A covariance matrix of $x$ has been estimated by the sum of its systematic part $\mathbf{A} \cdot \mathbf{A}^{\mathrm{T}}$, and noise part $\boldsymbol{\Lambda}$, i.e. $\underline{\boldsymbol{\Sigma}}^{\mathrm{x}} = \mathbf{A} \cdot \mathbf{A}^{\mathrm{T}} + \boldsymbol{\Lambda}$, in the model. Denoting the systematic part of the $p$-th variable in $x$ as $x^{\mathrm{s}}_p$, then the variance explanation ratio of the $p$-th variable in $x$, $x_p$, is given by $r_p = \mathrm{var}(x^{\mathrm{s}}_p) \cdot \mathrm{var}(x_p)^{-1}$. Hence, the average variance explanation ratio, $r_{\mathrm{avg}}$, can be defined by the mean of all the ratios.

$$r_p = \mathrm{diag}(\mathbf{A} \cdot \mathbf{A}^{\mathrm{T}})_p \cdot [\mathrm{diag}(\mathbf{A} \cdot \mathbf{A}^{\mathrm{T}})_p + \mathrm{diag}(\boldsymbol{\Lambda})_p]^{-1} \quad (38)$$

$$r_{\mathrm{avg}} = P^{-1} \cdot \textstyle\sum_{p \in P} r_p \qquad (39)$$

where $\mathrm{diag}(\cdot)_p \in \mathbb{R}^1$ and $\mathrm{diag}(\cdot)_p^{-1} \in \mathbb{R}^1$ denote the $p$-th diagonal element in the parenthesis and its inverse, respectively. Like the score plot of the PCA model, $\dim(z)$ of the FA model can be decided by plotting $r_{\mathrm{avg}}(l)$ for $l = 1, 2, \cdots, \dim(x)$. Notice that total variance explanation ratio in the PCA model is obtained from

$$\begin{aligned}
r^{\mathrm{PCA}} &= [\textstyle\sum_{p \in P}\mathrm{var}(x^{\mathrm{s}}_p)] \cdot [\textstyle\sum_{p \in P}\mathrm{var}(x_p)]^{-1} \\
&= [\textstyle\sum_{p \in P}\mathrm{diag}(\mathbf{A} \cdot \mathbf{A}^{\mathrm{T}})_p] \cdot [\textstyle\sum_{p \in P}\mathrm{diag}(\mathbf{A} \cdot \mathbf{A}^{\mathrm{T}})_p \\
&\quad + \mathrm{diag}(\boldsymbol{\Lambda})_p]^{-1} \\
&= \|\mathbf{A}\|_{\mathrm{F}}^2 \cdot \mathrm{tr}(\underline{\boldsymbol{\Sigma}}^{\mathrm{x}})^{-1} \qquad (40)
\end{aligned}$$

where $\|\cdot\|_{\mathrm{F}}$ denotes a Frobenius norm. Thus an important difference to find an appropriate number of latent variables between FA and PCA is that FA decides the number by the summation of all the elements' ratios while PCA determine it by the ratio of the two summations:

$$r^{\text{FA}} = P^{-1} \cdot \textstyle\sum_{p \in P}[\text{var}(x^s_p) \cdot \text{var}(x_p)^{-1}] \text{ vs.}$$
$$r^{\text{PCA}} = [\textstyle\sum_{p \in P}\text{var}(x^s_p)] \cdot [\textstyle\sum_{p \in P}\text{var}(x_p)]^{-1} \qquad (41)$$

They also indicate that PCA depends on variances of variables but FA does not.

## 3.2 Prediction

For a newly measured sample $x_n = [x^i_n; x^t_n]$, similar to PCA, a target vector is predicted via two steps: infer factor scores first, and then predict target score using the inference result.

$$\underline{z}_n \equiv \mathcal{E}[z|x^i_n] = \beta_n \cdot (x^i_n - \mu^i_n) \qquad (42)$$

$$\underline{x}^t_n \equiv \mathcal{E}[x^t_n|\underline{z}_n] = \mathbf{A}^t_n \cdot \underline{z}_n + \mu^t_n \qquad (43)$$

where $\beta_n \equiv \mathbf{A}^{i\ \text{T}}_n \cdot \mathbf{\Sigma}^{ii\ -1}_n$; $\mathbf{\Sigma}^{ii\ -1}_n = (\mathbf{A}^i_n \cdot \mathbf{A}^{i\ \text{T}}_n + \mathbf{\Lambda}^i_n)^{-1}$ should be calculated regarding to $\dim(x^i_n)$ using Eq. (15). Notice that PCA simply used orthogonal projection to infer the latent score and predict the target score as expressed in Eqs. (5) and (6); but FA utilizes conditional densities for these tasks. In fact Eq. (5) is the special case of Eq. (42) when using a mean centered calibration set, $\mu^i_n = 0$, and restrict the variances of the noise variables to have infinitesimal, i.e. $\lim_{\lambda p \to 0} \mathbf{A}^{i\ +}_n = \beta_n \ \forall \ p$. And Eq. (6) is identical to Eq. (43) for the mean centered data.

## 3.3 FA model of scale invariant

Contrast to the PCA model, the FA model is invariant to both scale-change and mean-shift. For rescaled x by scaling matrix $\mathbf{B} \equiv \text{diag}(\{b_p\}_{p \in P})$, i.e. $x \to \mathbf{B} \cdot x$, the rescaling does not affect the FA model structure, it only influences model parameters: $\mathbf{A} \to \mathbf{B} \cdot \mathbf{A}$, $\mu \to \mathbf{B} \cdot \mu$, and $\mathbf{\Lambda} \to \mathbf{B} \cdot \mathbf{\Lambda} \cdot \mathbf{B}$. For mean shifted $x$ by $m$, i.e. $x \to x + m$, it indicates simply $\mu \to m \cdot \mu$. When developing the FA model, we are free from the pretreatments both on the calibration set and the tested sample, which are essential in projection models.

## 3.4 Statistical tests

An estimation error of the input part in a new sample is given by $\underline{e}^i_n = x^i_n - \underline{x}^i_n$ where $\underline{x}^i_n = \mathbf{A}^i \cdot \underline{z}_n + \mu^i$. If the developed FA model is still valid to the sample, Mahalanobis squared norms of input error (MSNE), $\underline{e}^{i\ \text{T}}_n \cdot \mathbf{\Lambda}^{i\ -1}_n \cdot \underline{e}^i_n$, should follow the chi-square distribution with $\dim(x^i_n)$ degrees of freedom since the FA model assumed $\mathcal{N}(e^i: \mathbf{0}, \mathbf{\Lambda}^i)$. Thus if Eq. (38) is satisfied for the sample, the developed FA model would be adequate to the sample with the $\alpha$ level of significance (LoS). Under the condition, both types of the variance explanation ratios, elements based ratios in Eq. (38) and the average ratio in Eq. (39), are also expected to the sample. Furthermore, the process which generates $x^i_n$ is considered as in-control with $\alpha$ LoS if Eq. (45) is held since $\mathcal{N}(e^i: \mathbf{0}, \mathbf{I})$.

$$\underline{e}^{i\ \text{T}}_n \cdot \mathbf{\Lambda}^{i-1} \cdot \underline{e}^i_n \leq \chi^{-2}_{(1-\alpha;\ \dim(xni))} \qquad (44)$$

$$\underline{z}^{\text{T}}_n \cdot \underline{z}_n \leq \chi^{-2}_{(1-\alpha,\ L)} \qquad (45)$$

where $\chi^{-2}$ denotes the inverse of cumulative chi-square PDF. Equations (44) and (45) are closely related to the Q test (Jackson and Mudholkar, 1979), and Hotelling's $\mathcal{T}^2$ test (Hotelling, 1947), respectively; they are popular tests in statistical process control (SPC) and fault detection (MacGregor and Kourti, 1995; Wise and Gallagher, 1996). Note here that Q and $\mathcal{T}^2$ tests are given by

$$\text{Q:}\ x^{\text{T}} \cdot (\mathbf{I} - \mathbf{A} \cdot \mathbf{A}^{\text{T}}) \cdot x$$
$$\leq \theta_1 \cdot [(\mathcal{N}^{-1}_{s\ (1-\alpha)} \cdot (2 \cdot \theta_2 \cdot h_0^2)^{0.5} \cdot \theta_1^{-1})$$
$$+ (\theta_2 \cdot h_0 \cdot (h_0 - 1) \cdot \theta_1^{-2}) + 1]^{1/h0} \qquad (46)$$

$$\mathcal{T}^2\text{:}\ \|\mathbf{S}_S^{-0.5} \cdot x\|^2 \leq L \cdot (N^2 - 1) \cdot (N^2 - N \cdot L)^{-1} \cdot \mathcal{F}^{-1}_{(1-\alpha,\ L,N-L)} \qquad (47)$$

where $\mathcal{N}^{-1}_s$ denotes the inverse of cumulative $\mathcal{N}(0, 1)$; using the following SVD results of $\mathbf{X}$: $[\mathbf{A} = \{a_p\}_{p \in P}$, $\mathbf{S} = \text{diag}(\{s_p\}_{p \in P})$, $\mathbf{V}] = \text{SVDs}(\mathbf{X}, P)$, $\theta_j = \sum_{i=L+1,\ldots,P} s_i^{2 \cdot j}$ for $j \in \{1, 2, 3\}$, $h_0 = 1 - (2 \cdot \theta_1 \cdot \theta_3) \cdot (3 \cdot \theta_2^2)^{-1}$, $(\mathbf{S}_S^{-0.5})^{\text{T}} = \{s_p^{-1} \cdot a_p\}_{p \in L}$.

There are several advantages on the proposed tests over the conventional tests:
(1) They provide simpler forms than the existing tests.
(2) They work with the same measuring unit, i.e. squared Mahalanobis norm. All statistical inference should be based on the Mahalanobis measuring unit; however, the left-hand term in the Q test indicates the squared Euclidian norm of the estimation errors, i.e. $\|\underline{e}\|^2$, since $\underline{e} = x - \mathbf{A} \cdot \underline{z} = (\mathbf{I} - \mathbf{A} \cdot \mathbf{A}^{\text{T}}) \cdot x$, and $(\mathbf{I} - \mathbf{A} \cdot \mathbf{A}^{\text{T}})$ is an idempotent matrix. It is the reason why the right-hand term in the test is so complicated. But the left-hand term in the $\mathcal{T}^2$ test represents the squared Mahalanobis norm of principal components' scores with the normalizing constants $\{s_p^{-1}\}_{p \in L}$. Note here that the limiting condition of $L \cdot (N^2 - 1) \cdot (N^2 - N \cdot L)^{-1} \cdot \mathcal{F}^{-1}_{(1-\alpha;\ L,\ N-L)}$ when $N \approx \infty$ is $\chi^2_{(1-\alpha;\ L)}$; however, we do not set the condition for the upper control limit of Q.
(3) They can work on partial available elements among all the measurements; while Q and $\mathcal{T}^2$ tests should wait for all measurements to be gathered. This point may be the decisive merit of the tests when applying the tests to an on-line monitored process with some malfunctioning sensors or missing values.
(4) They provide statistical normality tests for all individuals both in $\underline{z}_n = \{\underline{z}_l\}_{l \in L}$ and $\underline{e}^i_n = \{\underline{e}_p\}_{p \in \dim(xni)}$. Notice that $\mathcal{N}(\underline{z}_n: \mathbf{0}, \mathbf{I})$ and $\mathcal{N}(\underline{e}^i_n: \mathbf{0}, \text{diag}\{\lambda_p\}_{p \in \dim(xni)})$ indicate $\{\mathcal{N}(\underline{z}_l: 0, 1)\}_{l \in L}$ and $\{\mathcal{N}(\lambda_p^{-0.5} \cdot \underline{e}_p: 0, 1)\}_{p \in P}$, respectively, where $p(\underline{z}_i|\underline{z}_j) = p(\underline{z}_i)$ and $p(\underline{e}_i|\underline{e}_j) = p(\underline{e}_i)$ for $i \neq j$. Therefore if the $p$-th sensor in $x^i_n$ is fault then Eq. (48) will not be satisfied with $\alpha$ LoS. And if Eq. (49) is not hold for the extracted $l$-th factor score in $\underline{z}_n$, $\underline{z}_l$, then the factor is considered responsible for the process' out-of-control condition with $\alpha$ LoS.

$$\underline{e}_p \in [\lambda_p^{0.5} \cdot \mathcal{N}^{-1}_{s\ (0.5 \cdot \alpha)}, \lambda_p^{0.5} \cdot \mathcal{N}^{-1}_{s\ (1-0.5 \cdot \alpha)}] \qquad (48)$$
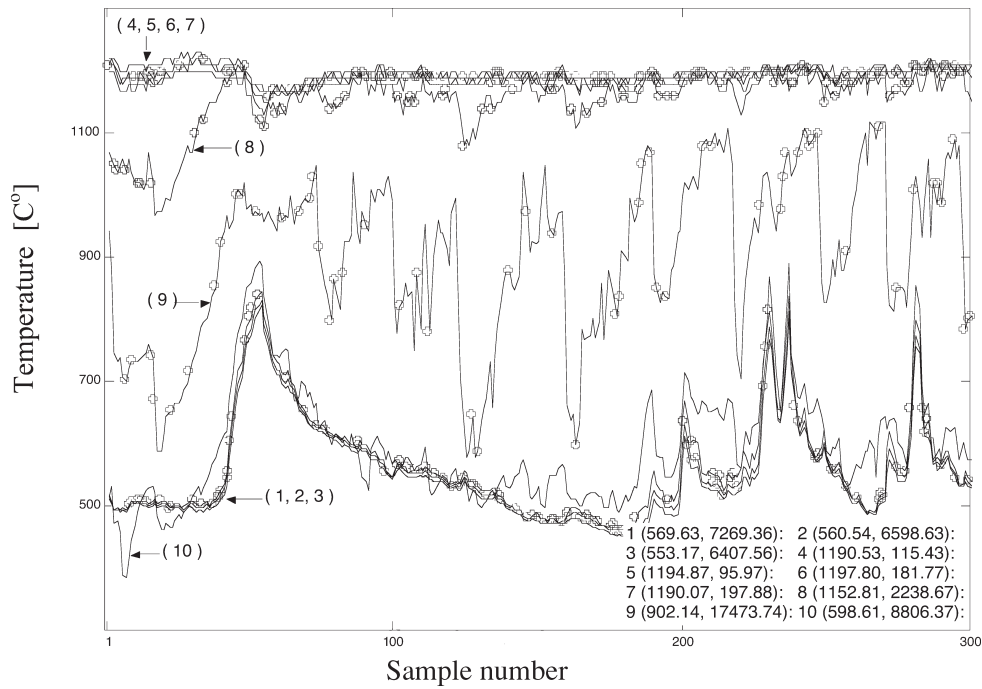
**Fig. 2** Liquid-fed ceramic melter (LFCM) data set made to have artificial 20% missing values at random

$$z_l \in [\mathcal{N}_s^{-1}{}_{(0.5 \cdot \alpha)}, \mathcal{N}_s^{-1}{}_{(1-0.5 \cdot \alpha)}] \tag{49}$$

These tests are comparable to the contribution plot used in conventional SPC.

## 4. Simulation and Results

Multivariate process data from liquid-fed ceramic melter (LFCM) data were employed to test the performance of the proposed method, where the data set are taken from the PLS_toolbox ver. 3.1 (Eigenvetor Research, Inc.). A slurry of nuclear waste and glass-forming chemicals were fed into the LFCM, and glass was periodically poured from the melter, resulting in time-dependent variations in the glass level. The LFCM was monitored with 10 thermocouples located in two thermo-wells within the glass pool (Stork *et al.*, 1997). Raw data set was treated to have 20% missing at random. The data pattern and missing values are shown in **Figure 2**, where the missing values are symbolized by '$\bigcirc$'. In the figure, $p(m, v)$ indicates that the $p$-th variable in the raw data set, without any missing value, has mean $m$ and variance $v$. As shown in the figure, every 10 variables show the strong correlation among the measured variables and also have remarkably different statistics of means and variances. If there were negligible missing data, the pretreatment step for PCA would be a simple task. However, 20% missing may not be a negligible quantity at all.

To develop a multivariate model from incomplete data set and predict the missing values in the set, we partitioned the overall incomplete set $\mathbf{X} = \{x_j\}_{j \in 300}$ into
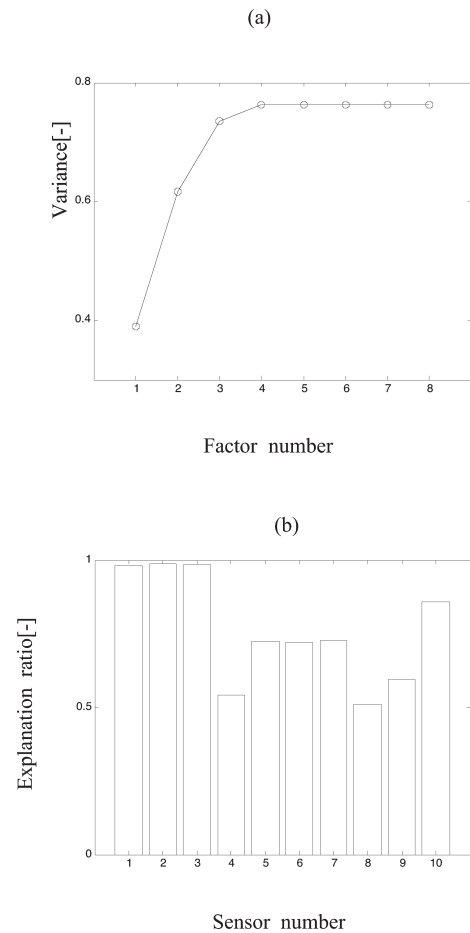
(a)



(b)



**Fig. 3** Plots of variance explanation ratios: (a) averaged variance explanation ratio, and (b) individual elements' variance explanation ratios
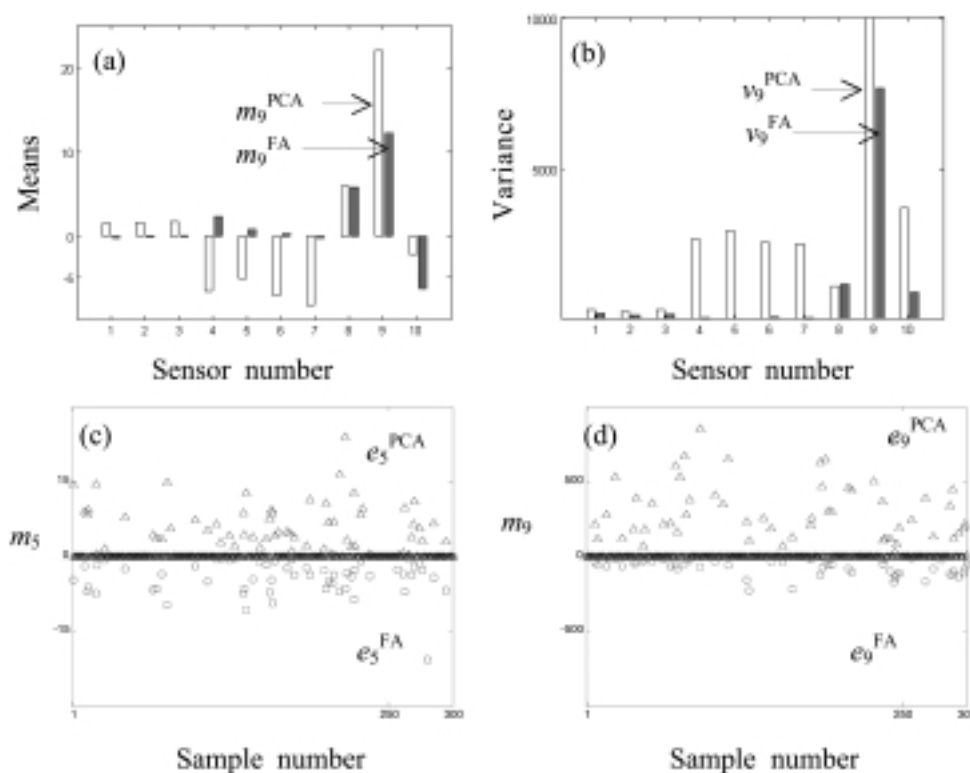
**Fig. 4** Comparisons of the prediction performance between the PCA based method and the FA based method. (a) Means of prediction errors of missing values by PCA and FA, (b) variances of prediction errors by PCA and FA, (c) minimal variance variable, $x_5$, and (d) maximal variance variable, $x_9$

two subsets: a calibration set $\mathbf{X}_{\text{cali}} = \{x_n\}_{n \in 250}$ and a test set $\mathbf{X}_{\text{test}} = \{x_n\}_{n=251, \ldots, 300}$. Multivariate models based on both PCA and FA are calibrated from the observed part in $\mathbf{X}_{\text{cali}}$, i.e. $\{x^{\text{o}}_n\}_{n \in 250}$. Using the developed models, the missing values in the $n$-th calibration sample $x^{\text{m}}_n$ by $x^{\text{o}}_n$ are estimated and the target values in the $n$-th test sample $x^{\text{t}}_n$ by $x^{\text{i}}_n$, are predicted.

Figure 3 shows the variance explanation ratios with an increased factor number and each sensor's variance explanation ratio. The appropriate number of factors is found to be four by plotting averaged variance explanation ratios as increasing the factor number from 1 to 8 in Figure 3(a). Under the four factor model structure, Figure 3(b) depicts each sensor's variance explanation ratios and 78% variances of individual elements are explained by the FA model. Notice that the ratios are arranged by $r_2 > r_3 \ldots > r_4 > r_8$; however, the order of the ratios does not depend on the variances of the variables but on the correlations between measurements and factors.

Prediction results for $\{x^{\text{m}}_n\}_{n \in 250}$ and $\{x^{\text{t}}_n\}_{n=251, \ldots, 300}$ using the PCA and the FA models are presented in **Figure 4**. Here, the means and the variances of the prediction errors estimated by PCA and FA for the missing values are shown in hollow bars and in gray bars, respectively. The mean and the variance of the $p$-th variable are shown as $m_p$ and $v_p$, respectively. These were calculated by (a) $m_p = N_{\text{pm}}^{-1} \cdot \sum_{n \in N\text{pm}} (x^{\text{m}}_{p\,n} - \underline{x}^{\text{m}}_{p\,n})$, and (b) $v_p = (N_{\text{pm}} - 1)^{-1} \cdot \sum_{n \in N\text{pm}} [(x^{\text{m}}_{p\,n} - \underline{x}^{\text{m}}_{p\,n}) - m_p]^2$, where $N_{\text{pm}}$ represents the number of missing elements in the $p$-th variable, and under bar denotes the estimate. Among all the variables, two extreme cases which are (c) minimal variance variable, $x_5$, and (d) maximal variance variable, $x_9$, are chosen for the detailed comparison. In order to simplify the comparison, ± absolute values are plotted in both (c) and (d), e.g. $e_{\text{PCA}} = \text{abs}(x^{\text{m}} - \underline{x}^{\text{m}}_{\text{PCA}})$ with a triangle symbol, $e_{\text{FA}} = -\text{abs}(x^{\text{m}} - \underline{x}^{\text{m}}_{\text{FA}})$ with a circle mark. As shown in the figure, The FA model gives better results for both biases and variances of the prediction errors than the PCA one. The results are straightforward because FA is devised to consider the probability density information of all the variables and its measuring unit is the Mahalanobis distance, while PCA is governed by Euclidian distance measure.

There are a calibrated multivariate model of the process and an incomplete new sample which is generated from the concerned process. Then there are two important questions which should be answered to monitor the process status when analyzing a new sample: the first question is whether the developed model is still valid to the sample or not, and the second question is if the model is proper to the sample, how we can decide the process condition from the sample, i.e.,
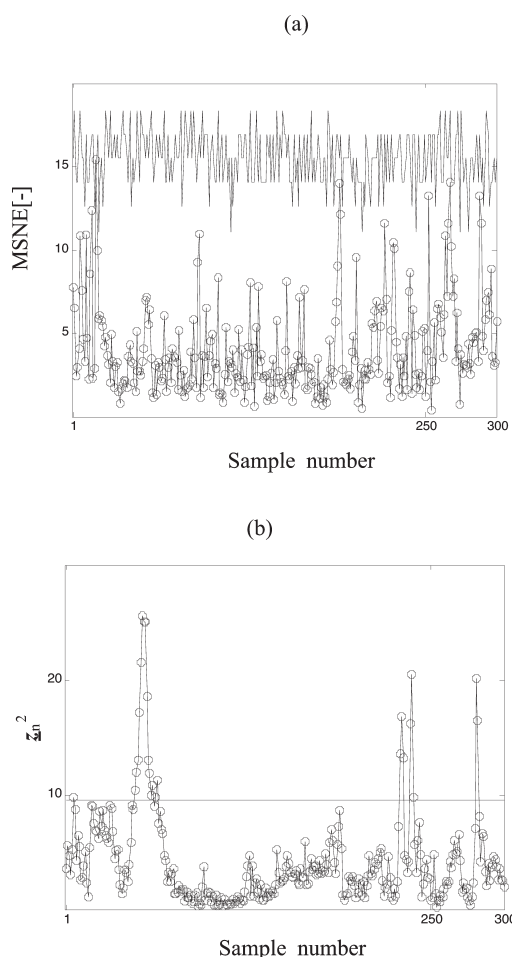
(a)



(b)



**Fig. 5** Multivariate control charts of the FA model,
(a) model error test, (b) process' in-control test

in-control or out-of-control. If there were no missing value in the tested sample, Q and $\mathcal{T}^2$ tests of the PCA model might be the answers to these questions. When the sample has any missing elements, these questions should be answered only by the available part in the sample, i.e. $x^i_n$.

When only partial information of the sample is given, **Figure 5** can be the answers for the first question and the second question, respectively. Figure 5 shows multivariate control charts of the FA model. Figure 5(a) shows the model error test of the tested samples using Mahalanobis squared norms of input error scores (MSNE), and Figure 5(b) shows in-control test using Mahalanobis squared norms of extracted factor scores. Chi-square upper control limits with solid lines were calculated regarding the 0.05% level of significance in both the charts, that is, 95% upper control limits. Notice that the 95% upper control limits in Figure 5(a) are varied according to the number of available measurements in the sample. However, the limit to check the in-control in Figure 5(b) is invariant to all the samples because of $\dim(z_j) = 4 \; \forall \; j$. As shown in Figure 5(a), both the available parts in calibration sam-

ples, $\{x^o_n\}_{n=1,...,250}$, and test samples, $\{x^i_n\}_{n=251,...,300}$, indicate that all samples are harmonized with the developed FA model with the 0.05% level of significance. Extracted factor scores from the FA model are enough to decide the current process condition. If there were some violations in the figure, we can also test which elements in the available part, $x^o_n$ or $x^i_n$, are responsible for these violations using Eq. (48). Therefore, the following general decision can come out that there is an out-of-control sample by Figure 5(b) and the propriety of the decision is confirmed by Figure 5(a).

**Conclusions**

A new calibration method of the FA model from an incomplete calibration set and a prediction method for the missing values in the calibration set and a newly measured sample are proposed. The proposed method gives better estimation results for the missing values than the well-known PCA based method. The results come from an underlying fundamental difference in each method. While PCA seeks to find a least-square sense optimal solution and a Euclidian distance is used as the measuring unit of similarities among data, FA aims to get a maximum likelihood optimal solution and the Mahalanobis distance is to be the unit. On the other hand, when they are applied to auto-scaled data, the PCA based method results similar to the FA based method since the Euclidian distance is converted to the Mahalanobis distance. However, the scaling is often impossible when the data has some missing values.

The probabilistic model such as FA is more suitable than the other projection models when a statistical decision should be made. All kinds of statistical decisions are based on statistics; more generally, it is the matter of the probability densities of variables within the model. But the model building and the decision making are completely separated in the projection model, e.g. PCA based process monitoring. In the probabilistic modeling approach, building the model and making the decision are done in one fold. Another good merit, which was not discussed in this article, of the probabilistic approaches is that it can be extended to mixture modeling, e.g. a mixture of Gaussians (Duda *et al.*, 2001), a mixture of factor analyzers (Ghahramani and Hinton, 1996), a mixture of probabilistic PCA (Tipping and Bishop, 1997), etc. In fact, there is no guarantee of assuming that all probabilistic density functions of the FA model would be Gaussians. But there is a good approximation technique of non-Gaussian density through the summations of several Gaussian densities, so called a mixture of Gaussian modeling (MOG). For a nonlinear function, there are two types of modeling techniques: (1) to model it by appropriate nonlinear function directly, or (2) to develop several local linear models for a nonlinear function, then approximate the original to the linear combinations of

the local models. Future work on the both topics, called a generic nonlinear modeling and a mixture modeling of FA, will be carried out.

## Appendix A: FA Model Calibration from Incomplete Data Set

FA model calibration from an incomplete data set is trained by the EM algorithm.

Set initial $\Theta$ then

E-step: $Q(\Theta|\Theta_t) = \mathcal{E}[\mathcal{L}(\mathbf{D}: \Theta)|\mathbf{X}^o: \Theta_t]$     (A1)

M-step: $\Theta_{t+1} = \arg_Q \max: Q(\Theta|\Theta_t)$     (A2)

Until $\Theta$ converges, iterate the E-step and the M-step:

The idea is realized as follows:

Set initial $\Theta = \{\mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Lambda}\}$ and let $\mathbf{A}' = [\mathbf{A}, \boldsymbol{\mu}]$.

Until $\Theta$ converge, iterate the followings:
For $n = 1, 2, \ldots, N$,

if $\dim(\boldsymbol{x}^o_n) \leq \dim(\boldsymbol{z})$, then $\boldsymbol{\Sigma}^{oo}{}_n^{-1} = (\mathbf{A}^o{}_n \cdot \mathbf{A}^o{}_n{}^T + \boldsymbol{\Lambda}^o{}_n)^{-1}$

if $\dim(\boldsymbol{x}^o_n) > \dim(\boldsymbol{z})$,
then $\boldsymbol{\Sigma}^{oo}{}_n^{-1} = \boldsymbol{\Lambda}^o{}_n^{-1} - \boldsymbol{\Lambda}^o{}_n^{-1} \cdot \mathbf{A}^o{}_n \cdot (\mathbf{I} + \mathbf{A}^o{}_n{}^T \cdot \boldsymbol{\Lambda}^o{}_n^{-1} \cdot \mathbf{A}^o{}_n)^{-1} \cdot \mathbf{A}^o{}_n{}^T \cdot \boldsymbol{\Lambda}^o{}_n^{-1}$

$\beta_n = \mathbf{A}^o{}_n{}^T \cdot \boldsymbol{\Sigma}^{oo}{}_n^{-1}$     (A3)

$\boldsymbol{z}_n = \beta_n \cdot (\boldsymbol{x}^o_n - \boldsymbol{\mu}^o_n): \boldsymbol{\Sigma}^{z|o} = \mathbf{I} - \beta_n \cdot \mathbf{A}^o{}_n: \underline{zz}_n = \boldsymbol{\Sigma}^{z|o}_n + \boldsymbol{z}_n \cdot \boldsymbol{z}_n{}^T$     (A4)

$\boldsymbol{z}'_n = [\boldsymbol{z}_n; 1]: \underline{zz}'_n = [\underline{zz}_n, \boldsymbol{z}_n; \boldsymbol{z}_n{}^T, 1]$     (A5)

$\underline{x}^m_n = \mathbf{A}^m_n \cdot \boldsymbol{z}_n + \boldsymbol{\mu}^m_n: \boldsymbol{\Sigma}^{m|z} = \boldsymbol{\Lambda}^m_n: \underline{xx}^m_n = \boldsymbol{\Sigma}^{m|z}_n + \underline{x}^m_n \cdot \underline{x}^m_n{}^T$     (A6)

$\underline{x}_n = [\boldsymbol{x}^o_n; \underline{x}^m_n]: \underline{xx}_n = [\boldsymbol{x}^o_n \cdot \boldsymbol{x}^o_n{}^T, \boldsymbol{x}^o_n \cdot \underline{x}^m_n{}^T; \underline{x}^m_n \cdot \boldsymbol{x}^o_n{}^T, \underline{xx}^m_n]$     (A7)

End for

$\mathbf{A}' = (\sum_{n \in N} \underline{x}_n \cdot \boldsymbol{z}'_n{}^T) \cdot (\sum_{n \in N} \underline{zz}'_n)^{-1}$     (A8)

$\boldsymbol{\Lambda} = N^{-1} \cdot \mathrm{diag} \sum_{n \in N} (\underline{xx}_n - \mathbf{A}' \boldsymbol{z}'_n \cdot \underline{x}_n{}^T)$     (A9)

$\mathbf{A} = \mathbf{A}'_{[:, 1:L]}: \boldsymbol{\mu} = \mathbf{A}'_{[:, L+1]}$     (A10)

End until

Final $\boldsymbol{\Theta} = \{\mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Lambda}\}$     (A11)

## Appendix B: Proof of Eq. (15) Using a Matrix Inversion Lemma

For any invertible $\mathbf{U}$, in this case $\mathbf{U} = \mathbf{I}$,

$$
\begin{aligned}
\mathbf{I} &= [\boldsymbol{\Lambda}^{-1} - \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1}] \cdot [\boldsymbol{\Lambda} + \mathbf{A} \cdot \mathbf{U} \cdot \mathbf{A}^T] \\
&= \mathbf{I} + \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot \mathbf{U} \cdot \mathbf{A}^T - \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \\
&\quad - \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot \mathbf{U} \cdot \mathbf{A}^T \\
&= \mathbf{I} + \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot [\mathbf{U} \cdot \mathbf{A}^T - (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \\
&\quad\quad - (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot \mathbf{U} \cdot \mathbf{A}^T] \\
&= \mathbf{I} + \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot [\mathbf{U} \cdot \mathbf{A}^T - (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{U}^{-1} \cdot \mathbf{U} \cdot \mathbf{A}^T \\
&\quad\quad - (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot \mathbf{U} \cdot \mathbf{A}^T] \\
&= \mathbf{I} + \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A} \cdot [\mathbf{I} - (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})^{-1} \cdot (\mathbf{U}^{-1} + \mathbf{A}^T \cdot \boldsymbol{\Lambda}^{-1} \cdot \mathbf{A})] \cdot \mathbf{U} \cdot \mathbf{A}^T \\
&= \mathbf{I} + \mathbf{0}
\end{aligned}
$$

## Literature Cited

Anderson, T. W.; An Introduction to Multivariate Statistical Analysis, 2nd ed., Wiley-Interscience, New York, U.S.A. (1984)

Dempster, A., N. Laird and D. Rubin; "Maximum Likelihood from Incomplete Data via the EM Algorithm," *J. Roy. Stat. Soc. B*, **39**, 1–38 (1977)

Duda, R. O., P. E. Hart and D. G. Stork; Pattern Classification, 2nd ed., Wiley-Interscience, New York, U.S.A. (2001)

Geladi, P.; "Some Recent Trends in the Calibration Literature," *Chemom. Intell. Lab. Syst.*, **60**, 211–224 (2002)

Ghahramani, Z. and G. E. Hinton; The EM Algorithm for Mixtures of Factor Analyzers, Technical Report CRG-TR-96-1, University of Toronto, Canada (1996)

Grung, B. and R. Manne; "Missing Values in Principal Component Analysis," *Chemom. Intell. Lab. Syst.*, **42**, 125–139 (1998)

Hotelling, H.; Multivariate Quality Control, Illustrated by the Air Testing of Sample Bombsights Techniques of Statistical Analysis, C. Eisenhart, M. W. Hastay and W.A. Wallis eds., McGraw-Hill, New York, U.S.A. (1947)

Jackson, J. and G. Mudholkar; "Control Procedures for Residuals Associated with Principal Component Analysis," *Techometrics*, **21**, 341–349 (1979)

Kalivas, J. H.; "Interrelationships of Multivariate Regression Methods Using Eigenvector Basis Sets," *J. Chemom.*, **13**, 111–132 (1999)

Kim, D. S. and I. Lee; "Process Monitoring Based on Probabilistic PCA," *Chem. Intell. Lab. Syst.*, **67**, 109–123 (2003)

MacGregor, J. F. and T. Kourti; "Statistical Process Control of Multivariate Processes," *Contr. Eng. Pract.*, **3**, 403–414 (1995)

Stork, C. L., D. J. Veltkamp and B. R. Kowalski; "Identification of Multiple Sensor Disturbances during Process Monitoring," *Anal. Chem.*, **69**, 5031–5036 (1997)

Tipping, M. E. and C. M. Bishop; Mixtures of Probabilistic Principal Component Analyzers, Technical Report NCRG/97/003, pp. 1–29, Aston University, Birmingham, U.K. (1997a)

Tipping, M. E. and C. M. Bishop; "Probabilistic Principal Component Analysis," *J. Roy. Stat. Soc. B*, **61**(Part 3), 611–622 (1997b)

Walczak, B. and D. L. Massart; "Dealing with Missing Data Part 1," *Chemom. Intell. Lab. Syst.*, **58**, 15–27 (2001)

Wise, B. and N. Gallagher; "The Process Chemometrics Approach to Process Monitoring and Fault Detection," *J. Process Contr.*, **6**, 329–348 (1996)