

Hybrid variable monitoring: An unsupervised process monitoring framework with binary and continuous variables [★]

Min Wang^a, Donghua Zhou^{b,a}, Maoyin Chen^a

^a*Department of Automation, Tsinghua University, Beijing 100084, China*

^b*College of Electrical Engineering and Automation, Shandong University of Science and Technology, Qingdao, 266590, China*

Abstract

Traditional process monitoring methods, such as PCA, PLS, ICA, MD *et al.*, are strongly dependent on continuous variables because most of them inevitably involve Euclidean or Mahalanobis distance. With industrial processes becoming more and more complex and integrated, binary variables also appear in monitoring variables besides continuous variables, which makes process monitoring more challenging. The aforementioned traditional approaches are incompetent to mine the information of binary variables, so that the useful information contained in them is usually discarded during the data preprocessing. To solve the problem, this paper focuses on the issue of hybrid variable monitoring (HVM) and proposes a novel unsupervised framework of process monitoring with hybrid variables including continuous and binary variables. HVM is addressed in the probabilistic framework, which can effectively exploit the process information implicit in both continuous and binary variables at the same time. In HVM, the statistics and the monitoring strategy suitable for hybrid variables with only healthy state data are defined and the physical explanation behind the framework is elaborated. In addition, the estimation of parameters required in HVM is derived in detail and the detectable condition of the proposed method is analyzed. Finally, the superiority of HVM is fully demonstrated first on a numerical simulation and then on an actual case of a thermal power plant.

Key words: Process monitoring, Healthy state data, Hybrid variables, Fault detection.

1 Introduction

Process monitoring is indispensable because it is the premise and guarantee for the safe and stable running of industrial systems [13,2,15,3,39,33]. In recent decades, a large number of data-driven approaches have been proposed for process monitoring [12,22,5,26,11,38,32,4,37,41]. However, most of them are highly based on continuous variables because they can't avoid involving Euclidean or Mahalanobis distance and can't be utilized for hybrid variables (containing continuous and binary variables) [37].

Among data-driven methods, principal component analysis (PCA) has received continuous attention once it was applied in process monitoring due to its effectiveness of data dimensionality reduction [21,12]. Based on PCA, dynamic PCA (DPCA) adopted the technol-

ogy of time lag shift to construct augmented matrix to mine time-related information [22]. Considering slowly changing in normal process, recursive PCA (RPCA) was proposed for adaptive process monitoring [27]. In order to capture nonlinear property, kernel PCA (KPCA) was developed [31,5]. Unlike PCA, partial least squares (PLS) and its variants pay much attention to quality-related fault [29,26]. To weaken the Gaussian hypothesis, independent component analysis (ICA) was proposed for process monitoring [25]. The Mahalanobis distance (MD) can also be directly used for process monitoring [19]. As understanding of the fault initiation becomes more and more thorough, the moving window methods also be proposed for incipient fault detection [18,32,30]. Considering the practical applicability in industrial processes, a large number of improved methods have been developed for multimode and nonstationary monitoring [40,42,17].

The aforementioned methods have made remarkable achievements in process monitoring, but almost all methods are based on Euclidean or Mahalanobis distance and are highly dependent on continuous variables.

[★] This paper was not presented at any IFAC meeting. Corresponding author: Donghua Zhou, Maoyin Chen.

Email addresses: m-wang18@mails.tsinghua.edu.cn (Min Wang), zdh@tsinghua.edu.cn (Donghua Zhou), mychen@tsinghua.edu.cn (Maoyin Chen).

However, the practical industrial processes sometimes have not only continuous variables, but also binary variables which may carry some useful information for process monitoring [37] and are usually deleted in the data preprocessing [14]. For hybrid variables, Langseth *et al.* used hybrid Bayesian networks for estimating human reliability [24]. Aguilera *et al.* developed the naïve Bayes (NB) and tree augmented naïve Bayes (TAN) models and applied to species distribution [1]. Zhu *et al.* considered the mixture of continuous and discrete variables in semantic model [43]. Talvitie *et al.* introduced a related model through employing an adaptive discretization approach for structure learning in Bayesian networks when there are both continuous and discrete variables [34]. Recently, Wang *et al.* utilized continuous and binary (two-valued) variables to detect the abnormalities of thermal power plant for the first time [37]. Then a more effective anomaly monitoring model named feature weighted mixed naïve Bayes model (FWMNB) was developed [36].

However, the hybrid variable approaches mentioned above are supervised methods and require both normal and fault data during training. Unfortunately, the systems in actual industrial processes are running without fault in most time, and the determinations of fault samples requires repeated research and careful discussion by experts, which are time-consuming and costly. So that the healthy state samples are usually available and it is difficult to collect sufficient fault instances, which is one of the reasons why monitoring methods only based on normal working condition data, such as PCA, PLS, ICA *et al.*, have attracted much attention. Process monitoring methods with hybrid variables only based on healthy state data are very urgently. Therefore, this paper focuses on hybrid variable process monitoring and proposes a novel unsupervised framework of process monitoring with hybrid variables named HVM which can simultaneously capture the process information of both continuous and binary variables. The main contributions are summarized as follows:

- (1) The article firstly focuses on hybrid variable monitoring only based on healthy state data. And a novel unsupervised framework of process monitoring with hybrid variables (continuous and binary variables) named HVM is proposed.
- (2) Under the unsupervised framework, the statistics and the monitoring strategy suitable for hybrid variables are firstly defined and the physical explanation behind the framework is elaborated. In addition, the expressions of parameters are derived in detail and the detectable condition is analyzed.
- (3) The effectiveness and efficiency of the proposed method is fully demonstrated first on a numerical

simulation and then on a practical fan system of ultra-supercritical power plant.

The remainder of the paper is organized as follows. The problem formulation and motivation are described in detail in Section 2. The framework of hybrid variable monitoring is introduced in Section 3. In Section 4, parameters learning and corresponding derivation are described. The fault form of hybrid variables is defined and the detectable condition is analyzed in Section 5. In Section 6, the effectiveness and efficiency of proposed framework is verified. Finally, conclusions are given in Section 7.

2 Problem formulation and motivation

With industrial processes becoming more and more complex and integrated, binary variables also appear in monitoring variables. For example, in Zhejiang Zheneng Zhongmei Zhoushan Coal and Electricity Co., Ltd. (Zhoushan Power Plant), Zhejiang Province, China, the number of monitoring variables in the No.1 power unit is about 17380, in which the number of binary variables among them is as many as 8820 [37]. In the fan system of No.1 power unit, 260 continuous variables and 495 binary variables are collected, where the number of binary variables is more than that of continuous variables [36]. The appearance of binary variables makes traditional monitoring approaches no longer applicable and process monitoring with hybrid variables more intractable. The binary variables are usually discarded during the data preprocessing because the traditional approaches mostly have applied Euclidean or Mahalanobis distance which can't be used to describe binary variables [14]. However, binary variables may carry some useful information for process monitoring [37,36].

The issue of supervised classification with hybrid variables has been paid attention to and investigated in other fields [24,1,43,34]. In process monitoring, Wang *et al.* have utilized continuous and binary variables for the anomaly detection of thermal power plant [37,36]. However, these approaches are supervised methods, which require both normal samples and fault instances to train the model. In practical processes, a lots of healthy state samples can be collected and it is difficult to obtain sufficient faulty samples. Therefore, this paper proposes a novel unsupervised framework of process monitoring with continuous and binary variables named HVM. HVM can simultaneously mine the information of both continuous and binary variables through a probabilistic framework. In HVM, the statistics of hybrid variables are computed with healthy state data and the control limit is determined by kernel density estimation (KDE) [28]. Then for the arriving sample \mathbf{x}_a , the statistic s_a can be computed with the same way of training. The

state of \mathbf{x}_a can be determined through the monitoring strategy. Finally the superiority of HVM is demonstrated through a numerical simulation and an actual case in the fan system of a thermal power plant.

3 Hybrid variable monitoring framework

3.1 Off-line statistics

Training data $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ are sampled under normal operating condition with n samples. $\mathbf{x}_i \in \mathbb{R}^d$ is the i th instance and contains d ($d = d_b + d_c$) features where d_b binary features and d_c continuous features are respectively collected. Let x^j be the j th variable. j_b and j_c mean the j_b th and j_c th variable of binary variables and continuous variables respectively. When the system is running in a steady state, the monitoring data tends to be stationary and with no trends [41]. Then the following assumptions are introduced.

Assumption 1 If x^j is a continuous variable (denoted as x^{j_c}), we suppose it obeys Gaussian distribution under normal condition, that is [37]

$$P_c(x^{j_c}; \boldsymbol{\theta}^{j_c}) = \mathcal{N}(x^{j_c}; \mu^{j_c}, \sigma^{j_c}), \quad (1)$$

where $\boldsymbol{\theta}^{j_c} = \{\mu^{j_c}, \sigma^{j_c}\}$, $\mathcal{N}(x^{j_c}; \mu^{j_c}, \sigma^{j_c})$ is the probability density function (pdf) defined as $\mathcal{N}(x^{j_c}; \mu^{j_c}, \sigma^{j_c}) = (2\pi)^{-1/2} (\sigma^{j_c})^{-1} \exp(-(x^{j_c} - \mu^{j_c})^2 2^{-1} (\sigma^{j_c})^{-2})$, μ^{j_c} and $(\sigma^{j_c})^2$ are the mean and corresponding variance of the j_c th variable.

Assumption 2 If x^j is a binary variable (denoted as x^{j_b}), the Bernoulli distribution is introduced as follows [9]:

$$P_b(x^{j_b}; \boldsymbol{\theta}^{j_b}) = (\eta^{j_b})^{x^{j_b}} (1 - \eta^{j_b})^{1-x^{j_b}}, \quad (2)$$

where $\boldsymbol{\theta}^{j_b} = \{\eta^{j_b}\}$, $P_b(x^{j_b}; \boldsymbol{\theta}^{j_b})$ is the distribution series (ds), η^{j_b} is the response probability which is defined as $\eta^{j_b} = P(x^{j_b} = 1)$.

Definition 1 In practical processes, variables are often correlated with each other. Then the occurrence probability of \mathbf{x}_i under normal condition is defined as

$$P(\mathbf{x}_i; \boldsymbol{\theta}) = \prod_{j_c=1}^{d_c} P_c(\mathbf{x}_i^{j_c}; \boldsymbol{\theta}^{j_c})^{\varphi^{j_c}} \prod_{j_b=1}^{d_b} P_b(x^{j_b}; \boldsymbol{\theta}^{j_b})^{\varphi^{j_b}}, \quad (3)$$

where $\boldsymbol{\theta} = \{\boldsymbol{\theta}^{j_c}, \boldsymbol{\theta}^{j_b}, \varphi^{j_c}, \varphi^{j_b}\}$, φ means the weight of the corresponding variable.

Affected by noise, there may be some outliers in data sampled under normal operating condition. Then the probability that \mathbf{x}_i belongs to \mathbf{X} can be obtained by

$$P(\mathbf{x}_i) = \tilde{\delta} P(\mathbf{x}_i; \boldsymbol{\theta}), \quad (4)$$

where $\tilde{\delta}$ is the prior normal probability, which represents the confidence level of the health state data and equals to $\tilde{\delta} = 1 - \delta$, δ is the significance level [10].

Proposition 1 $\forall \mathbf{x}_i \in \mathbf{X}$, \exists a positive decimal α ($0 < \alpha < 1$) to satisfy $\alpha \leq P(\mathbf{x}_i) < 1$.

Proof. For $\mathbf{x}_i \in \mathbf{X}$, suppose Assumption 1 holds and $\varphi^{j_c} > 0$ (which can be obtained by Definition 3, where $\mathcal{M}(x^j, x^{j'})$ is non-negative.), then

$$0 < \prod_{j_c=1}^{d_c} P_c(\mathbf{x}_i^{j_c}; \boldsymbol{\theta}^{j_c})^{\varphi^{j_c}} < 1. \quad (5)$$

Since the number of training samples n is an integer less than infinity, Assumption 2 is introduced, and $\varphi^{j_b} > 0$ (which can be obtained by Definition 3.), we have

$$0 \leq \prod_{j_b=1}^{d_b} P_b(x^{j_b}; \boldsymbol{\theta}^{j_b})^{\varphi^{j_b}} \leq 1. \quad (6)$$

Then $0 < P(\mathbf{x}_i) < 1$ for any $\mathbf{x}_i \in \mathbf{X}$. There must be a positive value ϱ that satisfies

$$0 < \varrho \leq P(\mathbf{x}_i) < 1. \quad (7)$$

The prior normal probability $0 < \tilde{\delta} < 1$, so that a positive decimal $0 < \alpha < 1$ can be found to satisfy $\alpha \leq P(\mathbf{x}_i) < 1$, where $\alpha = \varrho \tilde{\delta}$. \square

Remark 1 $P_c(x^{j_c}; \boldsymbol{\theta}^{j_c})$ and $P_b(x^{j_b}; \boldsymbol{\theta}^{j_b})$ are probability distributions (pdf or ds), which are fitted by training data. Thus the more \mathbf{x}_i deviates from the statistical characteristics of \mathbf{X} , the smaller $P(\mathbf{x}_i; \boldsymbol{\theta})$ is and the smaller $P(\mathbf{x}_i)$ is.

Definition 2 When $P(\mathbf{x}_i)$ of \mathbf{x}_i is obtained, then $f(\mathbf{x}_i)$ is computed as

$$f(\mathbf{x}_i) = \ln(P(\mathbf{x}_i)), \quad (8)$$

where $\ln(\cdot)$ is the natural logarithmic function.

Proposition 2 Compared to $P(\mathbf{x}_i)$, $f(\mathbf{x}_i)$ obtained in equation (8) is more sensitive to faulty instance.

Proof. According to **Proposition 1**, a lower bound α ($0 < \alpha < 1$) that satisfies $P(\mathbf{x}_i) \in [\alpha, 1)$ for normal data \mathbf{x}_i can be found. For a natural logarithmic function $f(\mathbf{x}_i) = \ln P(\mathbf{x}_i)$, $f(\mathbf{x}_i)$ monotonically increases and the derivative $\frac{\partial f(\mathbf{x}_i)}{\partial P(\mathbf{x}_i)} = \frac{1}{P(\mathbf{x}_i)}$ always satisfy that $\frac{1}{P(\mathbf{x}_i)} > 1$ when $0 < P(\mathbf{x}_i) < 1$. Fault data \mathbf{x}_f often deviates more from the statistical characteristics of \mathbf{X} and $0 < P(\mathbf{x}_f) < \alpha$. The detection performance is mainly reflected in the recognition ability of fault in the neighborhood $U(\alpha, \epsilon)$ of α , where $U(\alpha, \epsilon) = \{P(\mathbf{x}_i) | \alpha - \epsilon < P(\mathbf{x}_i) < \alpha + \epsilon\}$. Since $0 < \alpha - \epsilon < P(\mathbf{x}_i) < \alpha + \epsilon < 1$, so $f(\mathbf{x}_i)$ is more sensitive to faulty instance than $P(\mathbf{x}_i)$. The transformation of the natural logarithmic function is shown in Fig. 1. For the normal sample \mathbf{x}_1 and faulty sample \mathbf{x}_2 , $0 < P(\mathbf{x}_2) < \alpha < P(\mathbf{x}_1) < 1$ and $\alpha - P(\mathbf{x}_2) = P(\mathbf{x}_1) - \alpha$. $f(\mathbf{x}_1)$, $f(\mathbf{x}_2)$ and *threshold* are obtained from $P(\mathbf{x}_1)$, $P(\mathbf{x}_2)$ and α with natural logarithmic transformation, respectively. According to the properties of the natural logarithm function, $f(\mathbf{x}_1)$ -*threshold* $<$ *threshold* $- f(\mathbf{x}_2)$. \square

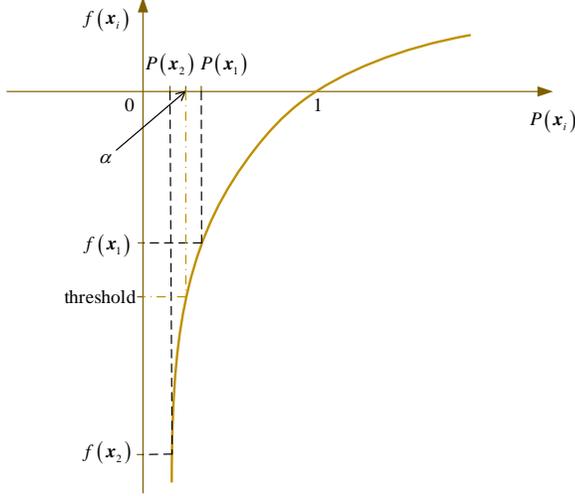


Fig. 1. Transformation of the natural logarithmic function.

According to equation (3), (4) and (8), $f(\mathbf{x}_i)$ can be written as

$$\begin{aligned} f(\mathbf{x}_i) &= \ln(\tilde{\delta}P(\mathbf{x}_i; \boldsymbol{\theta})) = \ln(\tilde{\delta}) + \ln(P(\mathbf{x}_i; \boldsymbol{\theta})) \\ &= \ln \tilde{\delta} + \ln\left(\prod_{j_c=1}^{d_c} P_c(x^{j_c}; \boldsymbol{\theta}^{j_c})^{\varphi^{j_c}} \prod_{j_b=1}^{d_b} P_b(x^{j_b}; \boldsymbol{\theta}^{j_b})^{\varphi^{j_b}}\right). \end{aligned} \quad (9)$$

Let $\Psi = \ln\left(\prod_{j_c=1}^{d_c} P_c(x^{j_c}; \boldsymbol{\theta}^{j_c})^{\varphi^{j_c}} \prod_{j_b=1}^{d_b} P_b(x^{j_b}; \boldsymbol{\theta}^{j_b})^{\varphi^{j_b}}\right)$, it can be learned that

$$\Psi = \sum_{j_b=1}^{d_b} \varphi^{j_b} \ln P_b(x^{j_b}; \boldsymbol{\theta}^{j_b}) + \sum_{j_c=1}^{d_c} \varphi^{j_c} \ln P_c(x^{j_c}; \boldsymbol{\theta}^{j_c}). \quad (10)$$

Considering equation (2), we have

$$\begin{aligned} \sum_{j_b=1}^{d_b} \varphi^{j_b} \ln P_b(x^{j_b}; \boldsymbol{\theta}^{j_b}) &= \sum_{j_b=1}^{d_b} \varphi^{j_b} \ln[(\eta^{j_b})^{x^{j_b}} (1 - \eta^{j_b})^{1-x^{j_b}}] \\ &= \sum_{j_b=1}^{d_b} \varphi^{j_b} [x^{j_b} \ln(\eta^{j_b}) + (1 - x^{j_b}) \ln(\tilde{\eta}^{j_b})] \\ &= \sum_{j_b=1}^{d_b} [\varphi^{j_b} x^{j_b} \ln \frac{\eta^{j_b}}{\tilde{\eta}^{j_b}}] + \sum_{j_b=1}^{d_b} \varphi^{j_b} \ln \tilde{\eta}^{j_b}, \end{aligned} \quad (11)$$

where $\tilde{\eta}^{j_b} = 1 - \eta^{j_b}$. According to equation (1), the following equation can be obtained that

$$\begin{aligned} \sum_{j_c=1}^{d_c} \varphi^{j_c} \ln P_c(x^{j_c}; \boldsymbol{\theta}^{j_c}) &= \sum_{j_c=1}^{d_c} \varphi^{j_c} \ln \mathcal{N}(x^{j_c}; \mu^{j_c}, \sigma^{j_c}) \\ &= \sum_{j_c=1}^{d_c} \varphi^{j_c} \ln[(2\pi)^{-1/2} (\sigma^{j_c})^{-1}] \\ &\quad + \sum_{j_c=1}^{d_c} \varphi^{j_c} [-(x^j - \mu^j)^2 2^{-1} (\sigma^j)^{-2}]. \end{aligned} \quad (12)$$

Substituting equation (11) and (12) into equation (9), $f(\mathbf{x}_i)$ is learned as

$$f(\mathbf{x}_i) = \boldsymbol{\tau}_i \cdot \tilde{\mathbf{x}}_i^T + \xi_i + \varepsilon_i, \quad (13)$$

where $\boldsymbol{\tau}_i = [\vartheta^1, \dots, \vartheta^{j_b}, \dots, \vartheta^{d_b}]$, $\vartheta^{j_b} = \varphi^{j_b} \ln \frac{\eta^{j_b}}{\tilde{\eta}^{j_b}}$, $\varepsilon_i =$

$$\sum_{j_c=1}^{d_c} \varphi^{j_c} \left[\ln((2\pi)^{-\frac{1}{2}} (\sigma^{j_c})^{-1}) - \frac{1}{2} (x_i^{j_c} - \mu^{j_c})^2 (\sigma^{j_c})^{-2} \right],$$

$$\tilde{\mathbf{x}}_i = [x_i^1, \dots, x_i^{j_b}, \dots, x_i^{d_b}], \xi_i = \ln(1 - \delta) + \sum_{j_b=1}^{d_b} \varphi^{j_b} \ln \tilde{\eta}^{j_b}.$$

$f(\mathbf{x}_i)$ obtained in equation (13) is negative. The statistics in process monitoring are often positive, and the judgment logic is generally that the statistics of the faulty data exceed the control limit. Thus for the collected training samples \mathbf{X} , the monitoring statistics \mathbf{s} are computed as

$$\begin{aligned} \mathbf{s} &= [s_1, \dots, s_i, \dots, s_n] \\ &= [f^2(\mathbf{x}_1), \dots, f^2(\mathbf{x}_i), \dots, f^2(\mathbf{x}_n)], \end{aligned} \quad (14)$$

where s_i is the statistic of \mathbf{x}_i .

3.2 On-line monitoring strategy

When the statistics \mathbf{s} of \mathbf{X} are obtained, the control limit s^{lim} can be got with the significance level δ by KDE [28], $\delta = 0.01$ in this paper. In online detection, the statistic s_a of arriving sample \mathbf{x}_a is computed by equation (13) and (14). Then the state of \mathbf{x}_a is determined through the monitoring strategy:

$$\begin{cases} \mathbf{x}_a \text{ is normal,} & \text{if } s_a < s^{\text{lim}}, \\ \mathbf{x}_a \text{ is faulty,} & \text{otherwise.} \end{cases} \quad (15)$$

Remark 2 Only continuous and binary variables are considered in this work. The Bernoulli distribution is introduced for binary variable which has only two values, where 0 and 1 can also denote two state such as high or low. It should be noted that the idea and the skills for binary variables in this work can be referenced to discrete variables with more than two values. Then the Bernoulli distribution should be replaced by the multinomial distribution and the subsequent processing of the model may also need to be adjusted.

4 Parameters learning

The model described in 3.1 mainly involves the estimation of parameters μ^j , σ^j , η^j , and φ^j . μ^j , σ^j and η^j can be obtained through maximum likelihood estimation (MLE) [6].

$$u^j = \sum_{i=1}^n x_i^j / n, \quad (16)$$

$$\sigma^j = \left\{ \sum_{i=1}^n (x_i^j - u^j)^2 \right\}^{1/2} (n-1)^{-1/2}, \quad (17)$$

$$\eta^j = \sum_{i=1}^n x_i^j / n, \quad (x_i^j \in \{0, 1\}). \quad (18)$$

In practical process, variables are usually correlated with others and variables that are more related to the other variables are more sensitive when abnormalities occur [37]. So each variable is assigned with the different feature weight φ^j [20]. The mutual information (MI) can capture the dependence of variables, both linear and non-linear [8], and is used to construct the feature weight φ^j which is defined as follows.

Definition 3 For the j th variable, the weight φ^j is defined as

$$\varphi^j = 1 + \frac{1}{d-1} \sum_{j=1, j \neq j'}^d \mathcal{M}(x^j, x^{j'}). \quad (19)$$

where $\mathcal{M}(x^j, x^{j'})$ is the MI of x^j and $x^{j'}$.

If x^j and $x^{j'}$ are continuous variables or both are binary variables, $\mathcal{M}(x^j, x^{j'})$ can be obtained by the definition of MI for continuous variables or discrete variables. However, x^j and $x^{j'}$ may include both continuous and binary variables. Then the auxiliary binary variable is constructed by Definition 4 for the continuous variable when the feature weight is computed.

Definition 4 If x^j is a continuous variable, x'^j is constructed as

$$x'^j_i = [x^j_i > \mu^j], \quad (20)$$

where $[\cdot]$ is Iverson brackets. If the condition $x^j_i > \mu^j$ is true, it returns 1, otherwise it returns 0.

Definition 4 makes it possible to characterize the correlation between hybrid variables. Then x'^j instead of x^j is used to compute MI. However, if x^j and $x^{j'}$ are continuous variables, $\mathcal{M}(x^j, x^{j'})$ and $\mathcal{M}(x'^j, x'^{j'})$ are not completely equivalent. The relationship between $\mathcal{M}(x^j, x^{j'})$ and $\mathcal{M}(x'^j, x'^{j'})$ is shown in Theorem 1.

Theorem 1 If x^j and $x^{j'}$ are continuous variables, x^j and $x^{j'}$ obey Gaussian distributions $\mathcal{N}(\mu^j, (\sigma^j)^2)$ and $\mathcal{N}(\mu^{j'}, (\sigma^{j'})^2)$ respectively, x^j and $x^{j'}$ are constructed by equation (20), the relationship between $\mathcal{M}(x'^j, x'^{j'})$ and $\mathcal{M}(x^j, x^{j'})$ is

$$\begin{aligned} \mathcal{M}(x'^j, x'^{j'}) &= \left(\frac{1}{\pi} \arcsin \rho + 0.5\right) \log\left(\frac{2}{\pi} \arcsin \rho + 1\right) \\ &+ \left(0.5 - \frac{1}{\pi} \arcsin \rho\right) \log\left(1 - \frac{2}{\pi} \arcsin \rho\right), \end{aligned} \quad (21)$$

where ρ is the correlation coefficient of continuous variables x^j and $x^{j'}$, and is expressed as $\rho = [1 - e^{-2\mathcal{M}(x^j, x^{j'})}]^{1/2}$.

Proof. See Lemma 1 and appendix A. \square

Lemma 1 [35] For continuous variables x^j and $x^{j'}$ that follow Gaussian distributions, $\mathcal{M}(x^j, x^{j'})$ is the MI of x^j and $x^{j'}$. ρ is the correlation coefficient between x^j and $x^{j'}$. Then

$$\rho = [1 - e^{-2\mathcal{M}(x^j, x^{j'})}]^{1/2}. \quad (22)$$

The Lemma 1 is proved in [7].

With Definition 4, the MI computation of hybrid variables is transformed to that of binary variables (or constructed binary variables). $\mathcal{M}(x^j, x^{j'})$ is defined as

$$\mathcal{M}(x^j, x^{j'}) = \sum_{x^j, x^{j'}} P(x^j, x^{j'}) \log \frac{P(x^j, x^{j'})}{P(x^j)P(x^{j'})}, \quad (23)$$

where $x^j, x^{j'}$ are binary variables or constructed binary variables. $P(x^j)$ is the probability of $x^j = \psi_{x^j}$, ψ_{x^j} is the indicative coefficient ($\psi_{x^j} = 1$ when $P(x^j = 1)$ is computed, and $\psi_{x^j} = 0$ otherwise), $P(x^j, x^{j'})$ is the joint probability of $x^j = \psi_{x^j}$ and $x^{j'} = \psi_{x^{j'}}$. $P(x^j)$ ($P(x^{j'})$) can be obtained in the same way can be computed by

$$P(x^j) = \psi_{x^j} \frac{\sum_{i=1}^n x^j_i}{n} + (1 - \psi_{x^j}) \left(1 - \frac{\sum_{i=1}^n x^j_i}{n}\right), \quad (24)$$

where x^j_i is the value at time i of x^j .

Proposition 3 For binary variables x^j and $x^{j'}$, $P(x^j, x^{j'})$ can be denoted as

$$\begin{aligned} P(x^j, x^{j'}) &= P(x^{j'} = \psi_{x^{j'}}) \varsigma^{\psi_{x^j} \psi_{x^{j'}}} (1 - \varsigma)^{\psi_{x^{j'}} - \psi_{x^j} \psi_{x^{j'}}} \\ &\times \varsigma'^{\psi_{x^j} - \psi_{x^j} \psi_{x^{j'}}} (1 - \varsigma')^{1 + \psi_{x^j} \psi_{x^{j'}} - \psi_{x^j} - \psi_{x^{j'}}}, \end{aligned} \quad (25)$$

where $P(x^j = 1 | x^{j'} = 1) = \varsigma$, $P(x^j = 1 | x^{j'} = 0) = \varsigma'$.

Proof. See appendix B. \square

Theorem 2 For binary variables x^j and $x^{j'}$, $P(x^j, x^{j'})$ is obtained as

$$\begin{aligned} P(x^j, x^{j'}) &= P(x^{j'} = \psi_{x^{j'}}) \{1 - \psi_{x^j} + (2\psi_{x^j} - 1) \\ &\times [\psi_{x^{j'}} \varsigma + (1 - \psi_{x^{j'}}) \varsigma']\}, \end{aligned} \quad (26)$$

where $\varsigma = \sum_{i=1}^n (x^j_i x^{j'}_i) (\sum_{i=1}^n x^{j'}_i)^{-1}$, $\varsigma' = (\sum_{i=1}^n x^j_i - \sum_{i=1}^n x^j_i x^{j'}_i) (n - \sum_{i=1}^n x^{j'}_i)^{-1}$.

Proof. See appendix C. \square

After $\mathcal{M}(x^j, x^{j'})$ is estimated, the weight φ^j of j th variable could be obtained through (19).

Remark 3 When the correlation between variables is not considered, that is $\varphi^j = 1$, all variables have the same weight.

5 Detectability analysis

5.1 Fault description

In multivariate statistical process monitoring, the fault model is usually described as

$$X^f = X + \Xi F, \quad (27)$$

where Ξ is the fault direction vector, F represents the fault magnitude vector[2,32]. The emergence of binary variables makes that the fault model described in equation (27) is no longer suitable. Thus the fault model of hybrid variables is defined as follows.

Definition 5 The fault model of hybrid variables (containing continuous and binary variables) is defined as

$$\mathbf{X}^f = \mathbf{X} + \boldsymbol{\Xi} \circ \mathbf{F}, \quad (28)$$

where \mathbf{X} is the healthy state data, $\boldsymbol{\Xi}$ means the fault direction matrix, \mathbf{F} represents the fault magnitude matrix, \circ is Hadamard product [16] which means the corresponding elements of $\boldsymbol{\Xi}$ and \mathbf{F} are multiplied, \mathbf{X}^f is the fault data.

Remark 4 When only continuous variables are monitored, fault can be described as multiplying with the direction vector and the amplitude vector. In fact, equation (27) is a special case of equation (28).

The fault model at time i in Definition 5 is

$$(\mathbf{X}^f)_i = (\mathbf{X})_i + \boldsymbol{\Xi}_i \circ \mathbf{F}_i, \quad (29)$$

where $(\mathbf{X}^f)_i, (\mathbf{X})_i, \boldsymbol{\Xi}_i, \mathbf{F}_i$ are $1 \times d$ vectors. If there is a fault at time i and it occurs on the j th ($1 \leq j \leq d$) variable, then $\boldsymbol{\Xi}_i^j = 1, \mathbf{F}_i^j$ means the corresponding fault amplitude, otherwise $\boldsymbol{\Xi}_i^j = 0$ or $\mathbf{F}_i^j = 0$. Faults can also appear on multiple variables.

Remark 5 If the j th variable is a binary variable, the fault amplitude \mathbf{F}^j must be 1 or -1 . \mathbf{F}^j must be -1 when $(\mathbf{X})^j = 1$, and \mathbf{F}^j must be 1 when $(\mathbf{X})^j = 0$.

5.2 Detectability conditions

According to the monitoring strategy (15), the state of \mathbf{x}_a is judged as fault if the statistic s_a exceed the control limit s^{lim} . The detectable condition is shown as the following Theorem.

Theorem 3 For the arriving sample \mathbf{x}_a , it can be judged to be faulty if and only if

$$P(\mathbf{x}_a; \boldsymbol{\theta}) < \tilde{\delta} e^{\tilde{s}}, \quad (30)$$

where $\tilde{s} = -\sqrt{s^{lim}}$, $P(\mathbf{x}_a; \boldsymbol{\theta}) = \prod_{j_c=1}^{d_c} \mathcal{N}(\mathbf{x}_a^{j_c}; \mu^{j_c}, \sigma^{j_c})^{\varphi^{j_c}}$
 $\prod_{j_b=1}^{d_b} [(\eta^{j_b})^{\mathbf{x}_a^{j_b}} (1 - \eta^{j_b})^{1 - \mathbf{x}_a^{j_b}}]^{\varphi^{j_b}}$, $\mathcal{N}(\mathbf{x}_a^{j_c}; \mu^{j_c}, \sigma^{j_c}) =$
 $(2\pi)^{-1/2} (\sigma^{j_c})^{-1} \exp(-(\mathbf{x}_a^{j_c} - \mu^{j_c})^2 2^{-1} (\sigma^{j_c})^{-2})$.

Proof. The fault occurs if the statistic s_a of \mathbf{x}_a exceeds the control limit s^{lim} . According to Proposition 1 and Proposition 2, the state of \mathbf{x}_a is judged as fault when $0 < P(\mathbf{x}_a) < \alpha$. Let $\ln^2 \alpha = s^{lim}$, we have

$$\alpha = e^{-\sqrt{s^{lim}}} \quad (31)$$

Then the following inequality can be obtained

$$0 < P(\mathbf{x}_a) < e^{\tilde{s}} \quad (32)$$

where $\tilde{s} = -\sqrt{s^{lim}}$. When the significance level δ is given, $\tilde{\delta} = 1 - \delta$. It can be learned that $P(\mathbf{x}_a) < \tilde{\delta} e^{\tilde{s}}$. According to Assumption 1 and 2, Theorem 3 is proved. \square

The procedure is summarized in Algorithm 1.

Algorithm 1 HVM

Off-line modeling:

Step 1 Identify continuous and binary variables.

Step 2 Estimate the means μ^j and the standard deviation σ^j for each continuous variable via (16) and (17).

Step 3 Estimate the response functions η^j for each binary variable via (18).

Step 4 Give the significance level δ and the confidence level $\tilde{\delta} = 1 - \delta$.

Step 5 Construct x^{jj} for each continuous variable x^j according to Definition 4.

Step 6 Estimate probability $P(x^j)$ and joint probability $P(x^j, x^{j'})$ via (24) and (26).

Step 7 Estimate MI $\mathcal{M}(x^j, x^{j'})$ between x^j and $x^{j'}$ via (23).

Step 8 Estimate weight φ^j via (19).

Step 9 Calculate $f(\mathbf{x}_i)$ of \mathbf{x}_i via (13).

Step 10 Calculate statistics \mathbf{s} of \mathbf{X} via (14).

Step 11 Calculate control limit s^{lim} through KDE.

On-line monitoring:

Step 12 Construct $\tilde{\mathbf{x}}_a$ through the arriving sample \mathbf{x}_a .

Step 13 Calculate ε_a and ξ_a .

Step 14 Calculate $f(\mathbf{x}_a)$ of \mathbf{x}_a via (13).

Step 15 Calculate statistic s_a of \mathbf{x}_a via (14).

Step 16 Determine the state of sample \mathbf{x}_a via (15).

6 Experimental verification

In this section, the superiority of HVM is demonstrated through two cases.

6.1 Numerical case

The fault model is considered as follows:

$$\mathbf{X}^f = \mathbf{X} + \boldsymbol{\Xi} \circ \mathbf{F} \quad (33)$$

where $\mathbf{X}, \boldsymbol{\Xi}, \mathbf{F}, \mathbf{X}^f \in \mathbb{R}^{n \times d}$. $\boldsymbol{\Xi} = \mathbf{0}$ in the normal working condition. \mathbf{X} contains 5 continuous variables (x^1, \dots, x^5) and 5 binary variables (x^6, \dots, x^{10}). 4000 normal samples are generated for training. Then 4000 instances are collected for verifying the effectiveness and efficiency of the proposed model. The fault is introduced from time 2001.

Two experiments are conducted. Under normal con-

Table 1
The distributions of continuous variables.

	Experiment I		Experiment II	
	normal(\mathbf{X})	fault (\mathbf{F})	normal (\mathbf{X})	fault (\mathbf{X}^f)
x^1	$\mathcal{N}(1.35, 0.66^2)$	$\mathcal{N}(0.15, 0.66^2)$	$\mathcal{N}(1.50, 0.76^2)$	$\mathcal{N}(0.55, 0.55^2)$
x^2	$\mathcal{N}(2.65, 0.80^2)$	$\mathcal{N}(0.05, 0.78^2)$	$\mathcal{N}(3.00, 0.68^2)$	$\mathcal{N}(2.55, 1.01^2)$
x^3	$\mathcal{N}(0.86, 0.66^2)$	$\mathcal{N}(0.10, 0.60^2)$	$\mathcal{N}(1.70, 0.85^2)$	$\mathcal{N}(2.20, 1.00^2)$
x^4	$\mathcal{N}(1.80, 0.90^2)$	$\mathcal{N}(0.15, 0.89^2)$	$\mathcal{N}(0.80, 1.01^2)$	$\mathcal{N}(1.45, 0.91^2)$
x^5	$\mathcal{N}(0.99, 0.55^2)$	$\mathcal{N}(0.30, 0.58^2)$	$\mathcal{N}(0.89, 0.64^2)$	$\mathcal{N}(1.30, 0.55^2)$

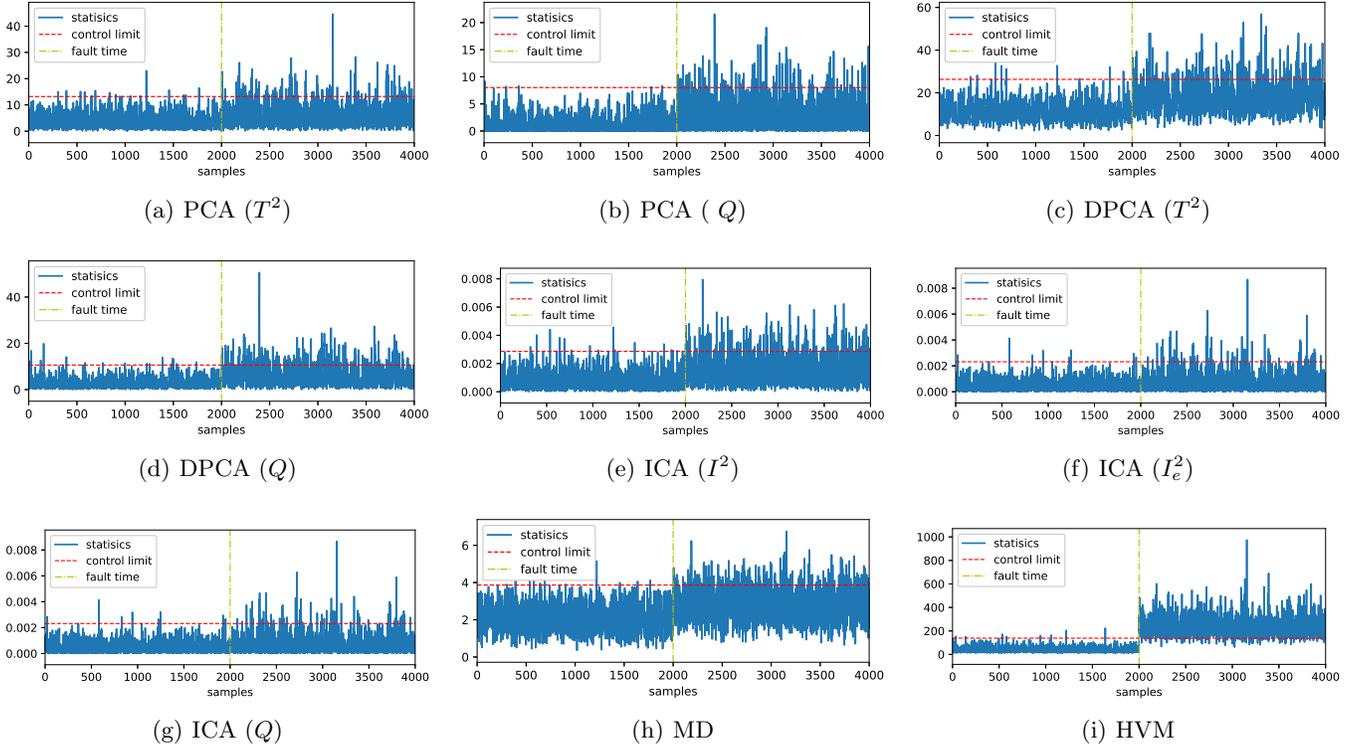


Fig. 2. Detection performance in experiment II.

dition in experiment I, the distributions of continuous variables and the values and ratios of binary variables are shown as normal (\mathbf{X}) in Table 1 and Table 2 respectively. A fault occurred from the time 2001, the continuous variables were disturbed by the Gaussian noises whose distributions is shown as fault (\mathbf{F}) in Table 1, the ratios of binary variables after fault arriving are listed as fault (\mathbf{F}) in Table 2. In experiment II, the process information carried by the binary variable is increased, and the difference in the distributions of continuous variables under normal and fault conditions is narrowed. The continuous variable distributions before fault occurring are assumed as normal (\mathbf{X}) of experiment II in Table 1, the

Table 2
The parameters of binary variables.

	Experiment I				Experiment II			
	normal (\mathbf{X})		fault (\mathbf{F})		normal (\mathbf{X})		fault (\mathbf{X}^f)	
	value	ratio	value	ratio	value	ratio	value	ratio
x^6	0	5	0	50	0	10	1	5
x^7	0	6	0	45	0	5	1	10
x^8	1	12	0	38	1	15	0	10
x^9	1	2	0	35	1	8	0	15
x^{10}	1	8	0	48	1	10	0	8

Table 3

The means of FARs and FDRs in the numerical study.

	PCA		DPCA		ICA			MD	HVM
	T^2	Q	T^2	Q	I^2	I_e^2	Q		
FAR _I	0.94	0.98	0.92	0.89	0.87	0.95	0.90	0.93	0.62
FDR _I	7.45	13.77	10.26	35.21	8.97	14.08	8.59	16.63	52.34
FAR _{II}	0.83	0.27	0.87	0.19	0.72	0.90	0.62	0.64	0.60
FDR _{II}	1.47	2.40	0.86	6.47	3.22	4.96	0.47	6.01	94.73

values and ratios of binary variables under normal condition are listed as normal (\mathbf{X}) of experiment II in Table 2. The values of binary variables after fault significantly changed which is depicted as fault (\mathbf{X}^f) of experiment II in Table 2. In order to make it more general, random jumps are added on binary variables and the adjustment ratio is shown in Table 2. Random jump means that the value changes at a time and recovers at the next moment. The distributions of continuous variables after fault are listed as fault (\mathbf{X}^f) of experiment II in Table 1.

According to the above parameters, 100 independent repeated experiments were conducted. Some classic methods, such as PCA [21,12], DPCA [22], ICA [25], MD [19] are used to verifying the effectiveness of HVM. For PCA and DPCA, the cumulative percent variance (CPV) is 0.80. Generally speaking, a larger CPV leads

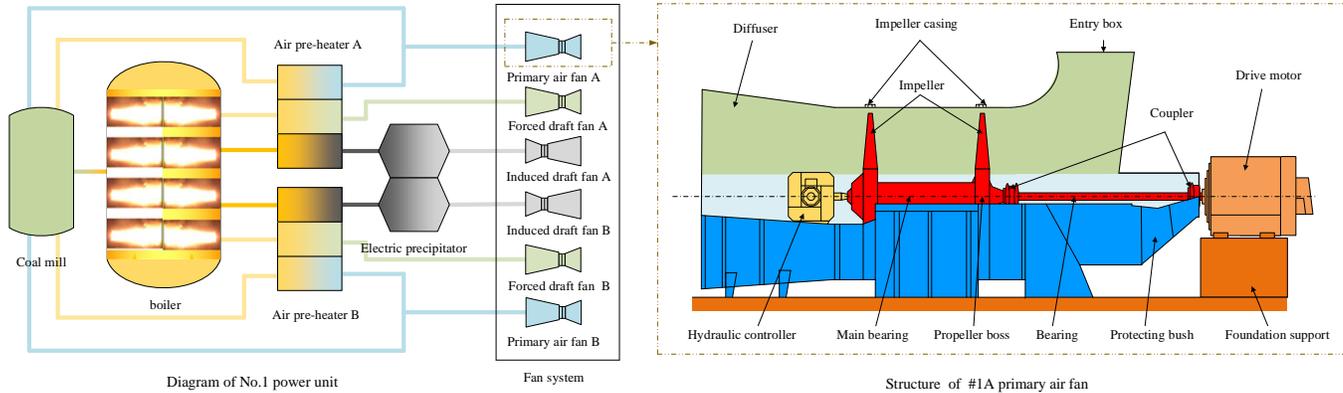


Fig. 3. Structure diagram and working condition of the primary air fan.

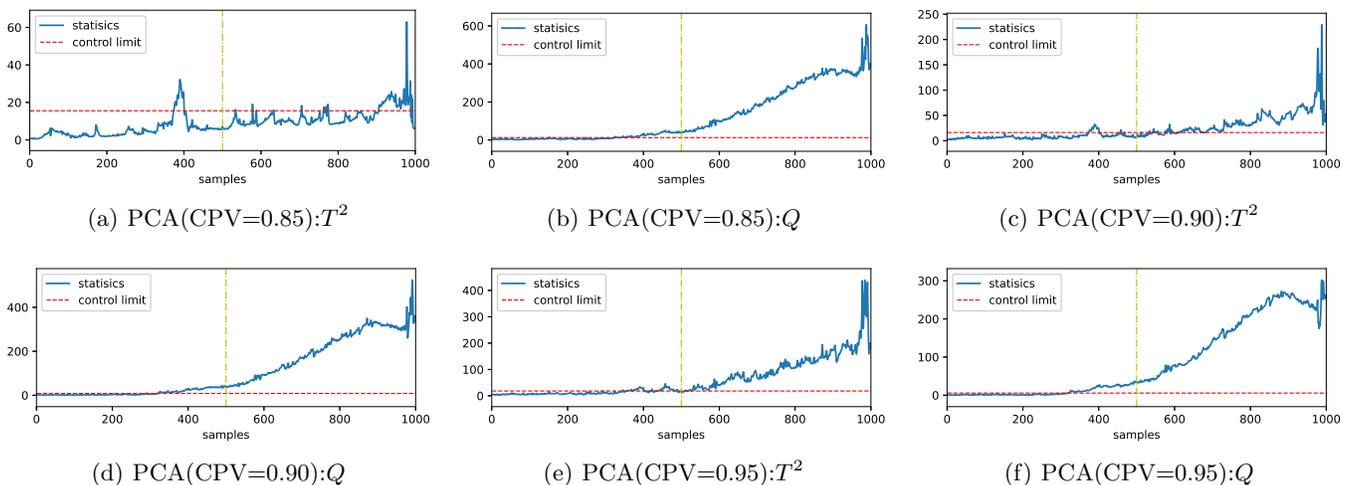


Fig. 4. Detection performance of PCA in the fan system.

to a larger number of principal components. The number of principal components is 4 in PCA when CPV is 0.80. In order to be consistent, the CPV is also 80% in DPCA. The number of principal components is 12 (the total dimension is 15). The time lag in DPCA is 2 [23]. The number of independent components (IC) in ICA equals to 3. The means of false alarm rates (FARs) and fault detection rates (FDRs) are listed in Table 3. The statistics of mentioned methods in experiment II are depicted in Fig. 2. In experiment I, the statistical characteristics of both continuous and binary variables are slightly different under normal and faulty conditions. The FDR of Q in DPCA is 35.21%. When the information carried in both continuous and binary variables is simultaneously mined, the FDR is improved to 52.34%. The difference in the distributions of continuous variables under normal and fault conditions is narrowed, and the difference of binary variables is more significant in experiment II. The best FDR of traditional methods with continuous variables is just 6.47%. But the FDR

of HVM is 94.73%, and the FAR is only 0.60%.

6.2 Fan system of the power plant

The ultra-supercritical thermal power plants have made great contributions to the development of society and still play a pivotal role in the current power system [30]. Efficient process monitoring is the foundation of continuous and stable operation for power plants. In this case, the effectiveness and efficiency of PVM is verified by an actual data collected from Zhoushan Power Plant, Zhejiang Province, China. The number of variables monitoring the No.1 power unit in Zhoushan power plant is more than 17380, and the number of binary variables among them is as many as 8820 [37]. In the fan system of the No.1 power unit, 260 continuous variables and 495 binary variables are collected, where the number of binary variables is more than that of continuous variables [36].

A vibration fault of the #1A primary air fan in the No.1 power unit occurred on September 3, 2017. The pri-

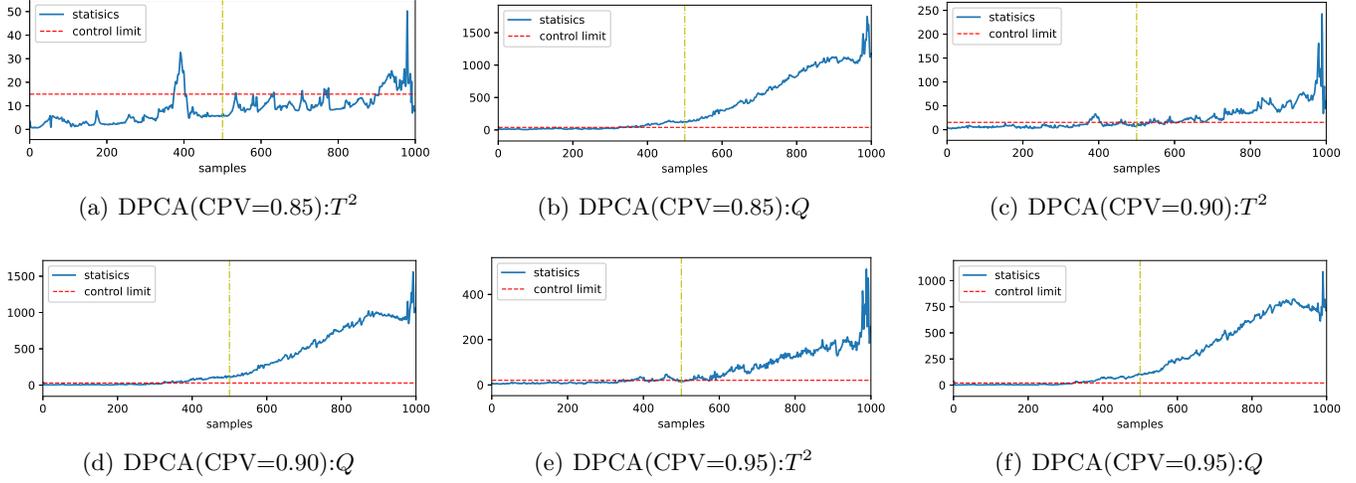


Fig. 5. Detection performance of DPCA in the fan system.

Table 4
The FARs and FDRs in the fan system

methods	PCA(CPV=0.80)		PCA(CPV=0.85)		PCA(CPV=0.90)		PCA(CPV=0.95)		DPCA(CPV=0.80)		DPCA(CPV=0.85)	
	T^2	Q	T^2	Q	T^2	Q	T^2	Q	T^2	Q	T^2	Q
FAR(%)	3.80	34.40	5.40	35.60	6.20	37.20	22.40	36.40	3.80	33.80	5.80	35.20
FDR(%)	19.40	100.00	19.80	100.00	83.80	100.00	97.60	100.00	17.20	100.00	20.20	100.00

methods	DPCA(CPV=0.90)		DPCA(CPV=0.95)		ICA(IC=10)			ICA(IC=15)			MD	HVM
	T^2	Q	T^2	Q	I^2	I_e^2	Q	I^2	I_e^2	Q		
FAR(%)	7.80	36.00	19.60	36.20	39.80	42.00	39.60	42.10	38.60	39.80	42.00	4.20
FDR(%)	86.60	100.00	95.60	100.00	100.00	100.00	99.80	100.00	100.00	100.00	100.00	100.00

Table 5
The number of principal components in PCA and DPCA.

methods	PCA				DPCA			
	80	85	90	95	80	85	90	95
number	2	3	5	7	2	3	5	8

mary air fan is the driving force for the transportation of pulverized coal, and provides hot air for the drying and oxygen for the combustion of pulverized coal. The working environment and structure diagram of the primary air fan are shown in Fig. 3. According to the recommendation of the practical engineers, 35 continuous variables and 35 binary variables are sampled every 5 seconds. 1000 instances under normal condition are used for modeling. 500 samples before and after the fault are collected respectively to test the effectiveness and efficiency.

For traditional monitoring models, PCA [21,12], DPCA [22], ICA [25] and MD [19] are adopted for process monitoring with continuous variables. Experiments

were conducted with CPV equals to 0.80, 0.85, 0.90, 0.95 respectively for PCA and DPCA, where T^2 and Q statistic are calculated. The number of principal components in PCA and DPCA is listed in Table 5. For DPCA, the time lag is 2 [23]. The FARs and FDRs of T^2 statistics keep increasing with the increase of CPV for PCA and DPCA. The best results appear on T^2 statistics of PCA and DPCA at CPV=0.9. The FAR and FDR of PCA with CPV=0.9 are 6.20% and 83.80% respectively. For DPCA, the FAR and FDR with CPV=0.9 are 7.80% and 86.70% respectively. The monitoring charts of PCA with CPV equals to 0.85, 0.90 and 0.95 are shown in Fig. 4. The statistics of DPCA when CPV is 0.85, 0.90 and 0.95 are depicted in Fig. 5. In ICA, IC=10 and IC=15 are considered. The results show that the monitoring performances of ICA are similar with different IC. The statistics of ICA when IC=15 is shown in Fig. 6(b), 6(c) and 6(g). The detection performance of MD can be seen in Fig. 6(a). The logarithmic statistics of I^2 , I_e^2 in ICA and MD are shown in Fig. 6(d), 6(e) and 6(f). The FDRs of MD and ICA are satisfactory, but the FARs is

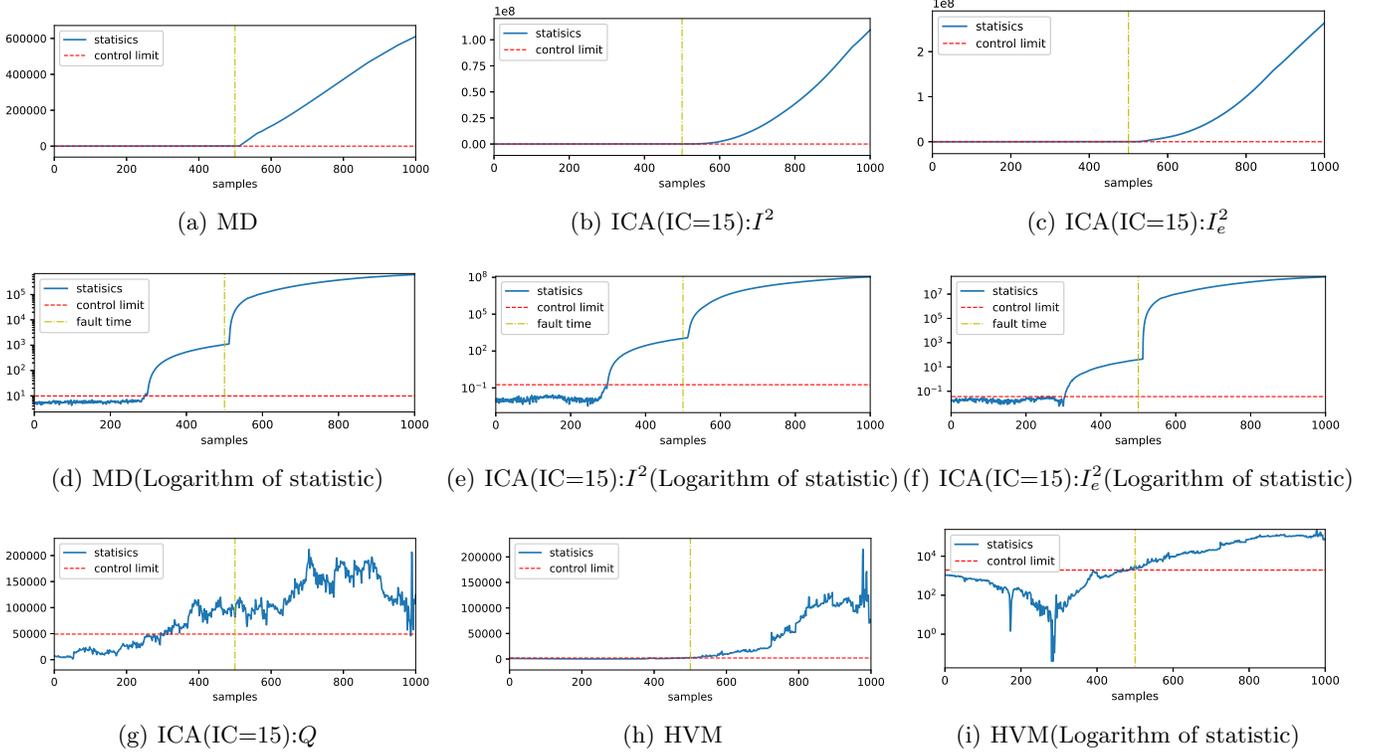


Fig. 6. Detection performance of MD, ICA and HVM in the fan system.

too high to be accepted. However, the FDR of HVM is 100% and FAR of HVM is 4.2% when both continuous and binary variables are utilized. Continuous variables contain current, air volume, vibration, temperature *etc.* of the fans. Binary variables mainly including control command signal, vibration over-limit signal, bearing vibration danger signal, moving blade position feedback signal, state signal, *etc.* are taken into consideration in HVM. Variables that are more strongly correlated with other variables tend to change easily when any other variable changes. In this case, the vibration-related variables have relatively larger weights. The detection performance of HVM is depicted in Fig. 6(h) and 6(i). The FARs and FDRs of all methods are listed in Table 4.

7 Conclusions

This paper focuses on the issue of hybrid variable monitoring only based on healthy state data and proposes a novel unsupervised process monitoring framework for hybrid variables named PVM. The statistics suitable for hybrid variables are defined and the physical explanation behind the framework is elaborated. In addition, the estimation of parameters is derived in detail and the detectable conditions of HVM is analyzed. Finally a numerical simulation and an actual case in the plant process of thermal power are utilized to verify the effectiveness and efficiency of the proposed model. Studies

demonstrate that HVM have the superiority when the information of both continuous and binary variables are effectively utilized.

Acknowledgements

This work was supported by the National Natural Science Foundation of China under Grant 62033008, 61873143.

A Proof of Theorem 1

For continuous variables x^j and $x^{j'}$, the joint probability function of x^j and $x^{j'}$ is

$$f(x^j, x^{j'}) = (2\pi\sigma^j\sigma^{j'})^{-1}(1-\rho^2)^{-1/2} \exp\{-(2-2\rho^2)^{-1} \times [(x^j - \mu^j)^2(\sigma^j)^{-2} + (x^{j'} - \mu^{j'})^2(\sigma^{j'})^{-2} - 2\rho(x^j - \mu^j)(x^{j'} - \mu^{j'}) (\sigma^j)^{-1}(\sigma^{j'})^{-1}]\}. \quad (\text{A.1})$$

Then $P(x^{j'} = 1, x^j = 1)$ is learned as

$$P(x^{j'} = 1, x^j = 1) = P(x^j > \mu^j, x^{j'} > \mu^{j'}) = \int_{\mu^j}^{\infty} \int_{\mu^{j'}}^{\infty} f(x^j, x^{j'}) dx^j dx^{j'} = \int_0^{\infty} \int_0^{\infty} f(y^j, y^{j'}) dy^j dy^{j'}, \quad (\text{A.2})$$

where $y^j = (x^j - \mu^j)(\sigma^j)^{-1}$ and $y^{j'} = (x^{j'} - \mu^{j'})\sigma^{j' -1}$. Since

$$f(y^j, y^{j'}) = (2\pi)^{-1}(1 - \rho^2)^{-1/2} \exp\{-(2 - 2\rho^2)^{-1} \times [(y^j)^2 + (y^{j'})^2 - 2\rho y^j y^{j'}]\}. \quad (\text{A.3})$$

Thus

$$\begin{aligned} P(x^{j'} = 1, x^j = 1) &= \int_0^\infty \int_0^\infty \{(2\pi)^{-1}(1 - \rho^2)^{-1/2} \\ &\times \exp\{-(2 - 2\rho^2)^{-1}[(y^j)^2 + (y^{j'})^2 - 2\rho y^j y^{j'}]\}\} dy^j dy^{j'} \\ &= \int_0^\infty \int_0^{\pi/2} \{(2\pi)^{-1}(1 - \rho^2)^{-1/2} r \\ &\times \exp\{-(2 - 2\rho^2)^{-1}(1 - \rho \sin 2\alpha)\}\} d\alpha dr \\ &= \int_0^{\pi/2} \{(2\pi)^{-1}(1 - \rho^2)^{-1/2}(1 - \rho \sin 2\alpha)^{-1}\} d\alpha \\ &= \int_0^{\pi/2} \{(2\pi)^{-1}(1 - \rho^2)^{-1/2} \\ &\times (1 + \tan^2 \alpha - 2\rho \tan \alpha)^{-1}\} d \tan \alpha \\ &= \frac{1}{2\pi} \arcsin \rho + 0.25. \end{aligned} \quad (\text{A.4})$$

In the same way, we have

$$\begin{aligned} P(x^{j'} = 0, x^j = 0) &= P(x^j \leq \mu^j, x^{j'} \leq \mu^{j'}) \\ &= \int_{-\infty}^{\mu^j} \int_{-\infty}^{\mu^{j'}} f(x^j, x^{j'}) dx^j dx^{j'} = \int_{-\infty}^0 \int_{-\infty}^0 f(y^j, y^{j'}) dy^j dy^{j'} \\ &= \int_{-\infty}^0 \int_{-\infty}^0 \{(2\pi)^{-1}(1 - \rho^2)^{-1/2} \exp\{-(2 - 2\rho^2)^{-1} \\ &\times [(y^j)^2 + (y^{j'})^2 - 2\rho y^j y^{j'}]\}\} dy^j dy^{j'}. \end{aligned} \quad (\text{A.5})$$

Since y^j and $y^{j'}$ are Gaussian distributions, it can be obtained that

$$f(z^j, z^{j'}) = f(-y^j, -y^{j'}) = f(y^j, y^{j'}). \quad (\text{A.6})$$

where $z^j = -y^j$, $z^{j'} = -y^{j'}$. Then

$$\begin{aligned} P(x^{j'} = 0, x^j = 0) &= \int_0^\infty \int_0^\infty \{(2\pi)^{-1}(1 - \rho^2)^{-1/2} \\ &\times \exp\{-(2 - 2\rho^2)^{-1}[(z^j)^2 + (z^{j'})^2 - 2\rho z^j z^{j'}]\}\} dz^j dz^{j'} \\ &= \frac{1}{2\pi} \arcsin \rho + 0.25. \end{aligned} \quad (\text{A.7})$$

The MI $\mathcal{M}(x^{j'}, x^j)$ of $x^{j'}$ and x^j ($x^{j'}$ and x^j are

constructed through equation (20)) is defined as

$$\mathcal{M}(x^{j'}, x^j) = \sum_{x^{j'}, x^j} P(x^{j'}, x^j) \log \frac{P(x^{j'}, x^j)}{P(x^{j'})P(x^j)}. \quad (\text{A.8})$$

Since x^j is a Gaussian process, it is obvious that $P(x^j = 1) = \int_{\mu^j}^\infty x^j dx^j = 1/2$. In the same way, we have $P(x^{j'} = 1) = \int_{\mu^{j'}}^\infty x^{j'} dx^{j'} = 1/2$. Then equation (A.8) is

$$\begin{aligned} \mathcal{M}(x^{j'}, x^j) &= \sum_{x^{j'}, x^j} P(x^{j'}, x^j) \log 4P(x^{j'}, x^j) \\ &= P(x^{j'} = 0, x^j = 0) \log 4P(x^{j'} = 0, x^j = 0) \\ &\quad + P(x^{j'} = 0, x^j = 1) \log 4P(x^{j'} = 0, x^j = 1) \\ &\quad + P(x^{j'} = 1, x^j = 0) \log 4P(x^{j'} = 1, x^j = 0) \\ &\quad + P(x^{j'} = 1, x^j = 1) \log 4P(x^{j'} = 1, x^j = 1), \end{aligned} \quad (\text{A.9})$$

Let $P(x^j = 0|x^{j'} = 0) = \lambda$, $P(x^j = 1|x^{j'} = 1) = \lambda'$. Since $P(x^{j'} = 0) = P(x^{j'} = 1) = 1/2$, then

$$P(x^{j'} = 0, x^j = 0) = \frac{1}{2}\lambda, \quad (\text{A.10})$$

$$P(x^{j'} = 0, x^j = 1) = \frac{1}{2}(1 - \lambda'), \quad (\text{A.11})$$

$$P(x^{j'} = 1, x^j = 0) = \frac{1}{2}(1 - \lambda), \quad (\text{A.12})$$

$$P(x^{j'} = 1, x^j = 1) = \frac{1}{2}\lambda'. \quad (\text{A.13})$$

According to equation (A.5) and (A.7), we have

$$\lambda = \lambda' = \frac{1}{\pi} \arcsin \rho + 0.5. \quad (\text{A.14})$$

Then

$$\begin{aligned} \mathcal{M}(x^{j'}, x^j) &= 2P(x^{j'} = 0, x^j = 0) \log 4P(x^{j'} = 0, x^j = 0) \\ &\quad + 2P(x^{j'} = 0, x^j = 1) \log 4P(x^{j'} = 0, x^j = 1) \\ &= \lambda \log 2\lambda + (1 - \lambda) \log 2(1 - \lambda) \\ &= \left(\frac{1}{\pi} \arcsin \rho + 0.5\right) \log\left(\frac{2}{\pi} \arcsin \rho + 1\right) \\ &\quad + \left(0.5 - \frac{1}{\pi} \arcsin \rho\right) \log\left(1 - \frac{2}{\pi} \arcsin \rho\right), \end{aligned} \quad (\text{A.15})$$

According to equation (A.15) and lemma 1, it is learned that

$$\begin{aligned} \mathcal{M}(x^{j'}, x^j) &= \left(\frac{1}{\pi} \arcsin \rho + 0.5\right) \log\left(\frac{2}{\pi} \arcsin \rho + 1\right) \\ &\quad + \left(0.5 - \frac{1}{\pi} \arcsin \rho\right) \log\left(1 - \frac{2}{\pi} \arcsin \rho\right), \end{aligned} \quad (\text{A.16})$$

where $\rho = [1 - e^{-2\mathcal{M}(x^j, x^{j'})}]^{1/2}$. \square

B Proof of Proposition 3

Let

$$P(x^j = 1|x^{j'} = 1) = \varsigma, P(x^j = 1|x^{j'} = 0) = \varsigma', \quad (\text{B.1})$$

it has

$$P(x^j = 0|x^{j'} = 1) = 1 - \varsigma, \quad (\text{B.2})$$

$$P(x^j = 0|x^{j'} = 0) = 1 - \varsigma'. \quad (\text{B.3})$$

Then we have

$$\begin{aligned} P(x^j = \psi_{x^j}|x^{j'} = \psi_{x^{j'}}) &= \varsigma^{\psi_{x^j} \psi_{x^{j'}}} (1 - \varsigma)^{\psi_{x^j} - \psi_{x^j} \psi_{x^{j'}}} \\ &\times \varsigma'^{\psi_{x^j} - \psi_{x^j} \psi_{x^{j'}}} (1 - \varsigma')^{1 + \psi_{x^j} \psi_{x^{j'}} - \psi_{x^j} - \psi_{x^{j'}}}. \end{aligned} \quad (\text{B.4})$$

Since

$$\begin{aligned} P(x^j, x^{j'}) &= P(x^j = \psi_{x^j}, x^{j'} = \psi_{x^{j'}}) \\ &= P(x^{j'} = \psi_{x^{j'}})P(x^j = \psi_{x^j}|x^{j'} = \psi_{x^{j'}}). \end{aligned} \quad (\text{B.5})$$

Thus Proposition 3 is proved. \square

C Proof of Theorem 2

Since

$$\begin{aligned} P(x^j = x_1^j, x^{j'} = x_1^{j'}) \dots P(x^j = x_n^j, x^{j'} = x_n^{j'}) \\ = \prod_{i=1}^n P(x^{j'} = x_i^{j'})P(x^j = x_i^j|x^{j'} = x_i^{j'}). \end{aligned} \quad (\text{C.1})$$

The likelihood function is

$$\begin{aligned} \ell(\varsigma, \varsigma') &= \prod_{i=1}^n P(x^{j'} = x_i^{j'})P(x^j = x_i^j|x^{j'} = x_i^{j'}) \\ &= \prod_{i=1}^n P(x^{j'} = x_i^{j'}) \prod_{i=1}^n P(x^j = x_i^j|x^{j'} = x_i^{j'}) \\ &= \varpi \varsigma^{\sum_{i=1}^n x_i^j x_i^{j'}} (1 - \varsigma)^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}} \varsigma'^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}} \\ &\times (1 - \varsigma')^{n + \sum_{i=1}^n x_i^j x_i^{j'} - \sum_{i=1}^n (x_i^j + x_i^{j'})}, \end{aligned} \quad (\text{C.2})$$

where ϖ is a constant. Then $\frac{\partial \ell(\eta, \eta')}{\partial \eta}$ can be obtained as

$$\begin{aligned} \frac{\partial \ell(\varsigma, \varsigma')}{\partial \varsigma} &= \varpi \varsigma^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}} (1 - \varsigma)^{n + \sum_{i=1}^n x_i^j x_i^{j'} - \sum_{i=1}^n (x_i^j + x_i^{j'})} \\ &[(-1)(1 - \varsigma)^{-1} (\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}) (1 - \varsigma)^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}} \\ &\varsigma^{\sum_{i=1}^n x_i^j x_i^{j'}} + (\sum_{i=1}^n x_i^j x_i^{j'}) \varsigma^{\sum_{i=1}^n x_i^j x_i^{j'}} \varsigma^{-1} (1 - \varsigma)^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}}]. \end{aligned} \quad (\text{C.3})$$

Let $\frac{\partial \ell(\varsigma, \varsigma')}{\partial \varsigma} = 0$, it can be obtained that

$$\varsigma = (\sum_{i=1}^n x_i^j x_i^{j'}) (\sum_{i=1}^n x_i^j)^{-1}. \quad (\text{C.4})$$

In the same way, $\frac{\partial \ell(\varsigma, \varsigma')}{\partial \varsigma'}$ is

$$\begin{aligned} \frac{\partial \ell(\varsigma, \varsigma')}{\partial \varsigma'} &= \varpi \varsigma^{\sum_{i=1}^n x_i^j x_i^{j'}} (1 - \varsigma)^{\sum_{i=1}^n x_i^{j'} - \sum_{i=1}^n x_i^j x_i^{j'}} \\ &[(\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}) \varsigma'^{-1} (1 - \varsigma')^{n + \sum_{i=1}^n x_i^j x_i^{j'} - \sum_{i=1}^n (x_i^j + x_i^{j'})} \\ &\times \varsigma'^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}} + (-1)(n + \sum_{i=1}^n x_i^j x_i^{j'} - \sum_{i=1}^n (x_i^j + x_i^{j'})) \\ &\times (1 - \varsigma')^{-1} (1 - \varsigma')^{n + \sum_{i=1}^n x_i^j x_i^{j'} - \sum_{i=1}^n (x_i^j + x_i^{j'})} \varsigma'^{\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}}]. \end{aligned} \quad (\text{C.5})$$

Let $\frac{\partial \ell(\varsigma, \varsigma')}{\partial \varsigma'} = 0$, ς' can be achieved as

$$\varsigma' = (\sum_{i=1}^n x_i^j - \sum_{i=1}^n x_i^j x_i^{j'}) (n - \sum_{i=1}^n x_i^{j'})^{-1}. \quad (\text{C.6})$$

According to equation (B.1), (B.2) and (B.3), we have

$$\begin{aligned} P(x^j = \psi_{x^j}|x^{j'} = \psi_{x^{j'}}) &= \{1 - \psi_{x^j} + (2\psi_{x^j} - 1) \\ &\times [\psi_{x^{j'}} \varsigma + (1 - \psi_{x^{j'}}) \varsigma']\}. \end{aligned} \quad (\text{C.7})$$

Hence, Theorem 2 is proven. \square

References

- [1] P.A. Aguilera, A. Fernández, F. Reche, and R. Rumí. Hybrid Bayesian network classifiers: Application to species distribution models. *Environmental Modelling and Software*, 25(12):1630–1639, 2010.
- [2] Carlos F. Alcalá and S. Joe Qin. Reconstruction-based contribution for process monitoring. *Automatica*, 45(7):1593–1600, 2009.
- [3] Hongtian Chen, Bin Jiang, Ningyun Lu, and Zehui Mao. Deep PCA based real-time incipient fault detection and diagnosis methodology for electrical drive in high-speed trains. *IEEE Transactions on Vehicular Technology*, 67(6):4819–4830, 2018.
- [4] Maoyin Chen and Jun Shang. Recursive spectral meta-learner for online combining different fault classifiers. *IEEE Transactions on Automatic Control*, 63(2):586–593, 2018.
- [5] Sang Wook Choi and In-Beum Lee. Nonlinear dynamic process monitoring based on dynamic kernel PCA. *Chemical Engineering Science*, 59(24):5897–5908, 2004.
- [6] Michael Collins. *Parameter estimation for statistical parsing models: Theory and practice of distribution-free methods*. Springer Netherlands, 2004.
- [7] G. A. Darbellay. *Predictability: An Information-Theoretic Perspective*, In: *Signal Analysis and Prediction*. Springer, 1998.
- [8] G.A. Darbellay and I. Vajda. Estimation of the information by an adaptive partitioning of the observation space. *IEEE Transactions on Information Theory*, 45(4):1315–1321, 1999.
- [9] Enric Junqué de Fortuny, David Martens, and Foster Provost. Wallenius Bayes. *Machine Learning*, 107(2):1–25, 2018.

- [10] Xiaogang Deng, Xuemin Tian, Sheng Chen, and Chris J. Harris. Deep principal component analysis based on layerwise feature extraction and its application to nonlinear process monitoring. *IEEE Transactions on Control Systems Technology*, 27(6):2526–2540, 2018.
- [11] Steven X. Ding, Ying Yang, Yong Zhang, and Linlin Li. Data-driven realizations of kernel and image representations and their application to fault detection and control system design. *Automatica*, 50(10):2615–2623, 2014.
- [12] Ricardo Dunia, S. Joe Qin, Thomas F. Edgar, and Thomas J. McAvoy. Identification of faulty sensors using principal component analysis. *AIChE Journal*, 42(10), 1996.
- [13] Zhiwei Gao and Steven X. Ding. Actuator fault robust estimation and fault-tolerant control for a class of nonlinear descriptor systems. *Automatica*, 43(5):912–920, 2007.
- [14] Zhiqiang Ge, Zhihuan Song, Steven X. Ding, and Biao Huang. Data mining and analytics in the process industry: the role of machine learning. *IEEE Access*, 5:20590–20616, 2017.
- [15] Zhiqiang Ge, Zhihuan Song, and Furong Gao. Review of recent research on data-based process monitoring. *Industrial and Engineering Chemistry Research*, 52(10):3543–3562, 2013.
- [16] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge: Cambridge University Press, 1985.
- [17] Yunyun Hu and Chunhui Zhao. Fault diagnosis with dual cointegration analysis of common and specific nonstationary fault variations. *IEEE Transactions on Automation Science and Engineering*, 17(1):237–247, 2020.
- [18] Ines Jaffel, Okba Taouali, Mohamed Faouzi Harkat, and Hassani Messaoud. Moving window KPCA with reduced complexity for nonlinear dynamic process monitoring. *Isa Transactions*, 64:184–192, 2016.
- [19] Hongquan Ji, Keke Huang, and Donghua Zhou. Incipient sensor fault isolation based on augmented Mahalanobis distance. *Control Engineering Practice*, 86:144–154, 2019.
- [20] Liangxiao Jiang, Lungan Zhang, Chaoqun Li, and Jia Wu. A correlation-based feature weighting filter for naive bayes. *IEEE Transactions on Knowledge and Data Engineering*, 31(2):201–213, 2019.
- [21] James V. Kresta, John F. Macgregor, and Thomas E. Marlin. Multivariate statistical monitoring of process operating performance. *Canadian Journal of Chemical Engineering*, 69(1):35–47, 1991.
- [22] Wenfu Ku, Robert H. Storer, and Christos Georgakis. Disturbance detection and isolation by dynamic principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 30(1):179–196, 1995.
- [23] Wenfu Ku, Robert H. Storer, and Christos Georgakis. Disturbance detection and isolation by dynamic principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 30(1):179–196, 1995.
- [24] Helge Langseth, Thomas D. Nielsen, Rafael Rumi, and Antonio Salmeron. Inference in hybrid Bayesian networks. *Reliability Engineering and System Safety*, 94(10):1499–1509, 2009.
- [25] Jong-Min Lee, ChangKyoo Yoo, and In-Beum Lee. Statistical process monitoring with independent component analysis. *Journal of Process Control*, 14(5):467–485, 2004.
- [26] Gang Li, S. Joe Qin, and Donghua Zhou. Geometric properties of partial least squares for process monitoring. *Automatica*, 46(1):204–210, 2010.
- [27] Weihua Li, H. Henry Yue, Sergio Valle-Cervantes, and S. Joe Qin. Recursive PCA for adaptive process monitoring. *Journal of Process Control*, 10(5):471–486, 2000.
- [28] Poovich Phaladiganon, Seoung Bum Kim, Victoria CP Chen, and Wei Jiang. Principal component analysis-based control charts for multivariate nonnormal distributions. *Expert Systems with Applications*, 40(8):3044–3054, 2013.
- [29] S. Joe Qin. Recursive PLS algorithms for adaptive data modeling. *Computers and Chemical Engineering*, 22(4):503–514, 1998.
- [30] Yihao Qin, Yayun Yan, Hongquan Ji, and Youqing Wang. Recursive correlative statistical analysis method with sliding windows for incipient fault detection. *IEEE Transactions on Industrial Electronics*, early access, 2021.
- [31] Bernhard Scholkopf, Alexander Smola, and Klaus-Robert Muller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998.
- [32] Jun Shang, Maoyin Chen, Hongquan Ji, and Donghua Zhou. Recursive transformed component statistical analysis for incipient fault detection. *Automatica*, 80:313–327, 2017.
- [33] Yabin Si, Youqing Wang, and Donghua Zhou. Key-performance-indicator-related process monitoring based on improved kernel partial least squares. *IEEE Transactions on Industrial Electronics*, 68(3):2626–2636, 2021.
- [34] Topi Talvitie, Ralf Eggeling, and Mikko Koivisto. Learning Bayesian networks with local structure, mixed variables, and exact algorithms. *International Journal of Approximate Reasoning*, 115:69–95, 2019.
- [35] A. A. Tsonis. Probing the linearity and nonlinearity in the transitions of the atmospheric circulation. *Nonlinear Processes in Geophysics*, 8(6):341–345.
- [36] Min Wang, Li Sheng, Donghua Zhou, and Maoyin Chen. A feature weighted mixed naive Bayes model for monitoring anomalies in the fan system of a thermal power plant. *IEEE/CAA J. Autom. Sinica*, 9(4):1–9, 2022.
- [37] Min Wang, Donghua Zhou, Maoyin Chen, and Yanwen Wang. Anomaly detection in the fan system of a thermal power plant monitored by continuous and two-valued variables. *Control Engineering Practice*, 102:104522, 2020.
- [38] Shen Yin, Xianwei Li, Huijun Gao, and Okyay Kaynak. Data-based techniques focused on modern industry: An overview. *IEEE Transactions on Industrial Electronics*, 62(1):657–667, 2015.
- [39] Wanke Yu, Chunhui Zhao, and Biao Huang. Moninet with concurrent analytics of temporal and spatial information for fault detection in industrial processes. *IEEE Transactions on Cybernetics*, early access, 2021.
- [40] Kai Zhang, Kaixiang Peng, and Jie Dong. A common and individual feature extraction-based multimode process monitoring method with application to the finishing mill process. *IEEE Transactions on Industrial Informatics*, 14(11):4841–4850, 2018.
- [41] Yinghong Zhao, Xiao He, Junfeng Zhang, Hongquan Ji, Donghua Zhou, and Michael G. Pecht. Detection of intermittent faults based on an optimally weighted moving average T2 control chart with stationary observations. *Automatica*, 123:109298, 2021.
- [42] Le Zhou, Jiaqi Zheng, Zhiqiang Ge, Zhihuan Song, and Shengdao Shan. Multimode process monitoring based on switching autoregressive dynamic latent variable model. *IEEE Transactions on Industrial Electronics*, 65(10):8184–8194, 2018.

- [43] Mingmin Zhu, Sanyang Liu, and Youlong Yang. Propagation in CLG Bayesian networks based on semantic modeling. *Artificial Intelligence Review*, 38:149–162, 2012.