

A Feature Weighted Mixed Naive Bayes Model for Monitoring Anomalies in the Fan System of a Thermal Power Plant

Min Wang, *Graduate Student Member, IEEE*, Li Sheng, *Member, IEEE*, Donghua Zhou, *Fellow, IEEE*, and Maoyin Chen, *Member, IEEE*

Abstract—With the increasing intelligence and integration, a great number of two-valued variables (generally stored in the form of 0 or 1) often exist in large-scale industrial processes. However, these variables cannot be effectively handled by traditional monitoring methods such as linear discriminant analysis (LDA), principal component analysis (PCA) and partial least square (PLS) analysis. Recently, a mixed hidden naive Bayesian model (MHNBM) is developed for the first time to utilize both two-valued and continuous variables for abnormality monitoring. Although the MHNBM is effective, it still has some shortcomings that need to be improved. For the MHNBM, the variables with greater correlation to other variables have greater weights, which can not guarantee greater weights are assigned to the more discriminating variables. In addition, the conditional probability $P(x_j | x_{j'}, y = k)$ must be computed based on historical data. When the training data is scarce, the conditional probability between continuous variables tends to be uniformly distributed, which affects the performance of MHNBM. Here a novel feature weighted mixed naive Bayes model (FWMNBM) is developed to overcome the above shortcomings. For the FWMNBM, the variables that are more correlated to the class have greater weights, which makes the more discriminating variables contribute more to the model. At the same time, FWMNBM does not have to calculate the conditional probability between variables, thus it is less restricted by the number of training data samples. Compared with the MHNBM, the FWMNBM has better performance, and its effectiveness is validated through numerical cases of a simulation example and a practical case of the Zhoushan thermal power plant (ZTPP), China.

Index Terms—Abnormality monitoring, Two-valued variables, Continuous variables, FWMNBM, Thermal power plant.

I. INTRODUCTION

WITH increasing intelligence and integration, a great number of two-valued variables (generally stored as 0 or 1 value) often exist in large-scale industrial processes. For instance, 17381 variables are monitored in the No. 1 generator unit of the Zhoushan thermal power plant (ZTPP),

This work was supported by the National Natural Science Foundation of China under Grant 62033008, 61873143. (Corresponding author: Donghua Zhou, Maoyin Chen.)

Min Wang, and Maoyin Chen are with the Department of Automation, Tsinghua University, Beijing 100084, China (e-mail: m-wang18@mails.tsinghua.edu.cn, mychen@tsinghua.edu.cn).

Li Sheng is with the College of Control Science and Engineering, China University of Petroleum (East China), Qingdao 266580, China (e-mail: shengli@upc.edu.cn).

Donghua Zhou is with the College of Electrical Engineering and Automation, Shandong University of Science and Technology, Qingdao, Shandong 266590, China, and also with the Department of Automation, Tsinghua University, Beijing 100084, China (e-mail: zdh@tsinghua.edu.cn).

where two-valued variables are more than 8820. These two-valued variables mainly include status monitoring variables and numerical range variables, such as control command signals and vibration over-limit signals, which switch from one state to the other with less influence from process fluctuation noise.

In order to insure the high safety and reliability of large-scale industrial processes, the problem of monitoring anomalies becomes more and more important [1]–[5]. The timely and accurate abnormal monitoring can effectively reduce waste of resources, economic losses, and even casualties [6]–[11]. Among a large number of monitoring methods, data-driven techniques have attracted much attention with the advantages of requiring less system information and prior knowledge than model-based and expert experience methods [12]–[19]. For example, principal component analysis (PCA) and its variants have been widely used in industrial processes [20], [21]. In order to detect quality-related faults, approaches based on partial least square (PLS) analysis have been proposed [22], [23]. When the training data contains both normal and abnormal working condition samples, linear discriminant analysis (LDA) has been utilized [24]. In addition, many other machine learning methods, such as K-nearest neighbors (KNN) [25], support vector machine (SVM) [26], etc., have also been applied in abnormal monitoring.

However, the fact that two-valued variables ubiquitously exist in large-scale industrial processes presents a challenge to traditional monitoring methods. It is well known that the above mentioned methods are strongly based on continuous variables and may be not suitable for two-value variables. For example, PCA, PLS, LDA, *etc.* obtain a subspace that is convenient for monitoring through decomposition and then construct statistics or hyperplanes. But these operations are based on Euclidean distance or Mahalanobis distance, which can't effectively mine the process information of two-valued variables. Two-valued variables are usually deleted during the data preprocessing stage [27], [28]. Recently, the mixed hidden naive Bayesian model (MHNBM) was proposed for the first time to combine both two-valued and continuous variables to improve monitoring performance [28]. Although MHNBM is effective, the variables with greater correlation to other variables have greater weights, which can not guarantee that greater weights are assigned to the more discriminating variables. Moreover, the conditional probability $P(x_j | x_{j'}, y = k)$ between x_j and $x_{j'}$ under $y = k$ must be computed based

on the historical data. When training data is scarce, the conditional probability between continuous variables tends to be uniformly distributed, which will affect performance.

Motivated by the above discussions, a model known as the feature weighted mixed naive Bayes model (FWMNBM) is proposed to overcome the shortcomings of MHNBM. In FWMNBM, the variables that are more correlated to the class have greater weights which results in variables with greater differences under different working conditions contribute more to the model. Meanwhile, FWMNBM can avoid calculating the conditional probability between variables such that it can still be used when there is not enough training data. In addition, a more effective consistent characterization technique is developed for the correlation of mixed variables, and the corresponding feasibility analysis is conducted. Compared with MHNBM, FWMNBM has better performance, and the effectiveness of FWMNBM is validated through the simulations of a numerical example and a practical vibration fault case.

In this paper, the remainder is organized as follows. Some preliminaries are briefly outlined in Section II. In Section III-A, FWMNBM is elaborated on. The estimation of parameters is introduced in Section III-B. In Section IV, the effectiveness of FWMNBM is verified. Finally, conclusions are drawn in Section V.

II. PRELIMINARY

In this section, the preliminary for FWMNBM is introduced. Let $\mathbf{X} = (\mathbf{x}_i)_{1 \leq i \leq n} \in \mathbb{R}^{n \times p}$ be the training data set which contains n instances. Here $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip}) \in \mathbb{R}^p$ is the sampled data at time i . From the assumption of independent and identical distribution, it can be determined that [29]:

$$P(\mathbf{X}|y = k) = \prod_{j=1}^p P(x_j|y = k), \quad (1)$$

where x_j is the j th variable of \mathbf{X} , y denotes the label variable and $y_i \in \{1, \dots, k, \dots, K\}$ is the class of instance \mathbf{x}_i .

For a continuous variable, suppose the variable obeys Gaussian distribution defined as

$$P(x_j|y = k) = \frac{1}{\sqrt{2\pi\sigma_{kj}^2}} \exp\left(-\frac{(x_j - \mu_{kj})^2}{2\sigma_{kj}^2}\right), \quad (2)$$

where μ_{kj} is the mean of the j th variable under the k class, and σ_{kj}^2 is the corresponding variance. If x_j is a two-valued variable, the Bernoulli distribution is introduced as follows [30]:

$$P(x_j|y = k) = \theta_{kj}^{x_j} (1 - \theta_{kj})^{1-x_j}, \quad (3)$$

where $\theta_{kj} = P(x_j = 1|y = k)$ is the response probability.

When \mathbf{X} only contains two-valued variables or continuous variables, the label of a new sampled data \mathbf{x}_{new} is the maximum posterior probability denoted by

$$P(y = k|\mathbf{X} \leftarrow \mathbf{x}_{new}) = \frac{P(\mathbf{X}|y = k)P(y = k)}{\sum_{k'=1}^K P(\mathbf{X}|y = k')P(y = k')}, \quad (4)$$

where $p_k = P(y = k)$ is the prior probability, $\mathbf{X} \leftarrow \mathbf{x}_{new}$ indicates that \mathbf{X} is replaced by \mathbf{x}_{new} .

III. MAIN ALGORITHM

A. FWMNBM

If two-valued variables and continuous variables exist in \mathbf{X} at the same time, let $\mathbf{X} = [\mathbf{X}_t, \mathbf{X}_c]$, where $\mathbf{X}_c \in \mathbb{R}^{n \times p_1}$ is the continuous data set, $\mathbf{X}_t \in \mathbb{R}^{n \times p_2}$ is the two-valued variable set, where $p = p_1 + p_2$. Under the assumption that the variables are independent, and assume that the continuous variable ($x_j \in \mathbf{X}_c$) obeys the Gaussian distribution as equation (2) which is denoted by $P_c(x_j|y = k)$, and two-valued variable ($x_j \in \mathbf{X}_t$) obeys the Bernoulli distribution as equation (3) which is denoted as $P_t(x_j|y = k)$. Then equation (1) can be written into

$$P(\mathbf{X}|y = k) = \prod_{j=1}^{p_1} P_c(x_j|y = k) \prod_{j=1}^{p_2} P_t(x_j|y = k). \quad (5)$$

After all parameters μ_{kj} , σ_{kj}^2 , θ_{kj} , p_k are estimated respectively based on the historical data, the posterior probability (4) of a new sample \mathbf{x}_{new} is

$$\begin{aligned} P(y = k|\mathbf{X} \leftarrow \mathbf{x}_{new}) \\ = \frac{P(y = k) \prod_{j=1}^{p_1} P_c(x_j|y = k) \prod_{j=1}^{p_2} P_t(x_j|y = k)}{\sum_{k'=1}^K P(y = k') \prod_{j=1}^{p_1} P_c(x_j|y = k') \prod_{j=1}^{p_2} P_t(x_j|y = k')}. \end{aligned} \quad (6)$$

In the actual process, the independence assumption of variables is too stringent. Among numerous approaches alleviating this assumption, feature weighting is a better choice. Inspired by [31], the model (5) could be amended as follows:

$$\begin{aligned} P(\mathbf{X}|y = k) &= P'(\mathbf{X}|y = k) = \prod_{j=1}^p P'(x_j|y = k) \\ &= \prod_{j=1}^{p_1} P'_c(x_j|y = k) \prod_{j=1}^{p_2} P'_t(x_j|y = k), \end{aligned} \quad (7)$$

where

$$P'(x_j|y = k) = P(x_j|y = k)^{FW_j}, \quad (8)$$

where FW_j is the feature weight of the j th feature. It should be noted that variables more correlated with the class should have greater weights, and uncorrelated attributes should have smaller contributions. Therefore, FW_j should reflect the correlation between the x_j and y as accurately as possible. The mutual information (MI) $MI(x_j, y)$ is used to characterize the correlation between x_j and y . $MI(x_j, y)$ can effectively describe the correlation between x_j and y , but it also contains some correlational information between x_j and other variables (such as $x_{j'}$) because variables are coupled. Then the average feature-feature intercorrelation is introduced to compute the feature weight [32]

$$CI_j = MI(x_j, y) - \frac{1}{p-1} \sum_{j=1, j \neq j'}^p MI(x_j, x_{j'}), \quad (9)$$

where $MI(x_j, x_{j'})$ is MI between the j th and j' th variables. The value of CI_j may be negative, but the weight must be non-negative. Thus, the final form of the weight of feature (FW_j) is transformed as

$$FW_j = 1/(1 + e^{-CI_j}). \quad (10)$$

The feature weights are normalized to satisfy

$$\sum_{j=1}^p FW_j = 1. \quad (11)$$

In order to effectively describe the correlation between two-valued and continuous variables, continuous variables are processed as follows.

Definition 1. If $x_j \in \mathbf{X}_c$, let $x'_j : \{x'_{1j}, x'_{2j}, \dots, x'_{nj}\}$ be the auxiliary two-valued variable corresponding to $x_j : \{x_{1j}, x_{2j}, \dots, x_{nj}\}$, which can be obtained through clipping [33] as

$$x'_{ij} = \begin{cases} 1, & x_{ij} > \sum_{k=1}^K (u_{kj}n_k)/n, \\ 0, & x_{ij} \leq \sum_{k=1}^K (u_{kj}n_k)/n, \end{cases} \quad (12)$$

where $n_k = \sum_{i=1}^n \pi_{ik}$. Then x'_j , instead of x_j , is used to build the correlation index. The feasibility analysis is given in Appendix A.

The MI between x_j and $x_{j'}$ can be computed by [34]:

$$MI(x_j, x_{j'}) = \sum_{x_j, x_{j'}} P(x_j, x_{j'}) \log \frac{P(x_j, x_{j'})}{P(x_j)P(x_{j'})}. \quad (13)$$

Then equation (6) can be transformed to

$$\begin{aligned} P(y = k | \mathbf{X} \leftarrow \mathbf{x}_{new}) &= [P(y = k) \prod_{j=1}^{p_1} P_c(x_j | y = k)]^{FW_j} \\ &\times \prod_{j=1}^{p_2} P_t(x_j | y = k)^{FW_j} \left[\sum_{k'=1}^K (P(y = k')) \right. \\ &\times \left. \prod_{j=1}^{p_1} P_c(x_j | y = k')^{FW_j} \prod_{j=1}^{p_2} P_t(x_j | y = k')^{FW_j} \right]^{-1}. \end{aligned} \quad (14)$$

$\sum_{k'=1}^K P(y = k') \prod_{j=1}^{p_1} P_c(x_j | y = k')^{FW_j} \prod_{j=1}^{p_2} P_t(x_j | y = k')^{FW_j}$ is constant for each k , then the label of new instance (\mathbf{x}_{new}) is

$$\begin{aligned} y_{new} &= \arg \max_k P(y = k | \mathbf{X} \leftarrow \mathbf{x}_{new}) \\ &= \arg \max_k [P(y = k) \prod_{j=1}^{p_1} P_c(x_j | y = k)]^{FW_j} \\ &\quad \times \prod_{j=1}^{p_2} P_t(x_j | y = k)^{FW_j} \\ &= \arg \max_k \{ \ln P(y = k) \\ &\quad + \sum_{j=1}^{p_1} FW_j \ln P_c(x_j | y = k) \\ &\quad + \sum_{j=1}^{p_2} FW_j \ln P_t(x_j | y = k) \}. \end{aligned} \quad (15)$$

Denote the estimations of the parameters μ_{kj} , σ_{kj}^2 , θ_{kj} , p_k as $\hat{\mu}_{kj}$, $\hat{\sigma}_{kj}^2$, $\hat{\theta}_{kj}$, \hat{p}_k , respectively. Suppose that these

parameters are estimated based on the historical data and feature weights are calculated. Then we have

$$\begin{aligned} y_{new} &= \arg \max_k \{ \ln \hat{p}_k + \sum_{j=1}^{p_2} FW_j [(\mathbf{x}_{new})_j \ln \hat{\theta}_{kj} \\ &\quad + (1 - (\mathbf{x}_{new})_j) \ln(1 - \hat{\theta}_{kj})] \\ &\quad + \sum_{j=1}^{p_1} FW_j \ln \left[\frac{1}{\sqrt{2\pi\hat{\sigma}_{kj}^2}} \exp\left(-\frac{((\mathbf{x}_{new})_j - \hat{\mu}_{kj})^2}{2\hat{\sigma}_{kj}^2}\right) \right] \}, \end{aligned} \quad (16)$$

where $(\mathbf{x}_{new})_j$ is the j th feature of \mathbf{x}_{new} .

Denoting that

$$\tilde{\mathbf{x}} = [(\mathbf{x}_{new})_1, (\mathbf{x}_{new})_2, \dots, (\mathbf{x}_{new})_{p_2}, 1], \quad (17)$$

$$\begin{aligned} \boldsymbol{\varphi}_k &= [FW_1 \ln \frac{\hat{\theta}_{k1}}{1 - \hat{\theta}_{k1}}, \dots, FW_{p_2} \ln \frac{\hat{\theta}_{kp_2}}{1 - \hat{\theta}_{kp_2}}, \\ &\quad \ln \hat{p}_k + \sum_{j=1}^{p_2} FW_j \ln(1 - \hat{\theta}_{kj})]^T, \end{aligned} \quad (18)$$

$$\phi_k = \sum_{j=1}^{p_1} FW_j \left[-\frac{1}{2} \ln(2\pi\hat{\sigma}_{kj}^2) - \frac{((\mathbf{x}_{new})_j - \hat{\mu}_{kj})^2}{2\hat{\sigma}_{kj}^2} \right], \quad (19)$$

we have

$$\begin{aligned} y_{new} &= \arg \max_k \left\{ \sum_{j=1}^{p_2} (\mathbf{x}_{new})_j FW_j \ln \frac{\hat{\theta}_{kj}}{1 - \hat{\theta}_{kj}} \right. \\ &\quad + \ln \hat{p}_k + \sum_{j=1}^{p_2} FW_j \ln(1 - \hat{\theta}_{kj}) \\ &\quad + \left. \sum_{j=1}^{p_1} FW_j \left(-\frac{1}{2} \ln(2\pi\hat{\sigma}_{kj}^2) - \frac{((\mathbf{x}_{new})_j - \hat{\mu}_{kj})^2}{2\hat{\sigma}_{kj}^2} \right) \right\} \\ &= \arg \max_k (\tilde{\mathbf{x}} \cdot \boldsymbol{\varphi}_k + \phi_k). \end{aligned} \quad (20)$$

B. Parameters Estimation

In this subsection, \mathbf{X} is used for parameter estimation. According to maximum likelihood estimation (MLE) [35], the prior probability can be given as

$$\hat{p}_k = \hat{p}_k^{MLE} = \sum_{i=1}^n \pi_{ik} / n, \quad (22)$$

where $\pi_{ik} = 1\{y_i = k\}$.

For two-valued variables, the response probability is represented by

$$\hat{\theta}_{kj} = \hat{\theta}_{kj}^{MLE} = \left(\sum_{i=1}^n x_{ij} \pi_{ik} \right) / \sum_{i=1}^n \pi_{ik}, \quad (23)$$

where x_{ij} is the i th sample of the x_j . In order to avoid a probability of 0 or 1, the estimated response probability and prior probability are truncated as following [36].

Assuming $\sum_i \pi_{iK} = \max_{1 \leq k \leq K} \sum_i \pi_{iK}$, then for $1 \leq k \leq K - 1$, we have

$$\hat{p}_k = \hat{p}_k^{DTE} = \max\{\xi, \min(\sum_{i=1}^n \pi_{ik} / n, 1 - \xi)\}, \quad (24)$$

where ξ is a small positive value ($\xi = 10^{-6}$ in this article). When $k = K$, we can get \hat{p}_K by

$$\hat{p}_K = \hat{p}_K^{DTE} = 1 - \sum_{k=1}^{K-1} \hat{p}_k. \quad (25)$$

Under the same assumption, we have $\theta_{Kj} = \max_{1 \leq k \leq K} (\sum_{i=1}^n x_{ij} \pi_{ik}) / \sum_{i=1}^n \pi_{ik}$, for $1 \leq k \leq K - 1$. Then the response probability is defined by

$$\hat{\theta}_{kj} = \hat{\theta}_{kj}^{DTE} = \max\{\xi, \min[(\sum_{i=1}^n x_{ij} \pi_{ik}) / \sum_{i=1}^n \pi_{ik}, 1 - \xi]\}. \quad (26)$$

When $k = K$, one has

$$\hat{\theta}_{Kj} = \hat{\theta}_{Kj}^{DTE} = 1 - \sum_{k=1}^{K-1} \hat{\theta}_{kj}^{DTE}, \quad (27)$$

where $\hat{p}_k^{DTE} \in [\xi, 1 - \xi]$, $\hat{\theta}_{kj}^{DTE} \in [\xi, 1 - \xi]$.

For continuous variables, the mean and the standard deviation (std) are respectively estimated as

$$\hat{\mu}_{kj} = \hat{\mu}_{kj}^{MLE} = (\sum_{i=1}^n x_{ij} \pi_{ik}) / \sum_{i=1}^n \pi_{ik}, \quad (28)$$

$$\hat{\sigma}_{kj} = \hat{\sigma}_{kj}^{MLE} = \sqrt{\sum_{i=1}^n [\pi_{ik} (x_{ij} - \hat{\mu}_{kj})]^2 / (\sum_{i=1}^n \pi_{ik} - 1)}. \quad (29)$$

The estimation of feature weights is mainly to estimate the MI (13). If x_j or $x_{j'}$ is a continuous variable, the corresponding auxiliary two-valued variable is used for computing MI. $P(x_j)$ and $P(x_{j'})$ can be calculated in a similar way, and this paper only shows the calculation of $P(x_j)$

$$\begin{aligned} \hat{P}(x_j = 1) &= \sum_{i=1}^n x_{ij} / n, \\ \hat{P}(x_j = 0) &= \sum_{i=1}^n \bar{x}_{ij} / n, \end{aligned} \quad (30)$$

where

$$\bar{x}_{ij} = \begin{cases} 1, & x_{ij} = 0, \\ 0, & x_{ij} = 1. \end{cases} \quad (31)$$

The estimate of $P(x_j, x_{j'})$ is shown in Theorem 1.

Theorem 1. For two-valued variables $x_j : \{x_{1j}, x_{2j}, \dots, x_{nj}\}$ and $x_{j'} : \{x_{1j'}, x_{2j'}, \dots, x_{nj'}\}$, $\hat{P}(x_j, x_{j'})$ can be estimated as

$$\begin{aligned} \hat{P}(x_j = 1, x_{j'} = 1) &= \hat{P}(x_{j'} = 1) \hat{\varphi}_n, \\ \hat{P}(x_j = 0, x_{j'} = 1) &= \hat{P}(x_{j'} = 1) (1 - \hat{\varphi}_n), \\ \hat{P}(x_j = 0, x_{j'} = 0) &= \hat{P}(x_{j'} = 0) \hat{\varphi}'_n, \\ \hat{P}(x_j = 1, x_{j'} = 0) &= \hat{P}(x_{j'} = 0) (1 - \hat{\varphi}'_n), \end{aligned} \quad (32)$$

where $\hat{\varphi}_n = \frac{\sum_{i=1}^n x_{ij} x_{ij'}}{\sum_{i=1}^n x_{ij}}$ and $\hat{\varphi}'_n = \frac{n + \sum_{i=1}^n x_{ij} x_{ij'} - \sum_{i=1}^n (x_{ij} + x_{ij'})}{n - \sum_{i=1}^n x_{ij'}}$.

The proof of Theorem 1 is shown in Appendix B. Compared with MHNBM, the introduction of feature weights and the technology of clipping effectively avoids directly calculating the conditional probability $P(x_j | x_{j'}, y = k)$ between variables. Equation (23), (28) and (29) are calculations for each variable. The feature weights can be obtained through two-valued or auxiliary two-valued variables. In this paper, all the estimates of probability are double-truncated according to the above method if necessary. The above analysis is illustrated by the following Algorithm.

Algorithm: FWMNBM

Offline modeling:

- Step 1. Divide the training data (\mathbf{X}, y) into two-valued variables \mathbf{X}_t and continuous variables \mathbf{X}_c .
 - Step 2. Construct the auxiliary two-valued variable for each continuous variable according to (12).
 - Step 3. Calculate the estimates of each probability via (30) and (32).
 - Step 4. Calculate the mutual information between variables and between variables and labels.
 - Step 5. Calculate the weight of the feature (FW_j).
 - Step 6. Estimate the response functions ($\hat{\theta}_{kj}$) and the prior probabilities (\hat{p}_k) of two-valued variables.
 - Step 7. Estimate the mean ($\hat{\mu}_{kj}$) and the standard deviation ($\hat{\sigma}_{kj}$) of continuous variables.
 - Step 8. Build the model with the estimated parameters.
- ##### Online detection:
- Step 9. Select the sampled data and construct vector $\tilde{\mathbf{x}}$ via (17).
 - Step 10. Calculate φ_k, ϕ_k via (18) and (19).
 - Step 11. Calculate $\tilde{\mathbf{x}} \cdot \varphi_k + \phi_k$ for every k , then $\arg \max_k (\tilde{\mathbf{x}} \cdot \varphi_k + \phi_k)$ is the predicted label.
-

IV. SIMULATION

In this section, the numerical cases of a numerical simulation example and a practical vibration fault case of ZTPP are utilized to validate the effectiveness of FWMNBM.

A. Numerical simulation

The numerical simulation data contains 5 continuous variables and 5 two-valued variables. The means of continuous variables are shown in Table I and corresponding stds are displayed in Table II. The two-valued variable values under different classes are depicted in Table III. In order to make the case more general, the two-valued variable values under different classes are randomly adjusted. The adjustment percentages are listed in Table III. For instance, some values of v_6 under normal working conditions, which are set as 0, are changed to be 1 after adjusting. Under each working condition, 1500 samples are randomly generated according to the parameters. The samples under normal 1 and fault 1 are used for training the model, and the other instances are used for testing.

The Gaussian naive Bayesian model (GNBM) is used for the continuous variables and the Bernoulli naive Bayesian

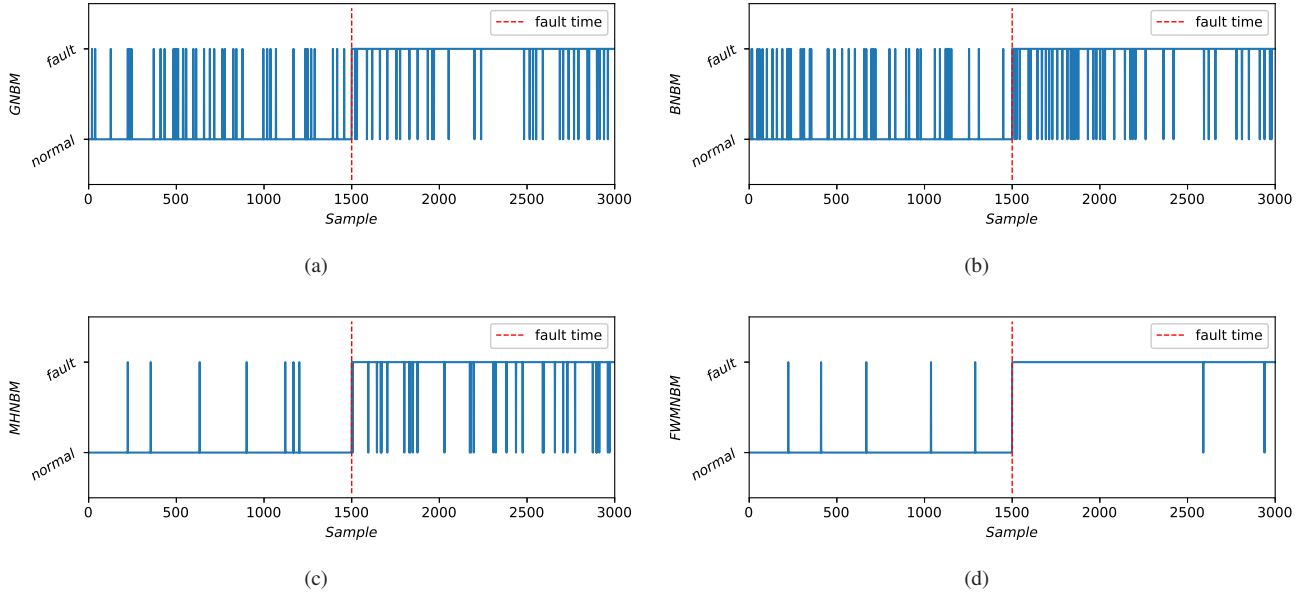


Fig. 1. The label results of different methods. (a)GNBM; (b)BNBM; (c)MHNBM; (d)FWMNBM.

TABLE I
THE PRESET MEANS OF CONTINUOUS VARIABLES

variables	means			
	normal 1	fault 1	normal 2	fault 2
v_1	0.00	3.50	0.32	3.25
v_2	0.00	4.50	0.28	4.40
v_3	0.00	3.20	0.24	3.12
v_4	0.00	2.20	0.01	2.15
v_5	0.00	0.80	0.00	0.80

TABLE II
THE PRESET STDS OF CONTINUOUS VARIABLES

variables	stds			
	normal 1	fault 1	normal 2	fault 2
v_1	1.50	1.00	1.49	1.25
v_2	1.60	2.50	1.56	2.55
v_3	0.80	1.70	0.88	1.73
v_4	2.00	1.80	2.00	1.75
v_5	1.40	2.70	1.40	2.70

TABLE III
THE VALUES AND ADJUSTMENT PERCENTAGE OF TWO-VALUED VARIABLES

variables	normal		fault	
	values	percentage	values	percentage
v_6	0	30%	1	30%
v_7	1	25%	0	25%
v_8	0	20%	1	20%
v_9	0	15%	1	15%
v_{10}	1	10%	0	10%

model (BNBM) is utilized to two-valued variables. That is only v_1, \dots, v_5 are used for build and test GNBM, and BNBM just utilize the information of v_6, \dots, v_{10} for modeling and verification. Different from GNBM and BNBM, MHNBM and FWHNBM are utilized for modeling and anomaly detection with both two-valued and continuous variables. The first 1500 samples of test data are normal data, and the rest are marked as faults. The test results of all above models for the testing

data are depicted in Fig. 1(a)-Fig. 1(d). There are a lot of false alarms and missing faults when only continuous or two-valued variables are used, which can be seen in Fig. 1(a) and Fig. 1(b). MHNBM and FWHNBM have better performance because they can simultaneously mine continuous and two-valued information at the same time. Compared to MHNBM, FWHNBM has the lower false alarm rate (FARs) and a higher fault detection rate (FDR) which are depicted in Fig. 2(e) and Fig. 2(f).

B. Actual data validation

A vibration fault of ZTPP is also used to illustrate the effectiveness of FWMNBM. At 11:35 on September 3, 2017, a hydraulic cylinder vibration fault of the primary air fan occurred, and it was recovered after 26 hours. The data, containing 495 two-valued variables and 260 continuous variables, is sampled every 5 seconds and collected from 11:35, September 1, 2017. A total of 53280 instances are collected for model training and testing.

The first 60% instances under normal conditions and first 60% fault instances are utilized for modeling, and the remaining data is used for testing. In this article, we used 35 two-valued variables and continuous variables respectively. The detailed variable selection process can refer to article [28]. In the traditional methods, LDA [24], decision trees (DT) [37], SVM [26], k-nearest neighbors (KNN) [25] are adopted to detect anomaly with the continuous variables. MHNBM and FWMNBM are used with both two-valued and continuous variables. The testing results of all methods are shown in Fig. 2(a)-Fig. 2(f).

Excepting for DT, the performance of other methods in terms of FDR are very satisfactory. DT has omission of fault and all methods have false alarms. In order to compare the performance of various methods, the FDRs and FARs of all methods are shown in Table IV. From the experimental results,

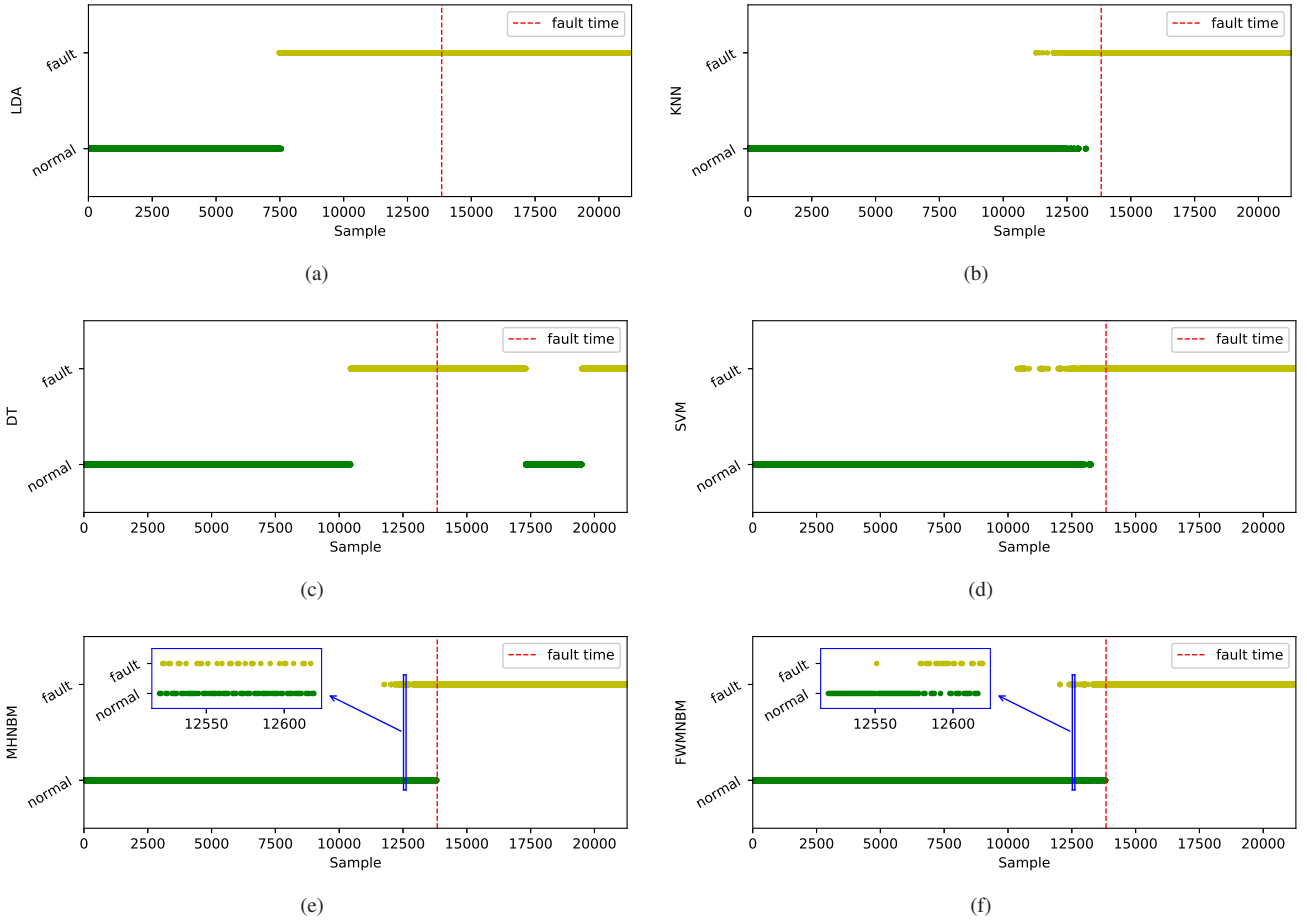


Fig. 2. The testing results. (a)LDA; (b)KNN; (c)DT; (d)SVM; (e)MHNBM; (f)FWMNBM.

TABLE IV
FARS AND FDRs OF ALL METHODS

Number	Methods	FARs(%)	FDRs(%)
1	LDA	45.65	100
2	KNN	7.21	100
3	DT	24.54	70.39
4	SVM	8.82	100
5	MHNBM	3.10	100
6	FWMNBM	2.45	100

the addition of two-valued variables can reduce the impact of parameter fluctuations before a fault occurs. Affected by anomaly evolution, some normal instances are misclassified into fault. However, MHNBM and FWMNBM effectively reduce FAR through combining multiple process data sources, because the advantages of both two-valued and continuous variables are taken into consideration. Among all methods, FWMNBM has the best detection performance.

V. CONCLUSION

A data-driven anomaly detection method called FWMNBM is proposed with both two-valued and continuous variables. For FWMNBM, the variables more correlated to class have greater weights, which makes the more discriminating variables contribute more to the model. At the same time,

FWMNBM can effectively avoid calculating the conditional probability between variables so that it can still be used when the amount of training data is not sufficient. In addition, a more effective consistent characterization method for the correlation of mixed variables is provided, and the corresponding feasibility analysis is conducted. The superior performance of FWMNBM is verified by the numerical cases of a numerical simulation example and an actual plant's case. Compared to traditional classical approaches, MHNBM and FWMNBM significantly improve the anomaly monitoring performance by increasing the information of two-valued variables. Furthermore, FWMNBM has more outstanding performance because greater weights are assigned to variables with greater difference under different working conditions.

APPENDIX A ANALYSIS OF DEFINITION 1

Definition 1 unifies the correlation analysis between variables containing both two-valued and continuous variables by the same standard. The correlation between two-valued variables and two-valued variables or between continuous variables and continuous variables can be effectively characterized, and the original two-valued variables do not change. Therefore, the rationality of Definition 1 can be proved when a quantitative relationship exists between the correlation index

of the auxiliary two-valued variables and that of original continuous variables.

Let $x_j \in X_c$ and $x_{j'} \in X_c$, where x'_j and $x'_{j'}$ are the corresponding auxiliary two-valued variables. For simplicity, we assume that x_j and $x_{j'}$ are standard normal distributions through standardization, that is, $E(x_j) = E(x_{j'}) = 0$, $D(x_j) = D(x_{j'}) = 1$.

The joint probability density function of x_j and $x_{j'}$ is

$$\begin{aligned} f(x_j, x_{j'}) &= \frac{1}{2\pi\sqrt{1 - [\rho(x_j, x_{j'})]^2}} \times \exp\left\{-\frac{1}{2(1 - [\rho(x_j, x_{j'})]^2)}\right. \\ &\quad \left. \times [(x_j)^2 + (x_{j'})^2 - 2\rho(x_j, x_{j'})x_jx_{j'}]\right\}. \end{aligned} \quad (33)$$

Since

$$\begin{aligned} E(x'_jx'_{j'}) &= P(x'_j = 1, x'_{j'} = 1) = P(x_j > 0, x_{j'} > 0) \\ &= \frac{1}{2\pi\sqrt{1 - [\rho(x_j, x_{j'})]^2}} \int_0^\infty \int_0^\infty \exp\left\{-\frac{1}{2(1 - [\rho(x_j, x_{j'})]^2)}\right. \\ &\quad \left. \times [(x_j)^2 + (x_{j'})^2 - 2\rho(x_j, x_{j'})x_jx_{j'}]\right\} dx_j dx_{j'} \\ &= \frac{1}{2\pi\sqrt{1 - [\rho(x_j, x_{j'})]^2}} \int_0^\infty \int_0^{\frac{\pi}{2}} r \exp\left\{-\frac{1}{2(1 - [\rho(x_j, x_{j'})]^2)}\right. \\ &\quad \left. \times [1 - \rho(x_j, x_{j'}) \sin 2\theta]\right\} d\theta dr \\ &= \frac{1}{2\pi\sqrt{1 - [\rho(x_j, x_{j'})]^2}} \int_0^{\frac{\pi}{2}} \frac{1}{1 - \rho(x_j, x_{j'}) \sin 2\theta} d\theta \\ &= \frac{1}{2\pi\sqrt{1 - [\rho(x_j, x_{j'})]^2}} \\ &\quad \times \int_0^{\frac{\pi}{2}} \frac{1}{\tan^2\theta - 2\rho(x_j, x_{j'}) \tan\theta + 1} d \tan\theta \\ &= \frac{1}{2\pi} \left\{ \frac{\pi}{2} + \arctan[\rho(x_j, x_{j'})/\sqrt{1 - [\rho(x_j, x_{j'})]^2}] \right\}, \end{aligned} \quad (34)$$

then

$$E(x'_jx'_{j'}) = \frac{1}{4} + \frac{1}{2\pi} \arcsin \rho(x_j, x_{j'}). \quad (35)$$

Let $P(x'_j = 1|x'_{j'} = 1) = \varphi$.

Note that

$$\begin{aligned} P(x'_j = 1) &= P(x_j > 0) = \frac{1}{2}, \\ P(x'_j = 0) &= P(x_j \leq 0) = \frac{1}{2}, \\ E(x'_j) &= 1 \times P(x'_j = 1) + 0 \times P(x'_j = 0) = \frac{1}{2}, \\ D(x'_j) &= E(x'_j - \frac{1}{2})^2 = E[(x'_j)^2] - [E(x'_j)]^2 = \frac{1}{4} \text{ and} \\ P(x'_{j'} = 1) &= P(x'_{j'} = 0) = \frac{1}{2}, \\ E(x'_{j'}) &= \frac{1}{2}, D(x'_{j'}) = \frac{1}{4}. \end{aligned}$$

Then we have

$$\begin{aligned} E(x'_jx'_{j'}) &= P(x'_j = 1, x'_{j'} = 1) \\ &= P(x'_j = 1)P(x'_j = 1|x'_{j'} = 1) = \frac{1}{2}\varphi, \end{aligned} \quad (36)$$

$$\begin{aligned} Cov(x'_j, x'_{j'}) &= E[(x'_j - \frac{1}{2})(x'_{j'} - \frac{1}{2})] \\ &= E(x'_jx'_{j'}) - \frac{1}{4} = \frac{1}{2}\varphi - \frac{1}{4}, \end{aligned} \quad (37)$$

$$\begin{aligned} \rho(x'_j, x'_{j'}) &= \frac{Cov(x'_j, x'_{j'})}{\sqrt{D(x'_j)D(x'_{j'})}} \\ &= 2\varphi - 1 = 4E(x'_jx'_{j'}) - 1 \\ &= 4 \times \left[\frac{1}{4} + \frac{1}{2\pi} \arcsin \rho(x_j, x_{j'}) \right] - 1, \end{aligned} \quad (38)$$

and then

$$\rho(x'_j, x'_{j'}) = \frac{2}{\pi} \arcsin \rho(x_j, x_{j'}). \quad (39)$$

APPENDIX B

PROOF OF THEOREM 1

Let $P(x_j = 1|x_{j'} = 1) = \varphi$, $\hat{P}(x_j = 1|x_{j'} = 1) = \hat{\varphi}_n$ and $P(x_j = 0|x_{j'} = 0) = \varphi'$, $\hat{P}(x_j = 0|x_{j'} = 0) = \hat{\varphi}'_n$.

Then

$$\begin{aligned} \hat{P}(x_j = 0|x_{j'} = 1) &= (1 - \hat{\varphi}_n), \\ \hat{P}(x_j = 1|x_{j'} = 0) &= (1 - \hat{\varphi}'_n). \end{aligned} \quad (40)$$

Since $P(x_j, x_{j'}) = P(x_{j'})P(x_j|x_{j'})$. It can be learned that

$$\begin{aligned} \hat{P}(x_j = 1, x_{j'} = 1) &= \hat{P}(x_{j'} = 1)\hat{\varphi}_n, \\ \hat{P}(x_j = 0, x_{j'} = 1) &= \hat{P}(x_{j'} = 1)(1 - \hat{\varphi}_n), \\ \hat{P}(x_j = 0, x_{j'} = 0) &= \hat{P}(x_{j'} = 0)\hat{\varphi}'_n, \\ \hat{P}(x_j = 1, x_{j'} = 0) &= \hat{P}(x_{j'} = 0)(1 - \hat{\varphi}'_n). \end{aligned} \quad (41)$$

In addition, we have [38]

$$\begin{aligned} &P(x_j = x_{1j}, x_{j'} = x_{1j'})P(x_j = x_{2j}, x_{j'} = x_{2j'}) \dots \\ &P(x_j = x_{ij}, x_{j'} = x_{ij'}) \dots P(x_j = x_{nj}, x_{j'} = x_{nj'}) \\ &= P(x_{j'} = x_{1j'})P(x_j = x_{1j}|x_{j'} = x_{1j'})P(x_{j'} = x_{2j'}) \\ &P(x_j = x_{2j}|x_{j'} = x_{2j'}) \dots P(x_{j'} = x_{ij'}) \\ &\times P(x_j = x_{ij}|x_{j'} = x_{ij'}) \dots P(x_{j'} = x_{nj'}) \\ &\times P(x_j = x_{nj}|x_{j'} = x_{nj'}) \\ &= \prod_{i=1}^n P(x_{j'} = x_{ij'}) \prod_{i=1}^n P(x_j = x_{ij}|x_{j'} = x_{ij'}). \end{aligned} \quad (42)$$

Then it can be obtained that

$$\begin{aligned} P(x_j = x_{ij}|x_{j'} = x_{ij'}) &= \varphi^{x_{ij}x_{ij'}} (1 - \varphi)^{x_{ij'} - x_{ij}x_{ij'}} \\ &\times \varphi'^{1+x_{ij}x_{ij'} - x_{ij} - x_{ij'}} (1 - \varphi')^{x_{ij} - x_{ij}x_{ij'}}. \end{aligned} \quad (43)$$

The likelihood function is

$$\begin{aligned} \ell(\varphi, \varphi') &= \prod_{i=1}^n P(x_{j'} = x_{ij'}) \prod_{i=1}^n P(x_j = x_{ij}|x_{j'} = x_{ij'}) \\ &= \varphi^{\sum_{i=1}^n x_{ij}x_{ij'}} (1 - \varphi)^{\sum_{i=1}^n x_{ij'} - \sum_{i=1}^n x_{ij}x_{ij'}} \varphi'^{n + \sum_{i=1}^n x_{ij}x_{ij'}} \\ &\quad \times \varphi'^{-\sum_{i=1}^n (x_{ij} + x_{ij'})} (1 - \varphi')^{\sum_{i=1}^n x_{ij} - \sum_{i=1}^n x_{ij}x_{ij'}}, \end{aligned} \quad (44)$$

where ϖ is a constant. The derivative $\frac{\partial \ell(\varphi, \varphi')}{\partial \varphi}$ of $\ell(\varphi, \varphi')$ with respect to φ is

$$\begin{aligned} & \frac{\partial \ell(\varphi, \varphi')}{\partial \varphi} \\ &= [\varpi \varphi^{n + \sum_{i=1}^n x_{ij} x_{ij'}} - \sum_{i=1}^n (x_{ij} + x_{ij'}) (1 - \varphi')^{\sum_{i=1}^n x_{ij} - \sum_{i=1}^n x_{ij} x_{ij'}}] \\ & \times \sum_{i=1}^n x_{ij} x_{ij'} \varphi^{\sum_{i=1}^n x_{ij} x_{ij'}} \varphi^{-1} (1 - \varphi)^{\sum_{i=1}^n x_{ij'} - \sum_{i=1}^n x_{ij} x_{ij'}} \\ & + [\varpi \varphi^{n + \sum_{i=1}^n x_{ij} x_{ij'}} - \sum_{i=1}^n (x_{ij} + x_{ij'}) (1 - \varphi')^{\sum_{i=1}^n x_{ij} - \sum_{i=1}^n x_{ij} x_{ij'}}] \\ & \times \varphi^{\sum_{i=1}^n x_{ij} x_{ij'}} [\sum_{i=1}^n x_{ij'} - \sum_{i=1}^n x_{ij} x_{ij'}] (1 - \varphi)^{\sum_{i=1}^n x_{ij'} - \sum_{i=1}^n x_{ij} x_{ij'}} \\ & \times (1 - \varphi)^{-1} (-1). \end{aligned} \quad (45)$$

Let $\frac{\partial \ell(\varphi, \varphi')}{\partial \varphi} = 0$, then

$$\hat{\varphi}_n = \frac{\sum_{i=1}^n x_{ij} x_{ij'}}{\sum_{i=1}^n x_{ij'}}. \quad (46)$$

The derivative $\frac{\partial \ell(\varphi, \varphi')}{\partial \varphi'}$ of $\ell(\varphi, \varphi')$ with respect to φ' is

$$\begin{aligned} & \frac{\partial \ell(\varphi, \varphi')}{\partial \varphi'} = [\varpi \varphi^{\sum_{i=1}^n x_{ij} x_{ij'}} (1 - \varphi)^{\sum_{i=1}^n x_{ij'} - \sum_{i=1}^n x_{ij} x_{ij'}}] \\ & \times [n + \sum_{i=1}^n x_{ij} x_{ij'} - \sum_{i=1}^n (x_{ij} + x_{ij'})] \\ & \times \varphi^{n + \sum_{i=1}^n x_{ij} x_{ij'}} - \sum_{i=1}^n (x_{ij} + x_{ij'}) \varphi'^{-1} (1 - \varphi')^{\sum_{i=1}^n x_{ij} - \sum_{i=1}^n x_{ij} x_{ij'}} \\ & + [\varpi \varphi^{\sum_{i=1}^n x_{ij} x_{ij'}} (1 - \varphi)^{\sum_{i=1}^n x_{ij'} - \sum_{i=1}^n x_{ij} x_{ij'}}] \\ & \times \varphi^{n + \sum_{i=1}^n x_{ij} x_{ij'}} - \sum_{i=1}^n (x_{ij} + x_{ij'}) (\sum_{i=1}^n x_{ij} - \sum_{i=1}^n x_{ij} x_{ij'}) \\ & \times (1 - \varphi')^{\sum_{i=1}^n x_{ij} - \sum_{i=1}^n x_{ij} x_{ij'}} (1 - \varphi')^{-1} (-1). \end{aligned} \quad (47)$$

Applying $\frac{\partial \ell(\varphi, \varphi')}{\partial \varphi'} = 0$, it can be learned that

$$\hat{\varphi}'_n = \frac{n + \sum_{i=1}^n x_{ij} x_{ij'} - \sum_{i=1}^n (x_{ij} + x_{ij'})}{n - \sum_{i=1}^n x_{ij'}}. \quad (48)$$

REFERENCES

- [1] Y. Yang, X. Shi, X. Liu, and H. Li, "A novel MDFA-MKECA method with application to industrial batch process monitoring," *IEEE/CAA J. Autom. Sinica*, vol. 7, no. 5, pp. 1446–1454, 2020.
- [2] Q. Zhu, "Latent variable regression for supervised modeling and monitoring," *IEEE/CAA J. Autom. Sinica*, vol. 7, no. 3, pp. 800–811, 2020.
- [3] A. Imtiaz, D. Aldo, and D. Yu, "Unsupervised anomaly detection based on minimum spanning tree approximated distance measures and its application to hydropower turbines," *IEEE Trans. Autom. Sci. Eng.*, vol. 16, no. 2, pp. 654–667, 2018.
- [4] D. Wu, L. Sheng, D. Zhou, and M. Chen, "Dynamic stationary subspace analysis for monitoring nonstationary dynamic processes," *Ind. Eng. Chem. Res.*, vol. 59, no. 47, pp. 20 787–20 797, 2020.
- [5] H. Ji, H. Xiao, J. Shang, and D. Zhou, "Incipient fault detection with smoothing techniques in statistical process monitoring," *Control Eng. Pract.*, vol. 62, pp. 11–21, 2017.
- [6] K. Zhong, M. Han, and B. Han, "Data-driven based fault prognosis for industrial systems: A concise overview," *IEEE/CAA J. Autom. Sinica*, vol. 7, no. 02, pp. 19–34, 2020.
- [7] G. Anagnostou, F. Boem, and S. Kuenzel, "Observer-based anomaly detection of synchronous generators for power systems monitoring," *IEEE Trans. Power Syst.*, vol. 33, no. 4, pp. 4228–4237, 2018.
- [8] X. Zhang and Z. Ge, "Automatic deep extraction of robust dynamic features for industrial big data modeling and soft sensor application," *IEEE Trans. Ind. Informat.*, vol. 16, no. 7, pp. 4456–4467, 2020.
- [9] H. Ji, K. Huang, and D. Zhou, "Incipient sensor fault isolation based on augmented mahalanobis distance," *Control Eng. Pract.*, vol. 86, pp. 144–154, 2019.
- [10] J. Shi, J. Sun, Y. Yang, and D. Zhou, "Distributed self-triggered formation control for multi-agent systems," *Sci. China Inf. Sci.*, vol. 63, no. 10, p. 209207, 2020.
- [11] Y. Wang, D. Zhou, M. Chen, and M. Wang, "Weighted part mutual information related component analysis for quality-related process monitoring," *J. Process Control*, vol. 88, pp. 111–123, 2020.
- [12] P. Y. Zhang, S. Shu, and M. C. Zhou, "An online fault detection model and strategies based on SVM-grid in clouds," *IEEE/CAA J. Autom. Sinica*, vol. 5, no. 2, pp. 445–456, 2018.
- [13] S. Yin, H. Gao, J. Qiu, and O. Kaynak, "Fault detection for nonlinear process with deterministic disturbances: A just-in-time learning based data driven method," *IEEE Trans. Cybern.*, vol. 47, no. 11, pp. 3649–3657, 2017.
- [14] Z. Ge, Z. Song, and F. Gao, "Review of recent research on data-based process monitoring," *Ind. Eng. Chem. Res.*, vol. 52, no. 10, pp. 3543–3562, 2013.
- [15] L. Yao and Z. Ge, "Scalable semisupervised GMM for big data quality prediction in multimode processes," *IEEE Trans. Ind. Electron.*, vol. 66, no. 5, pp. 3681–3692, 2019.
- [16] V. Agrawal, B. K. Panigrahi, and P. M. V. Subbarao, "Review of control and fault diagnosis methods applied to coal mills," *J. Process Control*, vol. 32, pp. 138–153, 2015.
- [17] Z. Gao, C. Cecati, and S. X. Ding, "A survey of fault diagnosis and fault-tolerant techniques—part I: Fault diagnosis with model-based and signal-based approaches," *IEEE Trans. Ind. Electron.*, vol. 62, no. 6, pp. 3757–3767, 2015.
- [18] K. Huang, H. Wen, H. Ji, L. Cen, X. Chen, and C. Yang, "Nonlinear process monitoring using kernel dictionary learning with application to aluminum electrolysis process," *Control Eng. Pract.*, vol. 89, pp. 94–102, 2019.
- [19] M. Wang, L. Sheng, D. Zhou, and M. Chen, "A feature weighted mixed naive Bayes model for monitoring anomalies in the fan system of a thermal power plant," *arXiv:2012.07230*, 2020.
- [20] R. Dunia, S. J. Qin, T. F. Edgar, and T. J. McAvoy, "Identification of faulty sensors using principal component analysis," *AIChE J.*, vol. 42, no. 10, pp. 2797–2812, 1996.
- [21] W. Li, H. H. Yue, S. Valle-Cervantes, and S. J. Qin, "Recursive PCA for adaptive process monitoring," *J. Process Control*, vol. 10, no. 5, pp. 471–486, 2000.
- [22] S. J. Qin, "Recursive PLS algorithms for adaptive data modeling," *Comput. Chem. Eng.*, vol. 22, no. 4-5, pp. 503–514, 1998.
- [23] N. Sheng, Q. Liu, S. J. Qin, and T. Chai, "Comprehensive monitoring of nonlinear processes based on concurrent kernel projection to latent structures," *IEEE Trans. Autom. Sci. Eng.*, vol. 13, no. 2, pp. 1129–1137, 2016.
- [24] Q. P. He, S. J. Qin, and W. Jin, "A new fault diagnosis method using fault directions in Fisher discriminant analysis," *AIChE J.*, vol. 51, no. 2, pp. 555–571, 2005.
- [25] Q. He and J. Wang, "Fault detection using the k-nearest neighbor rule for semiconductor manufacturing processes," *IEEE Trans. Semicond. Manuf.*, vol. 20, no. 4, pp. 345–354, 2007.
- [26] A. Widodo and B.-S. Yang, "Support vector machine in machine condition monitoring and fault diagnosis," *Mech. Syst. Signal Pr.*, vol. 21, no. 6, pp. 2560–2574, 2007.
- [27] Z. Ge, Z. Song, S. X. Ding, and B. Huang, "Data mining and analytics in the process industry: The role of machine learning," *IEEE Access*, vol. 5, pp. 20 590–20 616, 2017.
- [28] M. Wang, D. Zhou, M. Chen, and Y. Wang, "Anomaly detection in the fan system of a thermal power plant monitored by continuous and two-valued variables," *Control Eng. Pract.*, vol. 102, p. 104522, 2020.
- [29] M. Ozuysal, M. Calonder, V. Lepetit, and P. Fua, "Fast keypoint recognition using random ferns," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 3, pp. 448–461, 2010.
- [30] E. J. D. Fortuny, D. Martens, and F. Provost, "Wallenius Bayes," *Mach. Learn.*, vol. 107, no. 2, pp. 1–25, 2018.

- [31] L. Zhang, L. Jiang, C. Li, and G. Kong, "Two feature weighting approaches for naive Bayes text classifiers," *Knowl. Based Syst.*, vol. 100, no. may 15, pp. 137–144, 2016.
- [32] L. Jiang, L. Zhang, C. Li, and J. Wu, "A correlation-based feature weighting filter for naive Bayes," *IEEE Trans. Knowl. Data Eng.*, vol. 31, no. 2, pp. 201–213, 2019.
- [33] C. Ratanamahatana, E. Keogh, A. J. Bagnall, and S. Lonardi, *A novel bit level time series representation with implication of similarity search and clustering*. Springer Berlin Heidelberg, 2005.
- [34] Y. Wu and L. Liu, "Conditional entropy and mutual information in random cascading processes," *Physical Review D*, vol. 43, no. 9, pp. 3077–3079, 1991.
- [35] M. Collins, *Parameter estimation for statistical parsing models: Theory and practice of distribution-free methods*. Springer Netherlands, 2004.
- [36] N. Balakrishnan and V. B. Nevzorov, *Bernoulli Distribution*. U.K.: Springer, 2008.
- [37] Y. Sheng and R. Steven, "Decision tree-based methodology for high impedance fault detection," *IEEE Trans. Power Del.*, vol. 19, no. 2, pp. 533–536, 2004.
- [38] T. Zhang, D. Yue, Y. Gu, Y. Wang, and G. Yu, "Adaptive correlation analysis in stream time series with sliding windows," *Comput. Math. Appl.*, vol. 57, no. 6, pp. 937–948, 2009.