# Fault Warning Technology Based on Multivariate Statistical Analysis

**Hao Hu, Fuzhou Feng[*], Junfeng Han, Junzhen Zhu, Chao Song**
*Department of Vehicle Engineering, Army Academy of Armored Forces, Beijing, China*
*\*Corresponding Author.*

**Abstract: Traditional data-based modeling methods require obtaining many fault samples for fault warning. However, during the operation of the equipment, there is a large amount of normal operating state data collected while there is a small amount of failure sample data. Therefore, the reliability of failure samples is not high during the operation of the equipment. Multiple statistical analysis and state monitoring techniques, such as principal component analysis, can construct fault warning models with only data under normal working conditions. This article combines kernel transformation with principal component analysis to construct a kernel principal component analysis method suitable for small amounts of fault data conditions, thereby achieving effective early warning of equipment operation abnormalities and faults. This article proposes a fault identification method based on multivariate contribution rate graph to address the difficulty of identifying fault sources, achieving precise identification and localization of fault sources under abnormal working conditions. The research results of this article can lay the foundation for establishing an early warning model for equipment failures.**

**Keywords: Multivariate Statistical Analysis; Fault Warning; Principal Component Analysis; Data; Feature**

## 1. Introduction

The working status of a certain type of vehicle is quite complex, and early warning of its faults is of great significance for the health management, accident prevention, and normal use and maintenance of the vehicle. This article proposes a fault warning method based on PCA and applies it to vehicle fault warning. The data obtained from vehicles under different health conditions have different characteristics [1]. Zhang et al. used dynamic principal component analysis to obtain the time-varying characteristics of the system and used dissimilarity indicators to monitor the residual value of the system, thereby achieving the goal of fault detection and diagnosis of the system [2]. Han et al. used a combination of principal component analysis and multivariate state estimation to achieve effective prediction of induced draft fan faults [3]. Xu et al. applied the variable scale PCA method and K-means method to provide abnormal warning for load currents in the power grid and achieve fault localization [4]. Wu et al. used the KPCA method to monitor the relevant data of the precision rolling equipment and identified the main cause of the failure by drawing the contribution diagram of each parameter [5]. Guo et al. used principal component analysis to monitor chemical processes and verifies the correctness of the semiconductor generation process [6]. Yuan and Sun proposed a fault diagnosis method that combines nearest neighbor normalization of local information and principal component analysis, and successfully achieved monitoring of multimodal processes using $T^2$ and $SPE$ statistics [7]. Yao et al. pointed out that there are problems such as data anomalies and time delays in data-driven fault diagnosis technology, and provided research methods [8]. Wang et al. proposed a data-driven fault detection method for early warning of rolling bearing faults [9]. Huang et al. proposed a new method for precise separation of gearbox faults, and based on this, achieved real-time monitoring of the status of rotating machinery systems [10]. Yang et al. proposed using the integrated envelope spectrum peak factor to identify early faults in rolling bearings, and verified this method through experiments [11]. Zhang et al. used statistics $T^2$ and $SPE$ to

monitor the real-time operation status of the equipment [12]. Liu et al. proposed a multivariate statistical process monitoring method, which can monitor abnormal situations such as process disturbances and equipment failures in real time, while ensuring personnel safety and improving the efficiency and quality of the process [13]. Scholars have also applied the dynamic principal component analysis algorithm to chemical processes, extracting dynamic change information and achieving good results [14]. Zhang et al. utilized the principle of PCA to establish a new PCA based model for predicting roadbed settlement, and verified it in practical engineering [15]. This article uses PCA algorithm to construct a fault warning model, and based on this, proposes a fault warning model based on KPCA algorithm, thereby constructing a reasonable and efficient fault warning mechanism for equipment.

## 2. The Principle of the Model

### 2.1 PCA Data Dimensionality Reduction

PCA is a multivariate statistical analysis method proposed by Carl Pearson. Due to its simple algorithm and low parameter requirements, it is widely used in pattern recognition, fault diagnosis, and state monitoring. Its essence is to perform rotation and translation transformations on the original coordinate system, so that the origin of the new coordinate system is consistent with the center of gravity of the original data points. The first axis of the new coordinate system corresponds to the direction with the greatest change in the original data, the second axis of the new coordinate system corresponds to the direction with the second largest change in the original data, and the second axis is orthogonal to the first axis, and so on. Remove the axes that contain less information, leaving only $K(K < M)$ main axes. Then, $K$ main axes $(K < M)$ can well describe the changes in the original M-dimensional data, and the new K-dimensional space is the K-dimensional main hyperplane. It is possible to map information from the M dimension to the $K$ dimension, thereby achieving data dimensionality reduction.

The number of points in the collected data is m and the parameter variable is n, then the constructed sample set matrix $X_{m \times n}$ is:

$$X_{m \times n} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{bmatrix}$$

Each row in the matrix represents the size of each point at a certain time, and each column in the matrix represents the value at a certain point at different times.

The first step is to perform Z-score standardization on the dataset matrix $X_{m \times n}$.

$$z_{ij} = x^* = \frac{x_{ij} - \overline{x_j}}{s_j}, i = 1,2,\cdots,m; j = 1,2,\cdots n \quad (1)$$

In the formula: $\overline{x_j} = \frac{1}{m}\sum_{i=1}^{m} x_{ij}$,

$s_j = \sqrt{\frac{1}{m-1}\sum_{i=1}^{m}(x_{ij} - \overline{x_j})^2}$.

The standardized dataset matrix $Z$ is as follows:

$$Z = \begin{bmatrix} z_{11} & z_{12} & \cdots & z_{1n} \\ z_{21} & z_{22} & \cdots & z_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ z_{m1} & z_{m2} & \cdots & z_{mn} \end{bmatrix}$$

Step 2, calculate the correlation coefficient matrix $R$ of $Z$.

$$r_{kl} = \frac{s_{kl}}{s_k s_l} = \frac{\sum_{k=1}^{n}(z_{kl} - \overline{z_k})(z_{kl} - \overline{z_l})}{\sqrt{\sum_{k=1}^{n}(z_{kl} - \overline{z_k})\sum_{k=1}^{n}(z_{kl} - \overline{z_l})}} \quad (2)$$

In the formula: $\overline{z_k} = \frac{1}{m}\sum_{i=1}^{m} x_{ik}, \overline{z_l} = \frac{1}{m}\sum_{i=1}^{m} x_{il}$, $k$ and $l$ are the $k$-th and $l$-th columns of matrix $Z$, respectively, and $r_{kl} = r_{lk}, k = 1,2\cdots n$.

The obtained correlation coefficient matrix $R$ is:

$$R = \begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ r_{21} & r_{22} & \cdots & r_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ r_{n1} & r_{n2} & \cdots & r_{nn} \end{bmatrix}$$

Step 3, calculate the eigenvalues of the correlation coefficient matrix $\lambda$ and the eigenvector $v$. The characteristic equation of

matrix $R$ is $|\lambda I - R| = 0$. Perform orthogonal similarity transformation on matrix $R$. When the elements on the non diagonal of matrix $R$ are equal to or close to zero, the elements on the diagonal are the eigenvalues of matrix R, and the eigenvectors of matrix R are equal to the product of orthogonal similarity transformation matrices.

Step 4, the contribution rate $c_j$ of each monitoring parameter variable is calculated, and $\lambda_j$ Sorted to satisfy $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_n$.

$$c_j = \lambda_j / \sum_{j=1}^{n} \lambda_j \qquad (3)$$

In the formula: $\lambda_j$ is the eigenvalue of the correlation coefficient matrix.

Step 5, the cumulative contribution rate $C_k$ is calculated to determine the main monitoring parameters:

$$C_k = \sum_{j=1}^{k} \lambda_j / \sum_{j=1}^{n} \lambda_j \qquad (4)$$

In the formula: $k(k < n)$ is the first k principal elements.

In engineering, if $C_k \geq 85\%$, the first $k$ principal components can be regarded as maintaining the main information in the data as the main state monitoring parameters, thus achieving the goal of data dimensionality reduction.

## 2.2 Fault Warning Based on KPCA

The KPCA method can transform data from nonlinear space to linear space in high-dimensional feature space, which can process data more accurately and effectively solve the problem of model generalization. KPCA is a nonlinear extension based on PCA, which maps the input space to a high-dimensional feature space through nonlinear transformation, and performs principal component analysis on it to achieve efficient data processing. Kernel principal component analysis has great advantages in extracting nonlinear features and separating data, so it can be used for nonlinear analysis between parameters.

Assuming the sample set contains n samples $\{x_1, x_2, \cdots, x_n\} \in R^m$ (m is the quantity of components), by using a nonlinear function $\Phi(\bullet)$ to perform a nonlinear mapping,

mapping the original lower dimensions to higher dimensions.

The covariance matrix of $\Phi(x_i)$ is:

$$C = \frac{1}{n} \sum_{k=1}^{n} \Phi(x_i) \Phi^T(x_i) \qquad (5)$$

The kernel matrix is:

$$K_{ij} = K(x_i, x_j) \qquad (6)$$

In the mapping space, the $k$-th principal component of the sample is:

$$t_k = \sum_{i=1}^{n} a_i^k K(x, x_i) \qquad (7)$$

$$\alpha_i = \lambda_i^{-1/2} / h_i \qquad (8)$$

In the formula, $\alpha_i$ is the intermediate variable, $h_i$ is the $i$-th eigenvalue of the kernel matrix, and $\lambda_i$ is the $i$-th eigenvalue of the kernel matrix.

In order to reconstruct $\Phi(x)$ through $P$ principal components $t_k$ in the feature space, linear principal component analysis can be performed on it in the mapping space to obtain the reconstruction value:

$$\hat{\Phi}(x) = \sum_{k=1}^{p} t_k V^k \qquad (9)$$

Among them, $V^k$ is an eigenvector of the covariance matrix.

The Gaussian kernel function is as follows:

$$K(x, y) = \exp(-\|x - y\|^2 / c) \qquad (10)$$

In the formula, $x$ and $y$ are data samples; $c$ is a constant greater than 0.

The iterative formula for the original input data space $z$ is:

$$z_{i+1} = \sum_{i=1}^{N} \beta_i \exp(-\frac{\|z_i - x_i\|^2}{c}) x_i / \sum_{i=1}^{N} \beta_i \exp(-\frac{\|z_i - x_i\|^2}{c}) \qquad (11)$$

$\beta_i$ is an intermediate variable, and its formula is:

$$\beta_i = \sum_{k=1}^{p} t_k a_i^k \qquad (12)$$

Iterate the initial value $z_1 = x$ to obtain the vector $z$, which approximates the original input data space.

## 2.3 Multivariate Statistics of the Model

This model only requires monitoring the variables obtained under normal working conditions, utilizing the interrelationships between process variables and the

autocorrelation between variables for state monitoring. When abnormal situations occur, the multivariable statistics Hotelling-$T^2$ and prediction error SPE control chart are used to determine whether the system has faults.

The quantity A of the main elements is determined by cumulative variance contribution rate or cross validation method, where the $T^2$ statistic is defined as follows:

$$T^2 = [t_1, t_2 \cdots t_A]\Lambda^{-1}[t_1, t_2 \cdots t_A]^T, \Lambda^{-1} = \frac{1}{n-1}T^T T \quad (13)$$

The $T^2$ statistical control limit used to monitor whether a fault has occurred conforms to the $F$-distribution:

$$T^2_{A,n,\partial} \sim \frac{A(n^2-1)}{n(n-A)}F_{A,n-A,\partial} \quad (14)$$

The SPE statistic is calculated using the following formula:

$$SPE = \sum_{i=1}^{n} t_i^2 - \sum_{i=1}^{A} t^2 = \sum_{i=A+1}^{n} t_i^2 \quad (15)$$

The control limit of $SPE$ statistic is determined by weighting coefficient $g$ and degree of freedom $h$.

$$SPE_{k,\alpha} \sim g_k \chi^2_{h,\partial} \quad (16)$$

$$g_k = v_k / 2m_k, h_k = 2m_k^2 / v_k \quad (17)$$

$v_k$ is the mean of the prediction error for the square of time $k$, and $m_k$ is the variance of the prediction error for the square of time $k$.

## 2.4 Identification of Fault Source Variables

The existing machine learning based health management technologies require obtaining a large amount of sample data for diagnosis and prediction. In the early stages of device deployment and operation, due to the lack of effective fault sample information and the low reliability of fault samples in the early stages of device operation, it can lead to inaccurate fault prediction results and other issues. The data-driven fault prediction method proposed in this article can provide early warning for equipment without or with only a very small amount of fault samples. The core idea is to use multivariate statistical contribution maps for fault diagnosis, identify fault source variables, provide decision-making basis for fault localization, and achieve visual localization of equipment faults.

The existing machine learning based health management technologies require obtaining a large amount of sample data for diagnosis and prediction. In the early stages of device deployment and operation, due to the lack of effective fault sample information and the low reliability of fault samples in the early stages of device operation, it can lead to inaccurate fault prediction results and other issues. The data-driven fault prediction method proposed in this article can provide early warning for equipment without or with only a very small amount of fault samples. The core idea is to use multivariate statistical contribution maps for fault diagnosis, identify fault source variables, provide decision-making basis for fault localization, and achieve visual localization of equipment faults.

For statistics, contribution graphs can be used for fault diagnosis. According to the definition, $T^2$ can be expanded as follows:

$$T^2 = t_1^2 + t_2^2 + \cdots t_A^2 \quad (18)$$

The contribution of the $a$-th principal component $t_a$ to $T^2$ is:

$$C_{t_a} = t_a^2 / T^2 (a = 1, \cdots, A) \quad (19)$$

According to the definition of the principal component fraction, the contribution of the process variable $x_j$ of the $a$-th major component can be inferred backwards:

$$t_a = x p_a = [x_1, \cdots, x_m] \cdot \begin{bmatrix} p_{1,a} \\ \vdots \\ p_{m,a} \end{bmatrix} = \sum_{j=1}^{m} x_j p_{j,a} \quad (20)$$

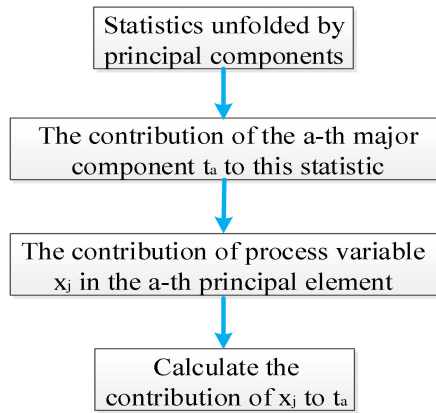The contribution of $x_j$ to $t_a$ is:

$$C_{t_a, x_j} = x_j p_{j,a} / t_a (a = 1, \cdots, A; j = 1, \cdots, m) \quad (21)$$

The $SPE$ contribution chart is simpler and more intuitive than the $T^2$ contribution chart. According to the definition of $SPE$ statistics, the contribution of each variable to $SPE$ is:

$$C_{SPE, x_j} = sign(x_j - \hat{x}_j) \cdot \frac{(x_j - \hat{x}_j)^2}{SPE} \quad (22)$$

In the formula, $sign(x_j - \hat{x}_j)$ can extract the positive and negative information of residuals. When using a contribution chart, the obtained variable contribution rate vector can be normalized to a vector with a modulus of 1, and then the contribution of each variable can be plotted using a bar graph. Process variables with a high contribution rate to statistics are affected by abnormal operating conditions, and

valuable fault information can be obtained by combining process knowledge. The calculation process of contribution is shown in Figure 1.



**Figure 1. The Calculation Process of Contribution**

## 3. Dataset Validation

### 3.1 Dataset Validation
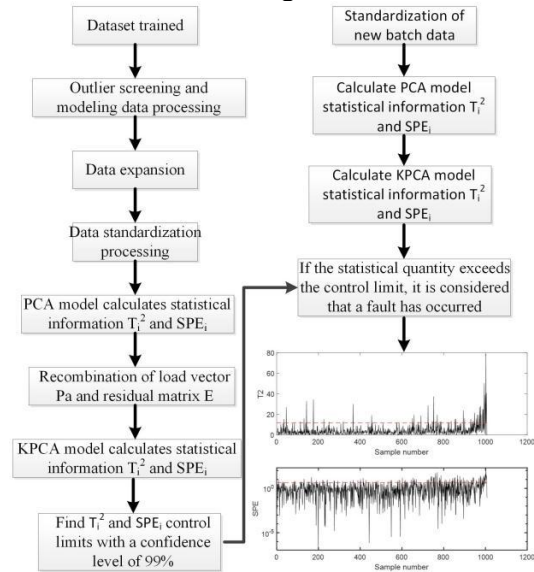
The state monitoring model works in the following way:

① For the collected data $x_{new}$, standardize the model data using the variance and mean of the data;

②Use weighted vector $\omega_i(i=1,2\cdots J)$ for data fusion, i.e. $x_{new,i} = x_{new} \times \omega_i$, to obtain $x_{new,1}, x_{new,2} \cdots x_{new,J}$;

③ The model information is called and the $T_i^2$ and $SPE_i$ statistics of $x_{new,i}(i=1,2\cdots J)$ are calculated in the model;

④ When the statistics in the model exceed the control limit, it indicates that the system has malfunctioned.

During the operation of the equipment, if a malfunction occurs or may occur, statistics such as $T^2$ and $SPE$ will exceed the corresponding control limits, and a warning will be given for the malfunction. The technical path for abnormal device operation warning is shown in Figure 2:

The vehicle data recorder can export bus data, as shown in Figure 3. The vehicle data recorder is a component of the vehicle's comprehensive electronic system, with the main mission of collecting, storing, and exporting vehicle status information during vehicle use. Through data query and analysis,

it provides data support for vehicle use, maintenance, and management.



**Figure 2. Technical Path for Abnormal Equipment Operation Warning**



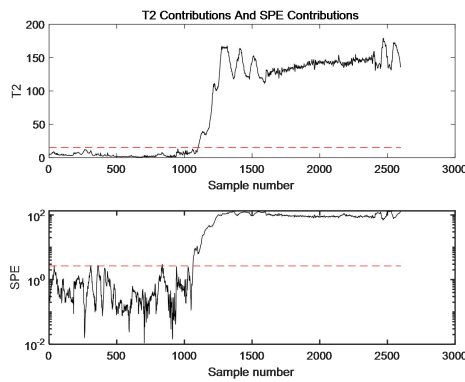**Figure 3. Vehicle Data Recorder**

Verify the proposed algorithm using real vehicle data, export relevant state data during stable vehicle operation, and select ten dimensional data such as intake temperature, intake pressure, oil pressure, and oil temperature from engine related data as state variables for analysis. The ten dimensional state variables are shown in Table 1:

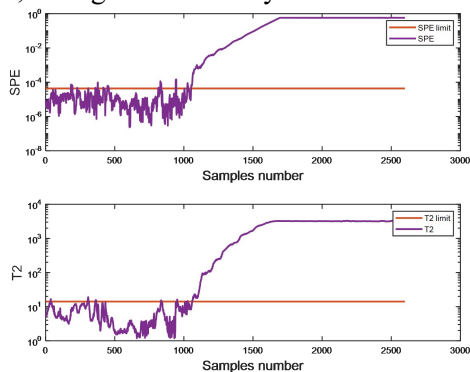**Table 1. State Variables used for Modeling a Certain type of Vehicle**

| Variable | Description of variables | Unit |
|---|---|---|
| $x_1$ | Intake air temperature | °C |
| $x_2$ | Intake air pressure | kPa |
| $x_3$ | Oil pressure | kPa |
| $x_4$ | Oilt engine-oil-temperature | °C |
| $x_5$ | Exhaust temperature | °C |
| $x_6$ | Engine speed | rpm |

| $x_7$ | Camshaft speed | rpm |
| $x_8$ | Fuel supply duration | °CA |
| $x_9$ | Fuel injection advance angle | °CA |
| $x_{10}$ | Coolant temperature | °C |

When the vehicle experiences gradual faults and abnormal working conditions during operation, the system will issue a warning. When a malfunction occurs or a safety risk is predicted in the system, the $T^2$ and SPE statistics of the model will exceed the corresponding control limits, thereby achieving early warning of abnormal vehicle operating conditions. The process of model warning when the engine oil temperature is abnormal is simulated. A simulated fault signal is added to the engine oil temperature data sequence. After inputting the data into the model, the verification results are shown in figure 4 and figure 5:
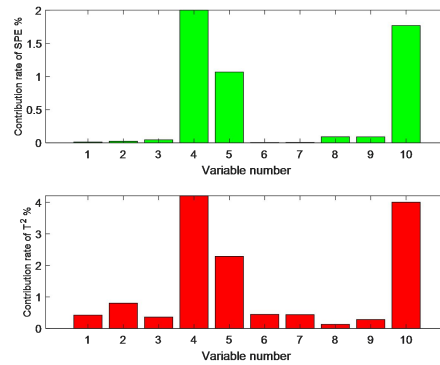


**Figure 4. Monitoring Results of PCA Model**

From the figure, it can be seen that when a fault occurs, the statistical values of $T^2$ and SPE both change, exceeding the control limit and triggering an alarm. In practical applications, this model can respond to faults and quickly trigger alarms in a short period of time, with good sensitivity and timeliness.



**Figure 5. Monitoring Results of KPCA Model**

## 3.2 Identification of Fault Sources

This article proposes a method for fault identification using multivariate contribution graphs. This method identifies the source variables of each fault one by one when a fault is detected. The figure 6 shows the statistical contribution graph. The variable index is consistent with the serial number of the running state variable. After a malfunction occurs, first locate the variables that have a significant contribution to the fault, followed by other variables that exceed the reference range. This method can accurately identify the fault source variables when there are abnormalities in the equipment status parameters, and can assist maintenance personnel in diagnosing and locating equipment faults.



**Figure 6. Identification of Fault Source Variables for Detecting Real Vehicle Faults Based on Statistical Contribution Graph**

## 4. Conclusion

This article establishes a state monitoring model based on multivariate statistical analysis and verifies the effectiveness of the proposed method using actual equipment data. Through experiments, it can be seen that the proposed method has high sensitivity and good real-time performance, and can meet the fault warning requirements under multiple working conditions. Aiming at the problem of difficult identification of fault sources, a method for fault identification using multivariate statistical contribution graphs is proposed to accurately identify the fault source variables under abnormal vehicle conditions and locate the fault source. The research conclusions and achievements can provide ideas and basis for the design of vehicle anomaly warning models.

## References

[1] YANG Dawei, ZHAO Yongdong, FENG Fuzhou, et al. Planetary Gearbox Fault Feature Extraction Based on Parameter Optimized Variational Mode Decomposition and Partial Mean of Multi-scale Entropy. Acta Armamentarii, 2018, 39(09): 1683-1691.

[2] ZHANG Cheng, DAI Xu-Nian, LI Yuan. Fault Detection and Diagnosis Based on Residual Dissimilarity in Dynamic Principal Component Analysis. ACTA Automatica Sinica, 2022, 48(1): 292-301.

[3] HAN Wanli，MAO Dajun，YIN Qimin. Induced Draft Fan Fault Warning based on PCA and Multivariate State Estimation Technique. Journal of Engineering for Thermal Energy and Power, 2020, 35(01): 92-97.

[4] XU Li, ZHU Pengdong, GU Hongjie, XU Wencai. Early Warning of Current Carrying Faults in Power Equipment Based on Variable Scale PCA. Electric Power Automation Equipment, 2012, 32(05): 147-151.

[5] WU Kai, SUN Yanguang, ZHANG Lin. Fault Diagnosis of Strip Breaking in Hot Strip Rolling Based on Kernel Principal Component Analysis. China Metallurgy, 2020, 30(11): 60-65.

[6] GUO Jinyu, WANG Xin, LI Yuan. Fault Detection in Chemical Processes Using Weighted Differential Principal Component Analysis. Journal of Chemical Engineering of Chinese Universities, 2018, 32(1): 183-192.

[7] YUAN Zhongshuai, SUN Sitong. Multimodal process fault monitoring of LNS-PCA based on local information. The Chinese Journal of Process Engineering. 2023, 23(5): 150-158.

[8] YAO Yuman, LUO Wenjia, DAI Yiyang. Research progress of data-driven methods in fault diagnosis of chemical process. Chemical Industry and Engineering Progress, 2021, 40(4): 1755-1764.

[9] WANG Qingfeng, WEI Bingkun, LIU Jiahe, MA Wensheng, XU Shujian. Research on Construction and Application of Data-driven Incipient Fault Detection Model for Rotating Machinery. Journal of Mechanical Engineering, 2020, 56(16): 22-30.

[10] HUANG Weiguo, LI Shijun, MAO Lei, et al. Research on Multi-source Sparse Optimization Method and Its Application in Compound Fault Detection of Gearbox. Journal of Mechanical Engineering, 2021, 57(7): 88-98.

[11] YANG Xinmin, GUO Yu, TIAN Tian, ZHU Yungui. Early fault detection index of rolling bearing based on integrated envelope spectrum. Journal of Vibration and Shock, 2023,42(10): 67-73.

[12] Zhang Cheng, Guo Qingxiu, Li Yuan, Gao Xianwen. Fault detection strategy based on difference of score reconstruction associated with principal component analysis. Control Theory & Application, 2019, 36(5): 774-782.

[13] Liu Qiang, Zhuo Jie, Lang Ziqiang, Qin S. Joe. Perspectives on data-driven operation monitoring and self optimization of industrial processes. Acta Automatica Sinica, 2018, 44(11): 1944-1956.

[14] Zhang Cheng, Guo Qingxiu, Li Yuan. Fault Detection Method Based on Principal Component Difference Associated With DPCA. Journal of Chemometrics, 2018, 33(4): 3082.

[15] Zhang Yan, Kuang Hewei. Settlement Prediction of Highway Subgrade with Reterance Vector Machine Based on Principal Component Analysis. Science Technology and Engineering, 2020, 20(1): 312.