Improved Batch Process Monitoring and Diagnosis Based on Multiphase KECA

Yongsheng Qi*, Yuan Wang*, Chenxi Lu*, Lin Wang*

* Inner Mongolia University of Technology, Huhhot, Inner Mongolia, 010051, China (e-mail: qyslyt@163.com).

Abstract: Multiple phases with transitions from phase to phase are important characteristics of many batch processes. The linear characteristics of batch processes are usually taken into consideration in the traditional algorithms while the nonlinearity is neglected. However, to monitor batch processes more accurately and efficiently, such process features are needed to be considered carefully. In this paper, a new similarity index based on KECA (kernel entropy component analysis) is defined for batch processes with nonlinear characteristics. A new phase division and monitoring method based on the proposed similarity index is brought forward simultaneously. First, nonlinear characteristics can be extracted in feature space via performing KECA on each preprocessed time-slice data matrix. Then phase division is achieved with the similarity change of the extracted feature information. By establishing a series of KECA models for transitions and steady phases, it reflects the diversity of transitional characteristics objectively and can preferably solve the stage-transition monitoring problem in multistage batch processes. Finally, in order to overcome the problem that the traditional contribution plot cannot be applied to the kernel mapping space, a nonlinear contribution plot diagnosis algorithm is proposed. Both results of simulation study and industrial application clearly demonstrate the effectiveness and feasibility of the proposed method.

Keywords: KECA; fault monitoring; fault diagnosis; batch process

1. INTRODUCTION

It is acknowledged that the multiplicity of operation stages is an inherent nature of many batch processes and each stage exhibits significantly different underlying behaviours. Up to now, multivariate statistical methods such as principal component analysis (PCA) and partial least square (PLS) have been successfully used in modelling multivariate continuous processes (Jiang et al.(2016)). However, it is difficult for PCA to reveal the changes of process correlations because it takes the entire batch data as a single object in modelling, neglecting the local behaviours within batch process stages. Thus, the unique process correlation information is not reflected. Furthermore, PCA is a kind of modelling method dealing with linear data sets, essentially. However, in industries, most batch processes possess nonlinear characteristics. Therefore, PCA is unsuitable for monitoring those nonlinear processes. Lu et al. (2004) developed a stage-based sub-PCA modelling method based on the fact that changes of the process correlations may relate to its stages diversity in multistage batch processes, which does not require fulfilling missing process observations and preserves the dynamic relationships.

However, their strict stage partition algorithm neglects the stage-to-stage transiting characteristics, which compromises the accuracy of sub-stage representative monitoring models. As a complement of sub-PCA method, Zhao *et al.* (2008) proposed a soft-transition multiple PCA (STMPCA) method which can identify and model both process phases and transitions between two neighboring stages. Wang *et al.* (2016) used fuzzy pattern recognition method calculating the

proximity degree of centre-of-gravity between two adjacent time-slice loading matrices along variable direction to achieve sub-stage division. Since changes of the process correlations may relate to its stages diversity in multistage batch processes, all these aforementioned methods do deep researches on this changing trend to figure out the changes in the internal operating mechanism of a process. In a word, these methods perform PCA on time slices to extract process features, generating time-slice loading matrices which are fed as the input to the existing clustering algorithms, or achieve the division and modelling of sub-stages, according to the similarity between two adjacent loading matrices. But methods based on PCA can only extract linear characteristics of processes, ignoring the nonlinear correlations among variables. Therefore, these methods mentioned above are not applicable to complex batch processes with nonlinear characteristics. Kernel entropy component analysis (KECA) method can extract nonlinear characteristics of batch processes effectively via performing nonlinear mapping with Renyi entropy. However, kernel mapping needs nonlinear functions which are commonly unknown, and most clustering algorithms require information about the loading matrix itself during the process of clustering, especially when calculating class centers. This means that in the sub-stage division of nonlinear batch processes, loading matrices cannot be fed as the input to the clustering algorithm directly. Therefore, a novel similarity index based on KECA is defined in the present article, and a nonlinear clustering algorithm based on KECA similarity is proposed simultaneously.

PCA is widely used in fault modelling and monitoring for batch processes. However, this kind of method is based on a

linearity assumption, which shows poor performance when utilized in batch processes involving complex nonlinearities. Scholkopf et al. (1998), therefore, proposed kernel principal component analysis (KPCA), where nonlinear data in the input space were equivalently transformed into linear data in the high-dimensional feature space via a nonlinear mapping. Robert Jenssen et al. (2010) proposed a new spectral data transformation method termed kernel entropy component analysis (KECA), which reveals angular structure related to the Renyi (quadratic) entropy of the input space data set and does not necessarily use the top eigenvalues and eigenvectors of the kernel matrix. Results demonstrate that PCs selected with KECA method show distinct angular structure. Thus, KECA shows better performance of fault detection than other methods. In our work, a new Cauchy-Schwarz (CS) divergence measurement used to describe the angular structure revealed by KECA is brought forward, which can preferably depict the similarity between different probability density distributions and distinguish the anomalies effectively.

When a fault is detected, further diagnosis is needed. There are different methods available for detection and diagnosis of batch processes, such as contribution plot and fault reconstruction. Yoon et al. (2001) proposed fault diagnosis method using contribution plot. Since the traditional contribution plot represents the percent contribution of the original measurement variables to the monitoring statistics, different monitoring methods need to derive the corresponding contribution quantity expressions. As for complex fault monitoring methods, such as kernel entropy learning method, it is difficult to construct a corresponding formula to calculate the contribution plot, which limits the application of the contribution plot method to a great extent. Yue et al. (2001) proposed a fault diagnosis method based on fault reconstruction. However, various history fault data are needed for this method while fault process variables can be located without fault data using contribution plot method. In view of the limitation of the contribution plot and the fault reconstruction methods, a novel Standard Vector Kernel Contribution Diagram (SV-KCD) method based on standard vector is introduced in the present article. The proposed method reconstructs the monitoring samples at fault time directly. This method is described quantitatively with histogram, being more intuitive and practical.

The article is organized as follows. Phase division algorithm based on KECA similarity index is outlined in Section 2. Fault monitoring for each phase/stage and SV-KCD fault diagnosis based on KECA are designed in Section 3. Simulation and Industrial application results are presented in Section 4. In Section 5, the conclusion is drawn ultimately.

2. PHASE DIVISION BASED ON KECA SIMILARITY

2.1 KECA

Assuming that the data set S: $x_1, ..., x_N$ are generated from an underlying probability density function p(x), the Renyi quadratic entropy is defined as

$$\widehat{V}(p) = -\log f(p) = -\log p^2(x)dx \tag{1}$$

and Parzen window density estimator is described as

operator, we then have

$$\hat{p}(x) = \frac{1}{N} \sum_{x_i \in D} k_\sigma(x, x_i)$$
(2)

where $k_{\sigma}(x, x_i) = \exp(-\frac{\|x - x_i\|}{2\sigma^2})$ is a Mercer function, σ is the parameter of the kernel function, or Parzen window. Using the sample mean approximation of the expectation

$$\hat{V}(p) = \frac{1}{N} \sum_{x_i \in D} \hat{p}(x_i) = \frac{1}{N^2} \sum_{x_i} \sum_{x_i} k_{\sigma}(x, x_i) = \frac{1}{N^2} I^{\mathsf{T}} K I$$
(3)

Where *K* is the $(N \times N)$ kernel matrix and *I* is an $(N \times 1)$ vector where each element equals one. To the end, the Renyi entropy estimator may be expressed in terms of eigenvalues and the corresponding eigenvectors of the kernel matrix, which may be eigen-decomposed as $K=EDE^{T}$ with *D* being a diagonal matrix storing the eigenvalues $\lambda_1, ..., \lambda_N$ and *E* being a matrix with the eigenvectors $e_1, ..., e_N$ as columns. Rewriting (3), yields

$$\hat{V}(p) = \frac{1}{N^2} \sum_{i=1}^{N} \left(\sqrt{\lambda_i} e_i^{\mathrm{T}} 1 \right)^2$$
(4)

Eq. (4) demonstrates that each eigenvalue and eigenvector makes different contributions to the Renyi quadratic entropy. Therefore, in the kernel entropy component analysis, the *l* number of eigenvalues which make the first *l* largest contributions to the Renyi entropy and their corresponding eigenvectors are selected, forming the PC matrices which are defined as $\varphi_{eca} = D_i^{\frac{1}{2}} E_i^{\mathrm{T}}$. Then the inner product of the data in the feature space is obtained $K_{eca} = \varphi_{eca}^{\mathrm{T}} \varphi_{eca}$.

2.2 A new similarity index based on KECA

Assuming that data set $X \in \mathbb{R}^N$ is projected onto the feature space F with KECA, yielding $\varphi = \{\phi_{(x_1)}, \dots, \phi_{(x_M)}\}$. Set Eq. (4) as the guideline of the direction of the KECA projection and the projection vector P is defined as

$$P = \frac{1}{\sqrt{\lambda_i}} \varphi e_i \tag{5}$$

To quantify the similarity of the two projection vectors, the similarity is defined as follows

$$D_{1} = diss(P_{1}, P_{2}) = \frac{4}{J} \sum_{j=1}^{J} (\lambda_{1}^{j} - 0.5)^{2} = \frac{4}{J} \sum_{j=1}^{J} (\lambda_{2}^{j} - 0.5)^{2}$$
(6)

Where P_1 and P_2 represent two different kernel entropy loading matrices. When λi approaches to 0.5, consider two loading matrices as similar whereas when λi approaches to 1 or 0, two matrices show a great difference.

2.3 Phase division and transition identification

Multiphase/multistage are divided into steady phases and transitions with the proposed similarity index. Since the

process characteristics are similar in each phase, a unified model can be established for the data in the same period.

The procedure is plotted in Fig.1 and the detailed description is given below.

1) Unfold three-way batch process data matrix along batch direction, and normalize the unfolded data matrix.

2) Perform KECA on each time-slice data matrix yielding loading matrix P_i which characterizes the correlation information between process variables.

3) For each $P_i = (I \times J)$, i = 1, 2, ..., K, calculate the similarity index $D_{1i}(k) = diss(P_i, P_i)$ (7)

where D_{1i} is fed as input samples to clusters.

4) Fuzzy C-means clustering (FCM) algorithm is used for phase division in the present article. According to the principle of maximum membership degree, the process is initially divided into C stages. A univariate control plot is used to monitor outliers of the maximum membership value for each stage. Since the samples detected as outliers mainly appears at the beginning and end of each phase, the starting and ending moments of a transition phase can be identified accordingly. Then, remove the transition process and the retaining part of the phases is identified as the corresponding stable stage. Here, the transition identification makes use of univariate statistical monitoring method to identify the transitions as outliers.

5) To set the control limits reasonably with the existence of outliers, the iterative calculation can be performed. In each iteration run, the control limits can be calculated, and the outliers are detected and removed from the reference set. Then in next run, the control limits can be re-calculated based on the updated reference set. Such steps are repeated until the control limits converge to a certain values. Since the control limits are calculated statistically instead of user-defined, the determination of the ranges of transitions is more reasonable and objective.

3. KECA BASED MONITORING AND DIAGNOSIS

3.1 CS statistic

It has been known that a distinct angular structure among the data set is led by KECA, where different clusters are distributed more or less in different angular directions. Therefore, an appropriate statistic is a must for fault monitoring. The *CS* divergence measurement between probability density functions corresponds to the cosine of the angle between kernel feature space mean vectors, which is able to express the angular structure.

The *CS* divergence is a measure of the "distance" namely similarity between two probability density functions $p_1(x)$ and $p_2(x)$, given by

$$CS = 1 - \cos \angle (M_k, M_k^i) = 1 - \sum_{j=1}^{l} \frac{m_{k,j}^{\mathrm{T}} m_{k,j}^i}{\|m_{k,j}\|} \|m_{k,j}^i\|}$$
(8)

in which $M_k = \frac{1}{I} \sum_{i=1}^{l} M_k^i$, $M_k^i = [m_{k,i}^i, m_{k,2}^i, ..., m_{k,l}^i]$. M_k^i is the PC matrix for batch *i* at sample time *k*, whereas M_k is the mean value of PC matrices of the total *I* batches at sample instance *k* and *l* is the number of PCs.

Monitoring models show high similarity under normal operating conditions and the value of CS statistic has been kept under control limit. Once the faults occurred, the value of CS statistic will increase above the control limit rapidly, whereas the similarity between two models will drop dramatically. The control limit R of CS statistic is calculated with kernel density estimation.



Fig. 1. Procedure of phase division

3.2 Steady phase modelling

Since the phase division is achieved, KECA models can be built for both steady phases and transitions. Now, the steady phase modelling procedure can be summarized as follows.

1) Reorganize the normalized data into a three-way data matrix and split it into K number of time-slice loading matrices. For each steady phase, unfold the 2-D time-slice loading matrix to $X_c(Ik_c \times J)$ where k_c is the number of samples included in phase C.

2) The kernel matrix K is calculated for the preprocessed time-slice matrix followed by the eigen-decomposition of each matrix K. Calculate the Renyi entropy corresponding to each eigenvalue with Eq.(4), and the l number of eigenvalues which make the first l largest contributions to Renyi entropy and their corresponding eigenvectors are selected, forming PC matrices.

3) Calculate *CS* statistics for each time interval k with Eq. (8) and R_l , the control limit of *CS* statistics, is calculated with the kernel density estimation.

3.3 Transition periods modelling

1) For each transition regions, unfold the preprocessed 2-D time-slice loading matrices along variable direction to matrix

 $X_m(Ik_m \times J)$ where k_m represents the number of samples included in phase *m*.

2) Build sliding weighted KECA models for input sample X_m . From the overall trend, process characteristics are usually similar to the previous steady phase at the beginning of transition periods. Then by going through some trajectories, processes transit to the next steady phase gradually. Thus, the weighted KECA model can be built for X_m and the PC matrix *M* is obtained as follows.

$$M = \lambda M_1 + (1 - \lambda)M_2 \tag{9}$$

In which λ is weighting coefficient, M_1 and M_2 represent the PC matrix of the steady phases before and after transition period, respectively.

3) Calculate *CS* statistics for each time interval k with Eq. (8) and R_2 , the control limit of *CS* statistics, is calculated with the kernel density estimation.

3.4 Fault diagnosis based on SV-KCD

Since it is impossible to find an inverse mapping from highdimensional feature space to low dimensional input space, the contribution formulation of the corresponding statistic cannot be deduced. Therefore, the traditional contribution plot cannot be applied in this paper.

To solve the problem above, a kernel space contribution plot method based on standard vectors (SV-KCD) is proposed in our work. This method not only preserves the advantage of simplicity and intuition of traditional contribution plot, but is free of the deduction of contribution formulation of the corresponding statistic. What's more, the fault samples are also not required. In a word, this method can be applied to any kernel mapping methods, such as KPCA, KICA etc. The schematic diagram of the SV-KCD method is shown in Fig.2. Supposing that the original data space is a 3-D space. Project the original data space to the high-dimensional feature space. Assuming three principal components are retained after projection. As Fig.2 (a) shows, at any sample time k, assuming normal data set is transformed into feature space. The projection data are gathered in a sphere with radius of r. After calculating CS statistics for new data set, one can readily realize that CS statistics located under the CS control limit.

Assuming that there is a central vector in the feature space, which is located at the centre of the sphere. The corresponding point of this vector in the original data space is depicted at point O in Fig.2. Set vector O as the standard vector. When a fault occurs at time K, each variable of the standard vector O is replaced by the corresponding variable of fault sample in the original data space at time K in turn. Then, fed as a new input sample, each vector with replaced variable is utilized in process monitoring. Thus, the contribution of each variable to the CS statistic can be described quantitatively, and whether the contribution is within a reasonable range can be known through the control limit, which provides more reasonable evidence for fault identification. As the Fig.2 (b) shows, it is readily to know whether the CS statistic goes beyond the control limit when variable x and y are replaced, respectively.

Now, finding out the standard vector at point O comes first. Since the inverse mapping of feature space to the original space cannot be obtained, it is impossible to find the standard vector O in the original space. However, from Eq. (8), it is transparent that the sample with the smallest *CS* value lies closer to the centroid vector O than any other sample. Therefore, the corresponding data sample in the original space can be regarded as standard vector at time k.



(b) Abnormal condition

Fig. 2. Schematic diagram of SV-KCD method



(a) FCM clustering result

(b) membership grades

(c) Transition ranges identification (d) Sketch map of the similarity

Fig.3 Phase division result using the proposed method



Fig.4 Monitoring results using MPCA for fault 5



Fig.5 Monitoring results using sub-PCA for fault 5



Fig.6 Monitoring and diagnosis results using KECA method

4. IILLUSTRATION AND DISCUSSION

4.1 Simulation validation

In the present simulation experiment, a total of 35 reference batches are generated using a simulator developed by the monitoring and control group of the Illinois Institute of Technology. A total of 10 process variables are selected to be monitored in this work. The duration of each batch is 400h, consisting of a preculture phase of about 45h and a fed-batch phase of about 355h. Fig.3 (a) shows the clustering result of the fed-batch phase using the proposed method, where the real fed-batch phase is sub-divided roughly into three main stages. So the whole process is divided into four primary stages as well as corresponding transition regions. The more

elaborate stage partition results emphasize the changes of process correlations rather than the physical operation, which will benefit making more detailed analyses of underlying process behaviours and establishing more appropriate monitoring models. It is clear that the phase division is consistent with the process nature. In fed-batch penicillin cultivation process, process nature changes with operation time. Such changes can be indicated with the trends of the similarity values Sim(k, phase c), as shown in Fig.3(b). Sim(k,c) is the membership grade between the kth time-slice data matrix and the cth cluster-centre KECA model. The values of Sim(k,c) changes gradually with the process operation. It becomes larger when the process approaching to phase c. During phase c, Sim(k,c) keeps large values which indicate high similarities between the time-slice data matrices and the current cluster-centre. When process operates far away from phase c, the similarities become small again. The gradual changes of Sim(k,c) values at the beginnings and ends of phases confirm the existence of transitions from phase to phase. The univariate statistical process control plots are utilized to identify the transition ranges. For each steady phase, the values of dissimilarity 1-Sim(k,c) $(k \in c)$ are plotted in Fig.3(c). As introduced in Section 3, the successive outliers at the beginning and the end of each phase indicate the transition ranges. After transition identification, each steady phase and transition range are modelled with the method described in Section 3. The membership grades are plotted in Fig. 3(d). Again, the transition attributes from phase to phase are shown clearly. Here are the results: steady phases (1~48), (70~188), (215~400), transition regions (49~69), (189~214). The models constructed using traditional MPCA, sub-PCA and the proposed method are then tested against monitoring of three different operating states batches.

For Fault 5, a linear decrease of slope 0.2% is imposed on the agitator power from time 100h until the end of the batch. Fig. 4 to 6 display the comparison of the detection results for Fault 5 using MPCA, sub-PCA and the proposed method, respectively. As depicted in Fig.6, the values of *CS* statistics increase sharply beyond the confidence limit right at time 47h when the fault is introduced, which is about 21h and 15h earlier than that of MPCA method and the sub-PCA method in *SPE* control plots, respectively, which also enables the operator to respond rapidly as soon as the abnormality occurs. Although the fault ends at 220h, the process correlation has been deteriorated and cannot return to the normal trajectory. So the *CS* values yield the decreasing trend around 220h but

still outside the normal boundary until the end of the batch. In T^2 control plots, the fault detection is achieved at 100h with sub-PCA method, which is about 55h slower than the result of proposed method in *CS* control plot. In addition, no abnormality in T^2 control plot is found with MPCA method. Analysis shows that the fault happened in the transition stage 1. Since the sub-PCA method divides phases into several sub phases with a kind of hard partition, ignoring the transition characteristics from phase to phase, faults are detected with an obvious time delay.

Once the abnormal condition is detected by the monitoring charts, the contribution plot obtained by using SV-KCD method is utilized to analyse the fault cause, which can indeed enhance the process understanding and improve the ability to fault detection and diagnosis, as exhibited in Fig.6. It transparently shows that the primary fault cause variable is variable 2 (Agitator power), which is well agreed with the real state. The comparison of fault detection performance with these 3 kinds of fault monitoring methods is exhibited in Table 1, the proposed method shows high efficiency to detect all kinds of faults and reaches the lowest false alarm rate among three methods which demonstrates that the proposed method can improve the reliability of the monitoring process to a certain extent.

4.2 Industrial data validation

The proposed method is applied to penicillin fermentation process monitoring of a pharmaceutical company in Hebei province, as is shown in this section. The fermentation system adopts SIEMENS PLC control system which can both detect and implement the control of temperature, pH, Aeration rate, agitator power etc. in real time. The duration of each batch is 212h, with a sampling interval of 4h. A total of 9 primary variables are monitored to reflect the cell growth and product synthesis in penicillin fermentation. 24 normal batches are selected as the initial modelling reference databases and a three-way data matrix $X (24 \times 9 \times 53)$ is obtained, simultaneously. The identification of steady phases and transition regions is achieved with the proposed method. There are three steady phases in total including sampling intervals (1-5), (9-21), (26-53). The transition range from phase I to II includes sampling interval 6-8, the transition range from phase II to phase III includes sampling interval 22-25.

Table 1. Monitoring results for three methods

Error rate of type I (%)			Error rate of type II (%)		
MPCA	sub-PCA	KECA	MPCA	sub-PCA	KECA
5.81	2.52	1.87			
1.54	0.79	0.67	10.46	45.08	6.6
5.71	1.43	1.67	6.97	0	0
7.71	2.8	4.1	12.33	0.9	0.78
3.84	1	1	22.4	40.7	12
3.13	1.36	1	7.9	39.9	2

Fig.7 depicted the monitoring result of fault batch 198. Due to mechanical reasons, aeration rate is reduced. However, the process returned to normal conditions with the timely adjustment of the operator. Since the final product is in line with the production requirements, the batch satisfies the final

product quality. In *CS* control plots, the proposed method is able to detect the fault at time 100h, which shows the effectiveness of detecting abnormalities and objectively reflect the efforts made by the operators to eliminate the faults after the failure occurs, and accurately reflect the quality of the final products, avoiding false alarms. Figure 7 shows the SV-KCD diagnosis results at steady phase *III* with the proposed method, which clearly demonstrates that the primary fault cause variable is variable 6 (Aeration rate).



Fig. 7. Monitoring and diagnosis results using KECA method

5. CONCLUSIONS

To overcome the irrationality of boundary data division between two adjoining clusters and nonlinear problems of transition processes, and improve the reliability and sensitivity of process monitoring, a novel multiphase KECA monitoring strategy has been proposed in our work. False alarm rate and missing alarm rate of online monitoring transitional data can be decreased with the proposed method, when processes transit from one stage to another. Both simulation and industrial application demonstrate that the diversity of characteristics in each phase is preferably reflected using our strategy. It also shows great values in solving fault detection problems widely existing in batch processes.

REFERENCES

- Zhao, C.H., Wang, F.L, Yao, Y, et al(2010). Phase-based statistical modeling, online monitoring and quality prediction for batch processes. *Acta Automation Sinica*, 36(3), 366 - 374.
- Jiang, Q.C, Yan, X.F(2016).Performance-Driven Distributed PCA Process Monitoring Based on Fault-Relevant Variable Selection and Bayesian Inference. *IEEE Transactions on Industrial Electronics*. 63(1), 377-386
- Lu, N.Y, Gao, F.R (2004). Sub-PCA modeling andon-line monitoring strategy for batch processes. *AIChE Journal*, 50(1), 255-259.
- Zhao, C.H., Wang, F.L, et al (2008). Improved batch process monitoring and quality prediction based on multiphase statistical analysis. *Ind Eng Chem Res*, 47(3), 835-849.
- Wang, P., Liu, X., Gao, X.J. (2016). Sub-stage PCA Modeling and On-line Monitoring for a Batch Process. *Journal of Beijing University of Technology*, 40(12), 1797-1803.
- Huang, J.P., Yan, X.F. (2017).Quality Relevant and Independent Two Block Monitoring Based on Mutual Information and KPCA. IEEE Transactions on Industrial Electronics, 64(8), 6518-6527
- Scholköpf, B., Smola, A., Miiller, K (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*. 10(5), 1299-1319
- JENSSEN, R (2010). Kernel entropy component analysis. IEEE Transactions on Pattern Analysis and Machine Intelligence, 32(5), 847-860.
- Yoon, S., MacGregor, J.F (2001). Fault Diagnosis with Multivariate Statistical Models Part I: Using Steady State Fault Signatures. *Journal* of Process Control, 11(4), 387-400
- Yue, H.H., Qin, S.J. (2001). Reconstruction-based fault identification using a combined index. Industrial and Engineering Chemity Rearch, 40(20), 4403-4414