

APPLICATION OF NON-GAUSSIAN BATCH MONITORING TO A BIOLOGICAL TREATMENT PROCESS

ChangKyo Yoo^{1,2,*}, Peter A. Vanrolleghem² and In-Beum Lee¹

¹*Dept. of Chemical Engineering, POSTECH, 790-784 Pohang, Korea*

²*BIOMATH, Ghent University, Coupure Links 653, B-9000 Gent, Belgium*

(*Corresponding author: ckyoo@postech.ac.kr)

Abstract: The batch biological wastewater treatment process poses an interesting challenge from the point of process monitoring characterized by non-stationary, batchwise, multiscale, and non-Gaussian characteristics. This contribution describes the monitoring on a pilot-scale sequencing batch reactor (SBR) using a batchwise multiway independent component analysis method (MICA) which can extract meaningful hidden information from non-Gaussian data. Given that independent component analysis (ICA) is superior to principal component analysis (PCA) to extract features from non-Gaussian data sets, the use of ICA may improve monitoring performance. The monitoring results of a pilot-scale SBR for biological wastewater treatment showed the power and advantages of MICA monitoring in comparison to conventional monitoring methods. *Copyright © 2005 IFAC*

Keywords: Batch monitoring, Multiway independent component analysis (MICA), Sequencing batch reactor (SBR), Wastewater treatment process

1. INTRODUCTION

Wastewater treatment process is widely known as the nonlinear, time-varying microorganisms. An efficient monitoring and modeling of such a process enhances the understanding of relevant biological phenomena and provides the basis for an operational monitoring, control and optimization strategy. As biological processes become more complex, the monitoring of biological processes is gaining importance to assess process performance and improve process efficiency and product quality. Early detection of faults can help avoid major breakdowns and incidents. In general, four tasks are involved in the process monitoring: (1) fault detection, which gives an indication that something is going wrong in the process; (2) fault identification (or diagnosis), which determines the root cause of the fault; (3) fault estimation, which assesses the size of the fault; and (4) fault reconstruction, which estimates the fault-free values. Fault detection is defined as a combination of process observations and measurements, data analysis and interpretation to detect abnormal features or effects and the isolation

of faults. Fault diagnosis involves the analysis of effects to identify aberrant variables and rank likely causes. Advice includes a synthesizing strategy to eliminate the causes and return the process to normal operating conditions (Olsson, G. and Newell, 1999).

Sequencing batch reactor (SBR) processes have demonstrated their efficiency and flexibility in the treatment of wastewaters with high concentrations of nutrient, nitrogen, phosphorous, and toxic compounds from domestic and industrial sources. A SBR has a unique cyclic batch operation, usually with five well-defined phases: fill, react, settle, draw and idle. Most of the advantages of SBR processes may be attributed to their single-tank designs and the flexibility that allows them to meet many different treatment objectives, and which is derived from the possibility of adjusting the duration of the different phases. But the SBR process is highly nonlinear, time-varying and subject to significant disturbances like hydraulic changes, composition variations and equipment failures. Small changes in concentrations or flows can affect effluent quality and microorganism growth. However, treatment

performance, the key indicator of process performance, is often only examined off-line in a laboratory. Even though operators are aware that there are some problems in treatment performance, they cannot quickly find out or predict what the causes are and when the problems will occur because most batch processes are run without any effective form of real-time on-line monitoring. Therefore, multivariate analysis and process monitoring of SBR are crucial to detect faults that can be corrected prior to completion of the batch or can be corrected in subsequent batches because it may take several days, week or even months for the biological process to recover from abnormal operation (Lee and Vanrolleghem, 2003).

Multivariate principal component analysis (MPCA) developed by Nomikos and MacGregor (1994) has been shown to be a powerful monitoring tool in many industrial batch processes. However, it has the shortcoming that the measurement variables of the batch process should be normally distributed. In this work, it is shown that multiway independent component analysis suggested by Yoo *et al.* (2004) can be used to overcome this drawback and obtain better monitoring performance.

2. THEORY

2.1 Multiway principal component analysis

Multivariate principal component analysis (MPCA) is used for the analysis and monitoring of batch process data. Batch data are typically reported in terms of batch numbers, variables and times. Therefore, batch processes are, by nature, leading to a 3-way matrix ($\underline{\mathbf{X}}(I \times J \times K)$) of data, where I is the number of batches, J is the number of variables and K is the number of times each batch is sampled. In a typical batch run, $j=1, 2, \dots, J$ variables are measured at $k=1, 2, \dots, K$ time intervals throughout the batch. There exists similar data on several ($i=1, 2, \dots, I$) similar process batch runs. This matrix can be decomposed using various three-way techniques, one of which is MPCA.

The three-way array $\underline{\mathbf{X}}$ can be unfolded in three ways, which give rise to the following two-dimensional matrices: 1) Batches \times variables at each time (time-wise unfolding), 2) Variables \times time for each batch (batch-wise unfolding), 3) Batches \times times for each variable (variable-wise unfolding). Time-wise unfolding is useful for analyzing the variability among samples, and batch-wise unfolding facilitates the analysis of the variability among batches by summarizing the information related to the measured variables and their variations over time. Variable-wise unfolding can be used to obtain information about the variability among the batch variables. In previous studies, the batch-wise unfolding method has been the most widely used method for analyzing batch process data. Its aim is model the differences of each batch run from a theoretical normal operating condition. In addition, it is suitable for deriving

estimate of final quality measures from process data that are usually not available until the batch terminates. Moreover, the majority of the nonlinear behavior of the process is eliminated by subtracting off the mean trajectory of each variable at each time. After the unfolded matrix has been mean-centered and scaled, PCA is performed. The results from PCA are the loading vectors, and the calculated scores for each batch. The loading vectors contain a weight for each variable at each time. In this paper, the batch-wise unfolding scheme in Fig. 1 is used. Therefore, MPCA is equivalent to performing ordinary PCA on a large two-dimensional matrix \mathbf{X} constructed by unfolding the three-way data in the manner shown schematically in Fig. 1.

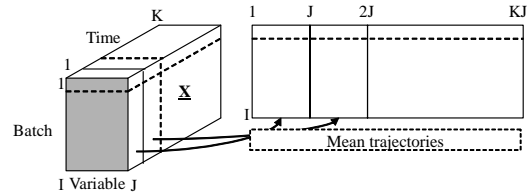


Fig. 1. Batchwise unfolding method for a three-way batch.

MPCA decomposes the three-way array $\underline{\mathbf{X}}$ into a summation of the product of a score t_r and a loading matrix P_r plus a residual array $\underline{\mathbf{E}}$ that is minimized in the least squares sense as follows:

$$\underline{\mathbf{X}} = \sum_{r=1}^R t_r \otimes P_r + \underline{\mathbf{E}} = \sum_{r=1}^R t_r p_r^T + \underline{\mathbf{E}} = \hat{\underline{\mathbf{X}}} + \underline{\mathbf{E}} \quad (1)$$

where \otimes denotes the Kronecker product ($\underline{\mathbf{X}} = \mathbf{t} \otimes \mathbf{P}$ is $\underline{\mathbf{X}}(i, j, k) = t(i)P(j, k)$), R denotes the number of principal components retained, t_r expresses the relationship among batches, p_r is related to variables and their time variation, $\underline{\mathbf{E}}$ is the residual matrix. The first expression in Eq. (1) gives the 3-D decomposition while the second expression displays the more common 2-D decomposition.

2.2 Independent component analysis (ICA)

What distinguishes ICA from other methods is that it looks for components that are both statistically independent and non-Gaussian. PCA is a dimensionality reduction technique in terms of capturing the variance of the data which is capable of extracting uncorrelated latent variables from correlated data, while ICA is designed to separate the independent components (ICs) that are independent and constitute the observed variables. Furthermore, PCA can only impose independence up to second order statistics information (mean and variance) while constraining the direction vectors to be orthogonal, whereas ICA has no orthogonality constraint and also involves higher-order statistics (Hyvärinen *et al.*, 2001). Hence, ICA may reveal more useful information in the non-Gaussian data than PCA (Hyvärinen *et al.*, 2001).

In the ICA algorithm, it is assumed that d measured variables x_1, x_2, \dots, x_d can be expressed as linear combinations of m ($\leq d$) unknown independent components s_1, s_2, \dots, s_m . The relationship between them is given by

$$\mathbf{X} = \mathbf{A}\mathbf{S} + \mathbf{E} \quad (2)$$

where $\mathbf{X} = [\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(n)] \in R^{d \times n}$ is the data matrix (in contrast to PCA, ICA employs the transposed data matrix.), $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m] \in R^{d \times m}$ is the unknown mixing matrix, $\mathbf{S} = [\mathbf{s}(1), \mathbf{s}(2), \dots, \mathbf{s}(n)] \in R^{m \times n}$ is the independent component matrix, $\mathbf{E} \in R^{d \times n}$ is the residual matrix, and n is the number of samples. Here, we assume $d \geq m$ (when $d=m$, the residual matrix, \mathbf{E} , becomes the zero matrix). The basic problem of ICA is to estimate both the mixing matrix \mathbf{A} and the independent components \mathbf{S} from only the observed data \mathbf{X} . Alternatively, one could define the objective of ICA as follows: to find a demixing matrix \mathbf{W} whose form is such that the rows of the reconstructed matrix $\hat{\mathbf{S}}$, given as

$$\hat{\mathbf{S}} = \mathbf{W}\mathbf{X} \quad (3)$$

become as independent of each other as possible (Hyvärinen *et al.*, 2001).

2.3 Multiway Independent Component Analysis (MICA)

The monitoring method based on MICA is similar to that based on MPCA. MICA is equivalent to performing ICA on a large two-dimensional matrix \mathbf{X} constructed by batchwise unfolding the three-way data matrix $\underline{\mathbf{X}}$. MICA decomposes the three-way array $\underline{\mathbf{X}}$ into a summation of the product of independent vectors \mathbf{s}_r and loading matrices \mathbf{A}_r plus a residual array $\underline{\mathbf{E}}$ so that the ICs \mathbf{s} become as independent of each other as possible:

$$\underline{\mathbf{X}} = \sum_{r=1}^R \mathbf{s}_r \otimes \mathbf{A}_r + \underline{\mathbf{E}} = \sum_{r=1}^R \mathbf{s}_r \mathbf{a}_r^T + \underline{\mathbf{E}} = \hat{\underline{\mathbf{X}}} + \underline{\mathbf{E}} \quad (4)$$

where \otimes denotes the Kronecker product ($\underline{\mathbf{X}} = \mathbf{s} \otimes \mathbf{A}$ is $X(i, j, k) = s(i)A(j, k)$) and R denotes the number of ICs retained. The \mathbf{S} and \mathbf{A} matrices in Eq. (4) can be equivalent to the loading matrix and score matrices by analogy with MPCA, *i.e.* \mathbf{S} can be regarded as the score matrix \mathbf{T} , and \mathbf{A} can be treated as the loading matrix \mathbf{P} . The i th elements of the independent vector \mathbf{s} correspond to the i th batch and summarize the overall variations in this batch with respect to the other batches over the entire history of the batch. The mixing matrix, \mathbf{A} , summarizes the time variations of the measured variables about their average trajectories. The elements of this matrix are the weights, which give the independent vectors \mathbf{s} for a batch when applied to each variable at each time interval within that batch (Yoo *et al.*, 2004).

Similar to MPCA, the key idea is to exploit the ability of MICA to extract features from three-way batch data by projecting the data onto a low-dimensional space that summarizes both the variables and their time trajectories. First, the three-

way matrix $\underline{\mathbf{X}}(I \times J \times K)$ is unfolded into a two-dimensional matrix, $\mathbf{X}(I \times JK)$ using a batchwise unfolding scheme. Second, the mean trajectory is removed from each variable and each time of the unfolded data matrix to remove the majority of the nonlinear behavior of the batch process. Third, the data matrix is normalized (*i.e.*, mean centered and standardized to unit variance). The normalized $\mathbf{X}(I \times JK)$ is then transposed, yielding the transposed matrix $\mathbf{X}_{normal}(JK \times I)$. Fourth, whitening is performed on $\mathbf{X}_{normal}(JK \times I)$ to acquire the uncorrelated whitened matrix $\mathbf{Z}_{normal} = \mathbf{Q}\mathbf{X}_{normal}$. Fifth, the matrices of \mathbf{A} , \mathbf{W} and \mathbf{S} are obtained using the FastICA algorithm. Sixth, the procedures for ordering and dimension reduction method of ICs are executed. The m rows of \mathbf{W} constitute a reduced matrix \mathbf{W}_d (deterministic part of \mathbf{W}), and the remainder of the rows of \mathbf{W} constitute a reduced matrix \mathbf{W}_e (excluded part of \mathbf{W}). Finally, the MICA model with the matrices \mathbf{W}_d , \mathbf{W}_e , $\hat{\mathbf{S}}_d$ and $\hat{\mathbf{S}}_e$ is constructed. Then, independent data vectors for a new batch k ($\mathbf{x}_{new}(k)$), $\hat{\mathbf{s}}_{newd}(k)$ and $\hat{\mathbf{s}}_{newe}(k)$, can be obtained by transformation through the demixing matrices \mathbf{W}_d and \mathbf{W}_e , *i.e.*, $\hat{\mathbf{s}}_{newd}(k) = \mathbf{W}_d \mathbf{x}_{new}(k)$ and $\hat{\mathbf{s}}_{newe}(k) = \mathbf{W}_e \mathbf{x}_{new}(k)$, respectively.

In MICA, two statistics are deduced from the process model in normal operation: the D -statistic for the systematic part of the process variation and the Q -statistic for the residual part of the process variation. The D -statistic for a batch k , also known as the I^2 statistic, is the sum of the squared independent scores and is defined as follows:

$$I^2(k) = \hat{\mathbf{s}}_{newd}(k)^T \hat{\mathbf{s}}_{newd}(k) \quad (5)$$

The Q -statistic for a batch k , also known as the SPE statistic, is defined as follows:

$$SPE(k) = \mathbf{e}(k)^T \mathbf{e}(k) = (\mathbf{x}(k) - \hat{\mathbf{x}}(k))^T (\mathbf{x}(k) - \hat{\mathbf{x}}(k)) \quad (6)$$

where $\hat{\mathbf{x}}$ can be calculated as follows:

$$\hat{\mathbf{x}} = \mathbf{Q}^{-1} \mathbf{B}_d \hat{\mathbf{s}} = \mathbf{Q}^{-1} \mathbf{B}_d \mathbf{W}_d \mathbf{x} \quad (7)$$

The confidence limits of the I^2 and SPE statistics in MICA can be obtained by kernel density estimation. Here, the I^2 value is used to detect faults associated with abnormal variations within an MICA model subspace, whereas the SPE value is used to detect new events that are not taken into account in an MICA model subspace (Yoo *et al.*, 2004).

3. RESULT AND DISCUSSION

3.1 Process description of the pilot-scale SBR system

The data used in this research were collected from a pilot-scale SBR system shown in Fig. 2. A fill-and-draw sequencing batch reactor (SBR) with a 80-liter working volume is operated in a 6h cycle mode and each cycle consists of fill/anaerobic (1h), aerobic (2h 30 min), anoxic (1h), re-aerobic (30min) and settling/draw (1h) phases. The hydraulic retention

time (HRT) and the solid retention time (SRT) are maintained at 12 hrs and 10 days, respectively. Six electrodes for pH, oxidation-reduction potential (ORP), dissolved oxygen (DO), temperature, conductivity and weight are connected to the individual sensors to check the status of the SBR, where a set of on-line measurements is obtained every one minute. The historical data set of the SBR consisted of 280 batches (70 days) for which 6 variables were measured at 300 time instants (Lee and Vanrolleghem, 2003).

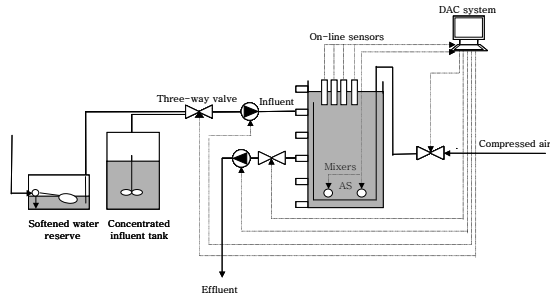


Fig. 2. Schematic diagram of the pilot-scale sequencing batch reactor.

3.2 Multivariate analysis of historical data set in SBR (MPCA and MICA)

Fig. 3 shows the monitoring result of all 280 batches of the SBR using the MPCA and MICA methods, where the dotted lines correspond to the 95 and 99% confidence limits. Five components of the MPCA model were selected by the cross-validation method. To ensure comparison of equivalent models, five ICs were selected for the MICA model. From this figure, we notice that the MICA plot shows characteristics dissimilar from the MPCA one. Compared to MPCA, MICA points to a lower number of abnormal batches in SBR. This difference can be explained by the density estimation of the SBR data. Fig. 4 (left) shows that the density estimate of the first score (t_1) in MPCA does not follow the Gaussian distribution but the ‘supergaussian distribution’ in which process variables take relatively more often values that are very close zero, where the probability density of the data is peaked in the middle and has heavy tails (large values far from zero). Thus, the T^2 and SPE charts of MPCA that are based on the assumption that the data are Gaussian distributed may cause a false result when it is used for SBR monitoring. This observation is the motivation of the MICA method because MICA is sensitive to modes whose influences on the measured variables follow a supergaussian distribution. Fig. 4(right) represents the loading plot of each variable of each time interval of the first IC. It shows the types of information that can be extracted when MICA is used in batch modeling. The loading plot obtained from MICA gives the history and identified important features of the SBR. From this figure, we notice that the DO, conductivity, and pH show large variations and have large influences during a batch, whereas ORP and weight show relatively small variations.

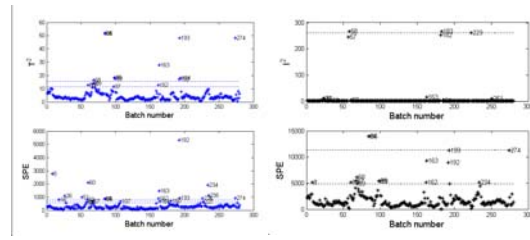


Fig. 3. Multivariate analysis of all 280 batches, (left) MPCA, (right) MICA.

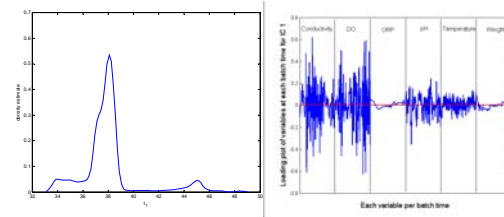


Fig. 4. The density estimate of MPCA and the variable loading plot of MICA. (left) Non-Gaussian distribution of the first principal score (t_1) obtained from MPCA, (right) Variable loading plot for the first independent score (i_1) obtained from MICA.

3.3 Batch monitoring of SBR (MPCA and MICA)

The MPCA and MICA models for the SBR monitoring were developed after an analysis of the historical SBR data set in Fig. 2. The MPCA model selected 143 batches to create a rather broad scope of normal batches, where 7 abnormal batches (batch number: 8,18,26,51,60,84,85) were excluded for the normal operating condition (NOC) model. The MICA model selected 146 batches, where 4 abnormal batches (batch number: 57, 58, 84, 85) were excluded for the normal NOC model. The test data set that consisted of the following 30 batches was projected onto the reduced MPCA and MICA model spaces. Fig. 5 shows the batch monitoring result by MPCA and MICA. While both of them could detect two abnormal batches (batch 12, 13), MPCA detected batch 9 as an abnormal batch while MICA left batch 9 as a normal batch. Actually, batch 9 is a normal batch. When MPCA is applied to non-Gaussian data, the T^2 chart of MPCA may suffer oversensitivity for normal batches, e.g., batch 9. As a data set deviates from a Gaussian distribution, the variance tends to increase and hence the T^2 statistic tends to decrease. Typically, this increases the false alarm rate of the MPCA in which a normal batch might be judged as a non-conforming one. Obviously, this deteriorates the reliability of the monitoring system.

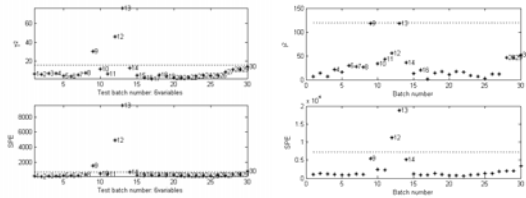


Fig. 5. Monitoring result of 30 test batches. (left) MPCA and (right) MICA. The dotted lines correspond to the 99% confidence limit.

4. CONCLUSION

This paper describes the application of a pilot-scale SBR monitoring using MICA which can extract meaningful hidden information from non-Gaussian data sets. The result showed a more powerful monitoring performance than the MPCA approach. Furthermore, the MICA method can be easily applied to most batch or fed-batch processes which have non-Gaussian distributed data.

ACKNOWLEDGEMENTS

This work was supported by the Post-doctoral Fellowship Program of the Korea Science & Engineering Foundation (KOSEF) and a Visiting Postdoctoral Fellowship of the Fund for Scientific Research-Flanders (FWO). And Lee J. is kindly thanked for his valuable discussion.

REFERENCES

- Hyvärinen, A., J. Karhunen, and E. Oja, (2001). *Independent component analysis*, John Wiley & Sons, INC., USA.
- Lee, D.S. and Vanrolleghem, P. A., (2003). Monitoring of a sequencing batch reactor using adaptive multiblock principal component analysis, *Bio&Bioeng*, **82**, 489-497.
- Nomikos, P. and J. F. MacGregor, (1994). Monitoring batch processes using multiway principal component analysis, *AIChE J.* **40**(8), 1361-1375.
- Olsson, G. and Newell, B. (1999) *Wastewater Treatment Systems: Modelling, diagnosis and Control*, IWA, UK.
- Yoo, C.K., J. Lee, P.A. Vanrolleghem and Lee, I.B., (2004). On-line monitoring of batch processes using multiway ICA, *Chemom. and Intel. Lab. Sys.* **71**(2), 151-163.

