# MODELING AND OPTIMIZATION OF BATCH PROCESS THROUGH WAVELET ANALYSIS AND MULTIVARIATE ANALYSIS

**Manabu Kano** [*] **Koichi Fujiwara** [*] **Shinji Hasebe** [*]
**Hiromu Ohno** [**]


[*] *Kyoto University, Kyoto 606-8501, Japan*
[**] *Kobe University, Kobe 657-0013, Japan*

Abstract: An efficient method is developed to build a batch process model and optimize its operating conditions including time-dependent operation profiles. In the proposed method, referred to as wavelet regression and optimization (WRO), important wavelet coefficients of operation profiles are selected as input variables of a statistical model, and then further dimensionality reduction is achieved through multivariate analysis. Then, on the basis of the developed model, optimal operation profiles are derived through wavelet reconstruction. In addition, WRO is integrated with an indicator variable technique for trajectory alignment. A case study of lysine production based on a semi-batch fermentation process demonstrates the superiority of the proposed method over the conventional multiway method. *Copyright © 2007 IFAC*

Keywords: Wavelet analysis, Multivariate analysis, Batch process, Operation profile optimization, Quality improvement

## 1. INTRODUCTION

More than ever, it has become crucial for industries to improve product quality and yield in a brief period of time, as product life cycles are getting shorter and international competition is getting keener. Batch processes play a particularly important role in producing specialty and high value-added products. Since batch processes are operated by following predetermined operation profiles, operation profile optimization is a key technology to offer strong competition to others.

To improve product quality and yield, good product ratio, and to optimize operation profiles, it is necessary to relate product quality with operating conditions. In the last decade or so, statistical approaches for building models and optimizing operating conditions have been investigated (Jaeckle and MacGregor, 1998; Kano et al., 2004). In most statistical approaches, partial least squares (PLS) or principal component regression (PCR) has been used. However, it is difficult in general to build a good model of a batch process through conventional multivariate analysis methods. Although multiway methods such as multiway principal component analysis (MPCA) and multiway PLS (MPLS) are available to build batch process models (Nomikos and MacGregor, 1994, 1995), such methods tremendously increase the number of input variables to cope with non-stationary operation profiles. The explosion of input variables would deteriorate the prediction performance of the derived models.

In the present research, to realize operation profile optimization of batch processes, a new method is proposed by integrating wavelet analysis and multivariate analysis. The usefulness of the

proposed method is demonstrated through a case study of lysine production based on a semi-batch fermentation process.

## 2. MODELING THROUGH MULTIVARIATE ANALYSIS

Multivariate analysis has played a very important role in process control and process monitoring. For example, in process control, multivariate analysis methods such as PLS have been widely used to develop softsensors that estimate key variables such as product quality from on-line measured process variables (Mejdell and Skogestad, 1991; Kano et al., 2000). On the other hand, in process monitoring, multivariate statistical process control (MSPC) based on PCA has been widely used (Jackson, 1959; Jackson and Mudholkar, 1979; Kourti and MacGregor, 1995; Kano et al., 2002).

Such multivariate analysis methods are suitable for two-dimensional data and cannot be applied to operation data obtained from a batch process, because a batch process is non-stationary and its operation data usually form a three-dimensional array (batches×variables×time). To apply multivariate analysis to batch process data, multiway methods such as multiway PCA (MPCA) and multiway PLS (MPLS) have been developed (Nomikos and MacGregor, 1994; Nomikos and MacGregor, 1995). In multiway methods, three-dimensional arrays are unfolded to two-dimensional arrays. For example, a three-dimensional array (batches×variables×time) can be unfolded to (batches×(variables×time)) when variables at different sampling points are regarded as different variables. In this type of approach, operation data in the $n$th batch are described by

$$\boldsymbol{x}_n = [x_{1,1} \cdots x_{1,T} \cdots x_{M,1} \cdots x_{M,T}]^T \quad (1)$$

where $x_{m,t}$ $(m = 1, 2, \cdots, M; \ t = 1, 2, \cdots, T)$ is the $m$th variable at the $t$th sampling point. Once batch process data are unfolded, conventional methods such as PCA and PLS can be applied to the data $\boldsymbol{X}$ consisting of $\boldsymbol{x}_n$ $(n = 1, 2, \cdots, N)$.

As clearly shown in Eq. (1), the number of input variables after unfolding becomes $T$ times larger than that of original input variables. The increase of input variables generally causes deterioration in the reliability or the estimation accuracy of the statistical model. One practical approach for solving this problem is to thin out data at a part of sampling points and to make the ratio of the number of samples, i.e., batches, to that of input variables higher. For example, Chu et al. (2004) proposed a bootstrapping-based variable selection method for the improvement of the quality estimation performance in a batch process. However, not only the estimation accuracy but also the operation profile reconstruction performance is critical to optimize operating conditions in a batch process.

## 3. WAVELET REGRESSION AND OPTIMIZATION (WRO)

In the present work, to realize operation profile optimization of batch processes, a new method based on wavelet analysis and multivariate analysis is proposed. Since the proposed method uses selected wavelet coefficients of operation profiles as input variables, it is referred to as wavelet regression and optimization (WRO). In particular, its modeling part is referred to as Wavelet Regression (WR).

### 3.1 Procedure

Consider a batch process that has $I$ input variables with non-stationary operation profiles, $L$ input variables with constant values, and $Q$ output variables including product quality. These variables are described as $\boldsymbol{u}_i$ $(i = 1, 2, \cdots, I)$, $\boldsymbol{s} \in \Re^L$, and $\boldsymbol{y} \in \Re^Q$, respectively.

First, wavelet decomposition of level $J_i$ is applied to operation profiles $\boldsymbol{u}_i$, and wavelet coefficients $\boldsymbol{a}_{J_i,i}, \boldsymbol{d}_{j,i}$ $(i = 1, 2, \cdots, I; \ j = 1, 2, \cdots, J_i)$ are calculated. Here, $\boldsymbol{a}$ and $\boldsymbol{d}$ are approximation coefficients and detail coefficients, respectively.

The wavelet coefficients are arranged in the form of column vectors

$$\boldsymbol{c}_i = \begin{bmatrix} \boldsymbol{a}_{J_i,i}^T & \boldsymbol{d}_{1,i}^T & \cdots & \boldsymbol{d}_{J_i,i}^T \end{bmatrix}^T \quad (2)$$

$$\boldsymbol{c} = \begin{bmatrix} \boldsymbol{c}_1^T & \boldsymbol{c}_2^T & \cdots & \boldsymbol{c}_I^T \end{bmatrix}^T \quad . \quad (3)$$

In addition, matrices $\boldsymbol{S} \in \Re^{N \times L}$, $\boldsymbol{C} \in \Re^{N \times P}$, and $\boldsymbol{Y} \in \Re^{N \times Q}$ are generated from the vectors $\boldsymbol{s}$, $\boldsymbol{c}$, and $\boldsymbol{y}$, respectively. Here, $N$ denotes the number of batches.

Second, unimportant columns (wavelet coefficients) are removed from the wavelet coefficient matrix $\boldsymbol{C}$, and the resultant matrix is denoted by $\boldsymbol{C} \in \Re^{N \times M}$ again. For example, several wavelet coefficients of $\boldsymbol{C}$ are judged to be important when they are useful for both reconstructing original signals and estimating outputs. The other wavelet coefficients of $\boldsymbol{C}$ are judged to be unimportant and removed.

It is crucial to select columns of $\boldsymbol{C}$ by taking into account both the estimation performance and the reconstruction performance, which might

be conflicting. Small wavelet coefficients are unimportant for reconstruction in general, but it might be useful for estimation. In the present work, the selection of wavelet coefficients is conducted within a two-objective optimization framework. An optimal solution, which can achieve a satisfactory compromise, is determined from Pareto optimal points that are derived by changing the wavelet decomposition level and the combination of wavelet coefficients. For efficiency, the importance of each column is ranked according to the sum of absolute value (elements of the column). In the selection procedure, higher priority is assigned to the column with larger sum. In addition, all approximation coefficients are selected in advance because approximation coefficients are necessary for reconstruction.

Finally, a statistical model is built by using an arbitrary modeling method after columns of the input data matrices $S$ and $C$ and the output data matrix $Y$ are mean-centered and scaled if necessary. The model is described by

$$ y = K^T \begin{bmatrix} s \\ c \end{bmatrix} + e \qquad (4) $$

if linear regression is used for modeling. Here, $K \in \Re^{(L+M)\times Q}$ is a regression coefficient matrix and $e \in \Re^Q$ is an error vector. In general, linear regression methods such as PLS and PCR, which can cope with collinearity, can realize satisfactory results. However, nonlinear modeling method can be used for WR if linear regression does not function effectively.

### 3.2 Trajectory Alignment

An assumption of WR as well as multiway methods is that all batches have equal duration. However, there are many situations in which the duration of batches are not the same. The differences may be caused by seasonal changes in environmental variables such as coolant temperature, batch-to-batch variations in quality and impurity concentrations of raw materials, arbitrary termination of batches by plant operators, and so on. To handle the problem of unequal duration, several techniques for synchronization or alignment of batch trajectories have been proposed.

One approach to aligning batch trajectories is dynamic time warping (DTW). DTW is a dynamic-programming-based pattern matching scheme and originated from the area of speech recognition. It works with pairs of patterns and is able to locally translate, compress, and expand the patterns so that similar features within the patterns are matched. DTW was used to detect

the onset of different growth phases or failures in a batch fermentation process (Gollmer and Postens, 1996) and to analyze and monitor a batch polymerization reactor (Kassidas et al., 1998). It was also suggested by Kassidas et al. (1998) to include the total time as an extra input variable. DTW has been proven to be useful for batch process monitoring. However, it is not suitable for modeling and optimization of batch process operation, because DTW aims to reconcile the timing differences among trajectories and is not able to capture the influence of the timing differences on the final product quality.

Another approach to aligning batch trajectories is the indicator variable technique (Nomikos and MacGregor, 1994). In this technique, a variable is selected to indicate the progress of the batch instead of time. This variable, referred to as indicator variable, should be monotonically increasing or decreasing in time and have the same starting and ending values for each batch. Any variable that changes monotonically during the batch can be used as an indicator variable, either directly measured variable or calculated from other measured variables. Then, the trajectories are plotted with respect to the indicator variable instead of time. Usually, a constant increment is selected along the indicator variable to indicate the progress of the batch. Synchronization or alignment is performed by retaining the points in the trajectories that corresponds to the incremental points of the indicator variable. In other words, all trajectories are resampled by interpolation techniques with respect to the indicator variable. The indicator variable technique assumes that such a variable exists and that process knowledge can be used to determine it. However, there may not exist a single indicator variable. In such a situation, an alternative is to develop a model between the process measurements and the local time stamps of each batch (Undey et al., 2003). The issues of trajectory alignment, especially focusing on the influence of interpolation, was discussed by Kourti (2003).

### 3.3 Profile Optimization

The authors have developed the DDQI (Data-Driven Quality Improvement) method and applied it to industrial processes (Kano et al., 2004; Kano et al., 2005). In this section, a method for improving product quality and optimizing operation profiles in a batch process is described. Although statistical models derived through WR could be linear or nonlinear, the procedure based on linear models is explained here for simplicity.

A statistical model of a batch process can be built through WR, and it is described as

$$\boldsymbol{y} = \boldsymbol{K}^T \boldsymbol{t} = \boldsymbol{K}^T \boldsymbol{V}_R^T \boldsymbol{z} \qquad (5)$$

where $\boldsymbol{z} = [\boldsymbol{s}^T \; \boldsymbol{c}^T]^T$. Each variable in $\boldsymbol{y}$ and $\boldsymbol{z}$ is assumed to be autoscaled. In addition, $\boldsymbol{K}$ denotes a regression coefficient matrix of PCR, $\boldsymbol{t}$ principal component scores, $\boldsymbol{V}_R$ a loading matrix, and $R$ is the number of principal components retained in the PCR model.

Given the desired product quality $\tilde{\boldsymbol{y}}$, the scores $\tilde{\boldsymbol{t}}$ that can achieve $\tilde{\boldsymbol{y}}$ are described by the following equation when $R \geq Q$.

$$\tilde{\boldsymbol{t}} = (\boldsymbol{K}^T)^+ \tilde{\boldsymbol{y}} + \text{null}(\boldsymbol{K}^T) \qquad (6)$$

where $\boldsymbol{A}^+$ denotes the pseudo-inverse matrix of $\boldsymbol{A}$, null$(\boldsymbol{A})$ the kernel of $\boldsymbol{A}$, and dim(null$(\boldsymbol{K}^T)$) = $R - Q$. The solution $\tilde{\boldsymbol{t}}$ cannot be determined uniquely from Eq. (6), but it can be optimized under an objective function and constraints. In DDQI, the objective function is optimized under the following three constraints: 1) the desired product quality is achieved, 2) the operating condition exists in the space spanned by principal components and also in the region where data exist, and 3) all operating condition variables exist within their upper and lower bounds.

If there is no solution that satisfies all constraints, i.e., the imposed specifications on quality are too severe, the operating condition that achieves as desired quality as possible should be determined under the second and the third constraints.

After the optimal solution $\tilde{\boldsymbol{t}}$ is derived, it is projected back onto the space spanned by $\boldsymbol{z}$.

$$\tilde{\boldsymbol{z}} = \boldsymbol{V}_R \tilde{\boldsymbol{t}} \qquad (7)$$

The derived $\tilde{\boldsymbol{z}}$ is inversely autoscaled, and consequently $\tilde{\boldsymbol{s}}$ and $\tilde{\boldsymbol{c}}$ are derived. Since $P - M$ wavelet coefficients of $\tilde{\boldsymbol{c}}$ were removed when the WR model was built, such wavelet coefficients are complemented with zero. The wavelet coefficients $\tilde{\boldsymbol{c}}_i$ representing the $i$th operation profile are extracted from the reconstructed vector $\tilde{\boldsymbol{c}}$. Finally, the optimal operation profile $\tilde{\boldsymbol{u}}_i$ can be derived by applying inverse wavelet transformation to $\tilde{\boldsymbol{c}}_i$.

## 4. CASE STUDY

The usefulness of WRO is demonstrated through a case study of lysine production based on a semi-batch fermentation process (Ohno and Nakanishi, 1978).

Lysine is an essential and economically important amino acid, which is used as food and feed supplements. It also has some pharmaceutical applications in the formulation of diets with balanced amino acid composition. The production of Lysine is a complex process that makes stringent demands on the sterility of both the plant and the raw ingredients. Since Lysine is produced by bacteriological fermentation, the environment where the fermentation takes place must be precisely controlled. In general, the quantity of the batch product is extremely large with the fermentation vessel having a capacity of hundreds m$^3$. In this process, substrate flow rate is manipulated by following the predefined profile.

### 4.1 Equal Batch Lengths

First, the lysine production process is operated with various operation profiles of substrate flow rate, and lysine production $y$ at the batch end $T_f = 40$ h is measured. The substrate flow rate is measured every 20 minutes, and its time-series is denoted by $\boldsymbol{u}$. WR and MPCR are used to build linear regression models. In WR, wavelet decomposition of level $J = 5$ is applied to the operation profiles $\boldsymbol{u}$ by using Daubechies wavelet ($N = 8$), and wavelet coefficients are calculated. In this case study, only approximation coefficients are selected as input variables. The dimension of inputs, i.e., the number of selected approximation coefficients, is only 18, while the dimension of the operation profile $\boldsymbol{u}$ is 121. After wavelet decomposition and input variable selection, a PCR model is built with five principal components. Through the above procedure, the input dimension is reduced from 121 to 5. In MPCR, five principal components are retained in the PCR models. Consequently, the dimension of the input is reduced from 121 to 5 in both methods. Data from 10 batches are used for modeling, and data from another 10 batches are used for validation. The results of modeling and validation via MPCR are shown in Fig. 1, and those via WR are shown in Fig. 2. The model construction results clearly show that WR is greatly superior to MPCR in prediction performance.

On the basis of the developed WR model, operation profiles that can achieve the desired production $\tilde{y} = 20, 25$, and 30 are derived. The operation profile cannot be uniquely determined because the number of input variables is larger than that of the output variable. To get operation profiles, minimum norm solutions, which satisfy null$(\boldsymbol{K}^T) = \boldsymbol{0}$ in Eq. (6), are derived. The operation profiles are shown in Fig. 3. The realized production is 21.8, 26.6, and 31.5, respectively. They are close to the desired production.

Fig. 1. Model construction results by MPCR. Validation of estimation accuracy using modeling data (top) and validation data (bottom).



Fig. 2. Model construction result by WR. Validation of estimation accuracy using modeling data (top) and validation data (bottom).



Fig. 3. Derived operation profiles for realizing $\tilde{y} = 20$, 25, and 30.

In the next step, operation profile optimization is executed for minimization of substrate consumption (case 1) and for maximization of lysine production (case 2). In case 1, the constraints are 1) the desired production $\tilde{y} = 30$ is achieved, 2) the solution is interpolated, 3a) liquid holdup is below $V_{max} = 1000$ kL, and 3b) substrate flow rate is not negative. In case 2, the first constraint in case 1 is removed, and the other constraints remain. To satisfy the second



Fig. 4. Optimal operation profiles for realizing $\tilde{y} = 30$ with minimum operation cost (Case 1) and for maximizing $y$ (Case 2).

constraint, Hotelling's $T^2$ statistic must be below its threshold. The optimization results are shown in Fig. 4. In Case 1, substrate consumption is reduced from 892 to 782, and lysine production of 30.2 agrees with its desired value of 30. In Case 2, lysine production is 35.7, and liquid holdup at the batch end is its maximum $V_{max}$.

The results of this case study show the usefulness of the proposed WRO.

### 4.2 Unequal Batch Lengths

WRO can be integrated with the indicator variable technique for trajectory alignment. To demonstrate the integration, the lysine production process is operated with various batch durations, and lysine production $y$ at the batch end $T_f = 38 \sim 42$ h is measured. The other settings are the same as the previous case study.

In this case study, the indicator variable technique is used for trajectory alignment. However, there is no candidate for the indicator variable in this lysine production process. Therefore, an artificial indicator variable is generated, which linearly increases from 0% at the batch start to 100% at the batch end. In addition, the total time is used as an extra input variable in $s$ to keep the information of batch lengths in the model, because any trajectory alignment technique looses the information of batch lengths.

First, the estimation accuracy of WR is validated and is compared with that of MPCR. The correlation coefficients between measurements and estimates of lysine production are 0.67 for MPCR and 0.98 for WR. The results clearly show that WR is greatly superior to MPCR in prediction performance.

Next, operation profile optimization is executed for minimization of substrate consumption under $\tilde{y} = 20$ (case 1) and for maximization of lysine production (case 2). The optimization results are shown in Fig. 5. The lysine production is 20.2 in

Fig. 5. Optimal operation profiles for realizing $\tilde{y} = 20$ with minimum operation cost (Case 1) and for maximizing $y$ (Case 2).

case 1 and 34.3 in case 2. The batch duration is 41.38 in case 1 and 41.44 in case 2. The results of this case study show that the integration of WRO and the indicator variable technique functions successfully.

## 5. CONCLUSIONS

A new data-driven method is proposed for modeling and optimization of batch processes by integrating wavelet analysis and multivariate analysis. The usefulness of the proposed method is demonstrated through a case study of the lysine production process.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] Chu, Y. H., Y. H. Lee, and C. H. Han (2004). Improved quality estimation and knowledge extraction in a batch process by bootstrapping-based generalized variable selection. *Industrial & Engineering Chemistry Research*, Vol. 43, pp. 2680-2690.

[2] Gollmer, K. and C. Postens (1996). Supervision of Bioprocesses Using a Dynamic Time Warping Algorithm. *Control Eng. Practice*, Vol. 4, pp. 1287–1295.

[3] Jackson, J.E. (1959). Quality Control Methods for Several Related Variables. *Technometrics*, Vol. 1, pp. 359–377.

[4] Jackson, J.E. and G.S. Mudholkar (1979). Control Procedures for Residuals Associated with Principal Component Analysis. *Technometrics*, Vol. 21, pp. 341–349.

[5] Jaeckle, C.M. and J.F. MacGregor (1998). Product Design through Multivariate Statistical Analysis of Process Data. *AIChE J.*, Vol. 44, pp. 1105–1118.

[6] Kano, M., K. Fujiwara, S. Hasebe, and H. Ohno (2004). Data-Driven Quality Improvement: Handling Qualitative Variables. *IFAC Symp. on Dynamics and Control of Process Systems (DYCOPS)*, CD-ROM, Cambridge, July 5–7.

[7] Kano, M., K. Fujiwara, S. Hasebe, and H. Ohno (2005). Product Quality Improvement Using Multivariate Data Analysis. *Preprints of the 16th IFAC World Congress*, CD-ROM, Tu-M03-TP/22, Prague, Czech Republic, Jul. 3–8.

[8] Kano, M., K. Miyazaki, S. Hasebe, and I. Hashimoto (2000). Inferential Control System of Distillation Compositions Using Dynamic Partial Least Squares regression. *J. Proc. Cont.*, Vol. 10, pp. 157–166.

[9] Kano, M., K. Nagao, H. Ohno, S. Hasebe, I. Hashimoto, R. Strauss, and B. R. Bakshi (2002). Comparison of Multivariate Statistical Process Monitoring Methods with Applications to the Eastman Challenge Problem. *Comput. Chem. Engng*, Vol. 26, pp. 161–174.

[10] Kassidas, A., J. F. MacGregor, and P.A. Taylor (1998). Synchronization of Batch Trajectories Using Dynamic Time Warping. *AIChE J.*, Vol. 44, pp. 864–875.

[11] Kourti, T. (2003). Multivariate dynamic data modeling for analysis and statistical process control of batch processes, start-ups and grade transitions, *J. Chemometrics*, Vol. 17, pp. 93–109.

[12] Kourti, T. and J. F. MacGregor (1995). Process Analysis, Monitoring and Diagnosis, Using Multivariate Projection Methods. *Chemometrics and Intelligent Laboratory Systems*, Vol. 28, pp. 3–21.

[13] Mejdell, T. and S. Skogestad (1991). Estimation of Distillation Compositions from Multiple Temperature Measurements Using Partial-Least-Squares Regression. *Ind. Eng. Chem. Res.*, Vol. 30, pp. 2543–2555.

[14] Nomikos, P. and J. F. MacGregor (1994). Monitoring Batch Processes Using Multiway Principal Component Analysis. *AIChE J.*, Vol. 40, pp. 1361–1375.

[15] Nomikos, P. and J. F. MacGregor (1995). Multiway Partial Least Squares in Monitoring Batch Processes. *Chemometrics and Intelligent Laboratory Systems*, Vol. 30, pp.97–109.

[16] Ohno, H. and E. Nakanishi (1978). Optimal Operating Mode for a Class of Fermentation. *Biotechnology & Bioengineering*, Vol. 20, pp. 625–636.

[17] Undey, C., S. Ertunc, and A. Cinar (2003). Online Batch/Fed-Batch Process Performance Monitoring, Quality Prediction, and Variable-Contribution Analysis for Diagnosis. *Ind. Eng. Chem. Res.*, Vol. 42, pp. 4645–4658.