

Improved Monitoring and Discrimination of Batch Processes using Correspondence Analysis

Shailesh R. Patel and Ravindra D. Gudi

Abstract— Correspondence analysis (CA) has recently been proposed as a superior alternative to PCA for the tasks of monitoring and early event detection [3]. Some of the key merits of CA include improved discrimination due to the inherent nonlinear scaling as well as the ability to capture and elegantly represent the joint variable-sample associations, which make it amenable to accommodate serial correlations in the data. In this paper, we propose an extension of the CA algorithm to address some of the key problems associated with monitoring of batch processes. The task of batch process supervision is fraught with problems of serial nonlinear correlations in addition to the need to accommodate variable run lengths of the batches. It is shown in the paper that the differential weighting of points in CA offers a unique advantage of representation of the row (sample) and column (variable) points on the same plot (“Bi-plot”), which reveals the prevailing association structure between them. This joint display of observations and variables facilitates the fault diagnosis step and helps to discriminate the batches based on their performance. The paper also uses the traditional DTW for batch synchronization and a simple Euclidean distance based future data filling method during on-line monitoring. The efficacy of the proposed offline and online monitoring method is validated through a simulation study of a penicillin fed-batch fermentation process.

I. INTRODUCTION

BATCH processes have vital importance in the production of low volume - high value products. Typically pharmaceutical, food, biochemical and semiconductor manufacturing are batch processes. Batch processes are finite duration process and products from batch processes are often manufactured in a series of steps. The variables have different dynamics in each stage i.e. the correlation structure prevailing between variables could change with time. This time varying correlation adds more complexity in batch monitoring over continuous process monitoring. The abnormal behavior of any stage leads to undesired quality of the final product. Monitoring the time evolution of an ongoing batch enables the detection of abnormal condition and discrimination helps to identify the cause of fault. Multiway principal component analysis (MPCA) and multiway projection to latent structures (MPLS) are multivariate statistical process control (MSPC) schemes that have been traditionally proposed for batch

process monitoring. As is well known, these MSPC schemes are representation-oriented, dimension reduction techniques that identify the reduced dimensions on the basis of maximizing the explained variance in the data. For the task of event detection and monitoring however, it has been shown that discrimination of variance is a more suitable requirement than its representation, and towards this requirement various supervised approaches to discrimination have been proposed in the case of continuous process monitoring. Supervised approaches however require labeling of class information which may not be straightforward to obtain. Towards this end, correspondence analysis has been shown to achieve discrimination due to a different approach of nonlinear scaling, and does not require supervised learning. Furthermore, it exhibits the ability to capture the dependencies between both row (samples) and columns (variables), rather than the individual dependencies in either the row or the column. This ability facilitates the capture of the joint association between samples and variables and their temporal changes (viz. serial correlation). Moreover, it also provides for a more suitable representation of these joint relationships, i.e. on a single plot called the biplot and therefore facilitates the analysis and diagnosis tasks.

Other important aspects of batch process monitoring that have received abundant attention in literature are batch unfolding and synchronization. Approaches to unfolding of batches include [7] batch-wise, variable-wise or hybrid versions of these. In terms of adherence to the key PCA assumption of statistical independence of the samples, the batch-wise unfolding is best suited. However, it has problems related to completion of batch records during online monitoring. While the variable-wise unfolding approach analyzes variable correlation and does not require the online batch record completion, it violates the assumption of statistical independence of samples and is therefore most susceptible to serial data correlation. Other important and useful approaches are based on both batch-wise and variable-wise unfolding and are reported in [6]. Approaches based on functional space approximation (FSA) [2] are oriented towards a more compact representation of the time trajectories and therefore simplify computational complexity involved during the subsequent steps of processing the unfolded matrix.

To address the problem of unequal batch durations, Nomikos and MacGregor [7] proposed the use of an indicator variable instead of time to represent the batch trajectories. Data synchronization methods have also been proposed in literature to make the batch trajectories equal.

Manuscript received September 22, 2008.

S. R. Patel is with the Honeywell Technology Solutions Lab, Bangalore - 560 076, India (e-mail: Shaileshkumar.Patel@Honeywell.com).

R. D. Gudi is with Research and Technology - Automation Lab in Honeywell Technology Solutions, Bangalore - 560 076, India & Chemical Engineering Department, Indian Institute of Technology Bombay, Powai, Mumbai - 400076, (phone: 91-22-25767204; e-mail: ravigudi@iitb.ac.in).

Kassidas et al. [5] proposed a method based on dynamic time warping (DTW) algorithm to synchronize the variable trajectories. Tomasi et al. [8] proposed the correlation optimized warping (COW) to handle the problem of unequal batch duration. The DTW based synchronization method aligns the events in time and helps to improve the control limit for Q and T^2 statistics.

This paper extends the application of correspondence analysis for discriminating and monitoring tasks in batch processes. For data unfolding, we propose a hybrid approach in which batch wise unfolding is first done in view of its statistical properties, followed by DTW for data synchronization. We next evaluate two approaches of analyzing the synchronized trajectories. In the first approach we propose an FSA on the synchronized trajectories to represent them more compactly and reduce the subsequent computation load. To aid the task of diagnosis, the resulting rows containing coefficients are then refolded back in a variable wise fashion. In the second approach, the synchronized trajectories are compared with the DTW reference trajectory and a variable-wise error vector is generated and used for further analysis. We next propose the use of Correspondence analysis on the resulting matrix with a view to represent the normal region of operation on the bi-plot and to better understand the sample variable relationships during different regimes of batch operation. Depending on the approach used, the matrix used for correspondence analysis would contain either the FSA coefficients or a single error term for each variable. Since we use a batch-wise unfolding approach first, during online monitoring, it is required to predict the future data of ongoing batch. For this we propose to use a simple Euclidean distance based approach to generate the data for remainder of the batch on the basis of its similarity with existing batch records. The proposed CA based approach for statistical model building and monitoring has been analyzed using the two approaches mentioned above and validated with data generated from simulations involving a representative antibiotic fermentation presented in literature [1].

The remainder of the paper is organized as follows. The relevant data preprocessing methods, DTW and FSA, are discussed briefly in section 2. A novel method for batch process monitoring, multiway -CA (MCA), is proposed in section 3. It also enlists the data unfolding method. A detailed algorithm for multiway CA is given in section 4. It shows that how data preprocessing step can be combined with proposed monitoring scheme along with detailed online and offline monitoring algorithm. Finally, the proposed method is applied on fed-batch penicillin fermentation process to validate the proposed monitoring scheme.

II. DATA PREPROCESSING

In this section relevant data preprocessing methods, DTW and FSA algorithm, are discussed.

A. Dynamic Time Warping (DTW)

Evolving from speech recognition literature, DTW nonlinearly warps the two trajectories by minimizing the distance using the principle of dynamic programming. Let T and R denote the multivariate trajectories of test and reference batch respectively; both are of dimensions $t \times N$ and $r \times N$ respectively, where t and r are the number of observations and N is the number of measured variables. If i and j denote the time indices of T and R respectively, then DTW finds a sequence O^* of K points on a $t \times r$ grid such that total distance measured between T and R trajectory is minimized.

$$O^* = \{p(1)p(2)\dots p(k)\dots p(K)\} \quad (1)$$

where $p(k) = [i(k), j(k)]$ is an ordered pair indicating a position on the grid. The objective function of DTW can be formulated as,

$$D(t,r) = \frac{\sum_{k=1}^K d[i(k), j(k)] \cdot w(k)}{N(w)} \quad (2)$$

where $D(t,r)$ is a normalized total distance between the two trajectories; $d = [i(k), j(k)]$ is the weighted local distance between the T and R trajectory. $N(w)$ is a normalization factor which make normalized total distance independent of the number of the path points. To solve the optimization problem, DTW imposes some constraints, which can be summarized as, (i) end point constraints, (ii) local continuity constraints, and (iii) global constraints. End point constraints put the restrictions on the first and last points. Local continuity constraints specify a set of allowable predecessor for each point. Global constraints define subsets of $t \times r$ grid to be the actual search space for the optimal path.

Many variants of DTW algorithm are presented in literature. Basically they all are classified as, (i) Symmetric DTW algorithm and (ii) Asymmetric DTW algorithm. Symmetric algorithm maps the time index i of T and j of R onto a common index ' k '. The optimal path will pass through all the points in both the trajectories. On the other hand, an asymmetric algorithm will map the time index of R on the time index of T or vice-versa. It implies that the path passes through all the points of reference trajectory, but it skips some of the points in the other, which could introduce the false alarm. For reasons discussed in Kassidas et al. [5], this paper also uses the symmetric DTW algorithm for data synchronization.

B. Functional Space Analysis

The function space analysis is based on the concept of orthonormal function approximation. The trajectories of process measurements in the batches are mapped onto the new feature parameters in the function space. In function space analysis, the trajectory of each variable in each batch run is represented as a function, $f(t)$. It is approximated in terms of an orthonormal set $\{\Phi_n\}$ of continuous function.

$$X = \begin{bmatrix} f_{1,1}(t) & f_{1,2}(t) & \cdots & f_{1,j}(t) \\ f_{2,1}(t) & f_{2,2}(t) & \cdots & f_{2,j}(t) \\ \vdots & \vdots & \vdots & \vdots \\ f_{l,1}(t) & f_{l,2}(t) & \cdots & f_{l,j}(t) \end{bmatrix} \quad (3)$$

$$f(t) \cong F_N(C_N, t) = \sum_{n=0}^{N-1} c_n \Phi_n(t) \quad (4)$$

$$C_N = \int f(t) \phi_n(t) dt \quad (5)$$

where $C_N = \{c_n\}_{n=0,1,\dots,N-1}$ are projections of $f(t)$ onto each basis function; N is the number of coefficients. The Legendre polynomial is used as an orthogonal basis function.

III. MULTIWAY CORRESPONDING ANALYSIS

Correspondence analysis is a multivariate statistical technique, which has an objective of dimensionality reduction along with discrimination. It projects the data such that the positions of the row and column points on low dimensional subspace are consistent with their associations.

A. Theory of Correspondence Analysis

The basic aim of correspondence analysis is to identify low-dimensional subspaces which best 'fit', or lie close to, a given set of point vectors. In CA, the closeness of a set of points to subspace is defined by weighted distance of given set of points to subspaces. CA assumes the data follows normal density distribution. This assumption leads that weighted Euclidean distance will follow chi-square (χ^2) statistic. CA decomposes a measure of row-column association, typically formulated as total χ^2 value, to find out low dimensional subspace. χ^2 value shows the deviation of points from the low dimensional subspace. It is testing whether a probability density conforms to some expected density. In order to test whether the deviations represent statistically significant departures from the expected frequencies, the following chi-square statistics is calculated,

$$\begin{aligned} \chi^2 &= \sum \frac{(\text{Observed Frequency} - \text{Expected Frequency})^2}{\text{Expected Frequency}} \\ &= (X - \bar{X}^T) D (X - \bar{X}^T) \\ &= n(P - \bar{P}^T) D (P - \bar{P}^T) \end{aligned} \quad (6)$$

where P and \bar{P} are the vectors of relative frequencies and of expected frequencies (average frequencies) respectively. X is row (or column) profile and \bar{X} is average row (or column) profile called as centroid of the row (or column) profile.

The solution to the problem of minimizing the weighted distances can be given by decomposition of inertia of row (or column) cloud, i.e. generalized SVD of the matrix $[(1/g)X - rc^T]$ [4]. $g = 1^T X 1$ is the grand sum of data matrix. The vectors r and c are the vectors of row and column sums of $[(1/g)X]$ respectively.

$$r = [(1/g)X]I \quad c = [(1/g)X]^T I \quad (7)$$

The total inertia of the row and column cloud can be shown to be the same (Greenacre [3]) and is given by the χ^2 value divided by g . The weight matrix D is chosen as diagonal matrix of the row sums (D_r) or the column sums (D_c). The generalized SVD of this matrix is defined as,

$$[(1/g)X - rc^T] = AD_\mu B^T \quad (8)$$

such that, $A^T D_r^{-1} A = I_{m \times m}$ and $B^T D_c^{-1} B = I_{n \times n}$

The generalized SVD of $[(1/g)X - rc^T]$ can also be realized as the regular SVD of appropriately scaled matrix X , as explained below. Define matrix P as,

$$P = D_r^{-1/2} [(1/g)X - rc^T] D_c^{-1/2} \quad (9)$$

Then the regular SVD of the matrix P gives the required singular vectors. A and B define the principal axes for the column and the row cloud respectively. The respective coordinates of the row and the column profiles with respect to their own principle axes are given by,

$$F = D_r^{-1} A D_\mu \quad \text{and} \quad G = D_c^{-1} B D_\mu \quad (10)$$

In general, major part of the inertia can be explained by retaining only first k ($k \ll m, n$) principal axes corresponding to largest singular values.

B. Data unfolding and Scaling

The archived batch data usually occur as stacked into a three way array data matrix. The data should be unfolded into two way matrix for further analysis. There are three different ways to cut the three way array data into slices and two different ways to arrange this slices, one below the other or side by side, which will give total six different types arrangements. Among these six arrangements, two unfolding methods are quite popular in the literature. Nomikos and MacGregor unfold the data in $NB \times (NV \times NK)$ where NB is number of batches, NV is number of variables and NK is samples points. It brings an advantage in terms of scaling. The scaling will remove the mean (average) trajectory (column mean) from the data. The monitoring of resultant scaled data will be equivalently checking the deviation around the mean. Wold et al. [9], unfold the data as $(NB \times NK) \times NV$. This type of unfolding will bring an advantage to find the correlation between the samples and variables. As CA offers an advantage to access the sample-variable relationship; the second unfolding is more suitable for CA analysis. But it brings the disadvantage in terms of scaling. Scaling in this case will remove the column mean that is grand mean of variable trajectory. The monitoring of scaled data will check the deviation from this grand mean, which will increase the false alarm rate.

As explained in previous section, correspondence analysis works on frequency vector, which in turns requires positive entries. The FSA algorithm, while mapping the original variable trajectory onto new feature vector, could result in non-positive features. The proposed algorithm in this paper scales the data to convert them into positive

values. To take advantage of both unfolding methods, the paper proposes a two tier unfolding method.

- 1) The three way array coefficient matrix is unfolded in the same way as the approach of Nomikos and MacGregor [7]. The data are scaled by subtracting each coefficient by its minimum value along the column and dividing by its average range.
- 2) The unfolded matrix is rearranged into a form described by Wold et al. [9], by placing the coefficients data for batch one below the other i.e. $(NB \times NC) \times NV$, where NC is number of coefficients.

C. Statistical Control charts for CA

Correspondence Analysis builds a statistical model of the process from the historical process data. To detect the fault, Detroja et al. [3] define Q and T^2 statistic analogous to PCA. The measurement vector x is converted into frequency vector by dividing by its row sum. The resultant vectors are projected onto the principal axes to obtain the row coordinates in lower dimension.

$$X = [x_1, x_2, \dots, x_j] \quad (11)$$

The row sum of this measurement vector, r is given by

$$r = \sum_{j=1}^J x_j \quad (12)$$

and the new row coordinates can be obtained as

$$f = \begin{bmatrix} \frac{1}{r} x^T G D_\mu^{-1} \\ r \end{bmatrix} \quad (13)$$

Q-statistic for CA

Any significant deviation in the direction of $n-k$ PCs (corresponding to smallest singular values), is also indicative of fault. This space can be monitored more robustly by using the Q statistic.

The Q -statistic for CA is defined as,

$$Q = \left[Bf - \left(\frac{1}{r} x - c \right) \right]^T \left[Bf - \left(\frac{1}{r} x - c \right) \right] \quad (14)$$

The T^2 -statistic for CA is defined as,

$$T^2 = f^T D_\mu^{-2} f \quad (15)$$

where D_μ contains first k -largest singular values, which were retained in CA model. The control limit for Q (and T^2) statistic is taken as $(1-\alpha)$ percentile, α level of significance, of SPE (and T^2) value for training data.

IV. MULTIWAY CA ALGORITHM

This section contains a brief description of algorithm steps for both off-line and on-line monitoring

A. Offline Monitoring

Approach -A

Step A: Data Preprocessing

Step 1: Data Scaling

The data are scaled by dividing each variable with its average range.

Step 2: Data Synchronization

The batch trajectories are synchronized using the symmetric DTW algorithm. The synchronized trajectories are mapped onto orthogonal basis function (Legendre functions). All the coefficients are tabulated in a three way array theta matrix.

Step B: Process Monitoring Scheme

Step 1: Data Unfolding

The three way theta matrix is unfolded by two tier unfolding method. Let \mathbf{X} be the resultant two way array theta matrix.

Step 2: Multi-way Correspondence Analysis

The grand sum (g) of data matrix (X) is calculated by summing all matrix elements. Each element is divided by the grand sum to convert into corresponding frequency. The masses (D_{rn}) and weights (D_c) are calculated as,

$$D_r = \text{diag}(r), \text{ and } D_c = \text{diag}(c) \quad (16)$$

The total inertia matrix of row (or column) cloud is calculated as $[(1/g)X - rc^T]$.

The inertia matrix is decomposed by generalized SVD to get the principal axes (co-ordinates) for the row and column profiles. $[(1/g)X - rc^T] = AD_\mu B^T$

$$F = D_r^{-1} AD_\mu \text{ and } G = D_c^{-1} B D_\mu$$

Only first k principal components are retained which are sufficient to explain most of the process inertia. Thresholds for both T^2 and Q statistics are calculated as per Eq. 14 & 15.

Approach -B

Approach-B differs from Approach-A in the data preprocessing step. After data scaling (step A.1) and synchronization (Step A.2), a new step is added to calculate the dissimilarity measure, A dissimilarity measure, Euclidean distance, is calculated for all batch trajectories with reference to the reference batch trajectory obtained from DTW algorithm. Thus, it summarizes the variable trajectory information into a single scalar value, distance. This would result into a two way array, which will eliminate the need of unfolding (Step B.1) and makes the subsequent procedure computationally efficient. The calculation procedure from step B.2 onwards are similar to Approach A.

Dissimilarity measure calculation

The dissimilarity measure is calculated as distance between the batch B_i and reference batch B_{ref} . The sum of squared distance for each variable is calculated over the time. $d_{B_i} = [d_{i1} \ d_{i2} \ \dots \ d_{iJ}]$ where d_{B_i} is the distance vector for batch B_i , which contains the measure of departure of batch from the reference batch. The distance vectors for all batches are tabulated as batch \times variables. Let D be the resultant distance matrix. The method assumes that the distance value for each variable follows Gaussian distribution and violation will indicate the fault.

B. Online Monitoring

The on-line implementation of the Multi-way based monitoring scheme is similar to the off-line implementation with one important difference in the on-line data case, the prediction of the future behavior of the batch trajectory up to its expected end is required. The difficulty is total length of new batch is not known. One has to predict how different the new batch is with respect to reference batch in order to fill future data. To know which point of reference trajectory best represents the progress of the new batch up to the current time, Kassidas et al (1998) proposed a method which is based on DTW. The method poses a new DTW problem at each time instant to predict the age of batch, which is computationally intensive. In our work, a new approach is proposed to fill the future data. The propose approach uses historian batch data to predict the evolution of new batch in future. The method assumes that the behavior of the process will remain in ‘normal’ state. It finds a batch from the historian data which has a similar dynamics as new batch up to current instant based on distance criteria. The total Euclidean distance between B_{NEW} and B_i , $i=1,2,\dots,I$ is calculated up to current time instant. The batch B_K with lowest distance value is selected as closet batch to B_{NEW} . The batch B_K data from the current time instant to its end is taken as future profile for the B_{NEW} . Once the batch run data are filled by above method, the monitoring steps outlined in offline monitoring section is followed to monitor the time evolution of ongoing batch. The procedure is repeated as soon as the new observation comes.

V. CASE STUDY

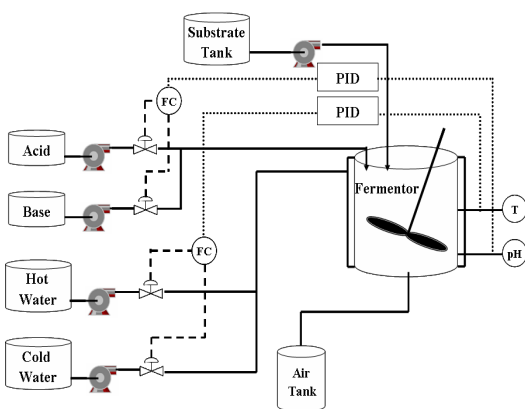


Fig. 1. Penicillin fed batch fermentation process diagram

The application of proposed monitoring scheme is illustrated here with a fed-batch penicillin fermentation process. The data are generated from PenSim v2.0 simulator (available at the web site: <http://www.chee.iit.edu/~cinar/software.html>) developed by Prof. Ali Cinar, Illinois Institute of Technology, Chicago, U.S.A. The fermentation process is shown in Figure 4. 55 sets of normal batch data are generated from PenSim v2.0 simulator by changing the initial conditions. The total duration of batch process is

varying from 330 hrs to 390 hrs. Out of total 17 variables 8 variables are used to build statistical model. Four faulty batches were simulated by introducing step and ramp change in aeration rate agitator power respectively at 100th time instant.

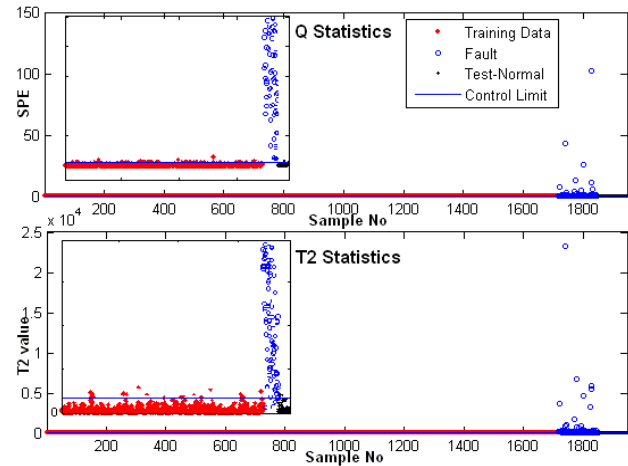


Fig. 2. Offline monitoring- Q and T^2 Statistics - Approach A

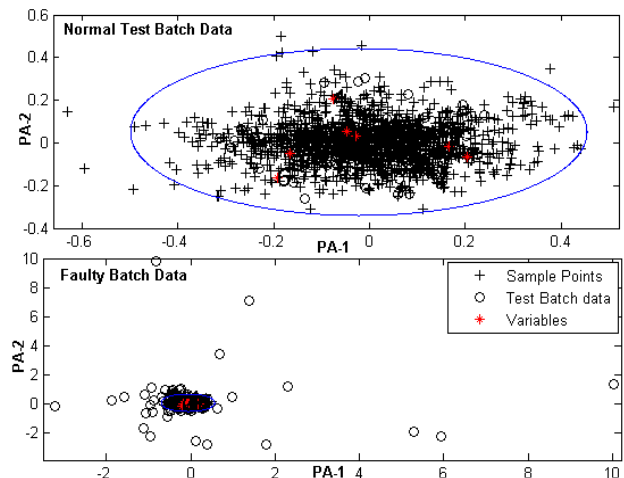


Fig. 3. Offline monitoring – Biplot - Approach A

A. Offline Monitoring Results

The results of the offline multiway-CA framework (Approach-A) on normal and faulty test batch data are shown in figures 2 & 3. The faulty batches are detected, as SPE and T^2 value for these batches crosses the control limits in both Q and T^2 charts respectively (Fig 2). The biplot presentation for approach-A is depicted in Figure 3. The faulty batch data projection (Fig 3 subplot 1) are found to fall outside the cluster formed by normal batches, where as, the normal test batch data are projected (Fig 3 subplot 2) on the cluster as expected. Figures 4 & 5 show the offline monitoring result for Approach –B. As per the expectation, all faulty batch projections violate Q as well as T^2 control limit while normal bath projection falls within limits.

B. Online Monitoring Results

The online monitoring result of a faulty batch by CA approach is depicted in Figure 6 & 7. The fault is introduced

in variable 1 at 100th time instant. The fault is detected on the Q and T^2 plot at 110th time instant as it crossed the control limit (Fig 6 & 7 subplot 1). The biplot presentation of new batch can be seen in the subplot (Fig 6 & 7 subplot 2). As the time progress, the projection of the new batch (marked with 'encircled pink +') falls outside the normal cluster (Fig 7). The angle between variable-1 (marked with '1') and new batch projection is close to zero and remains the same as time progress. The projection of new batch data moves in the direction shown by the black arrow. The biplot analysis suggests that the cause of fault is variable 1, which is in line with expectation.

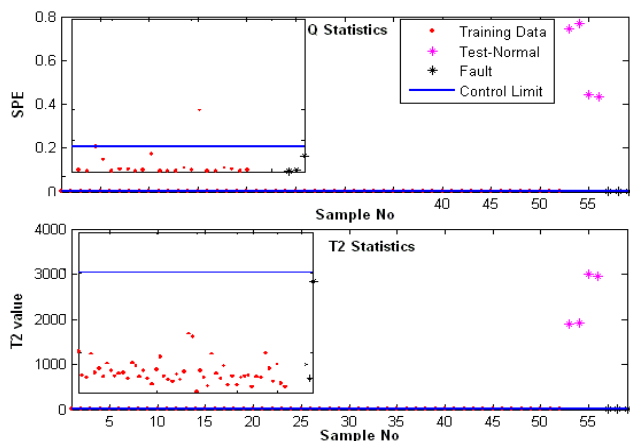


Fig. 4. Offline monitoring- Q and T^2 Statistics - Approach B

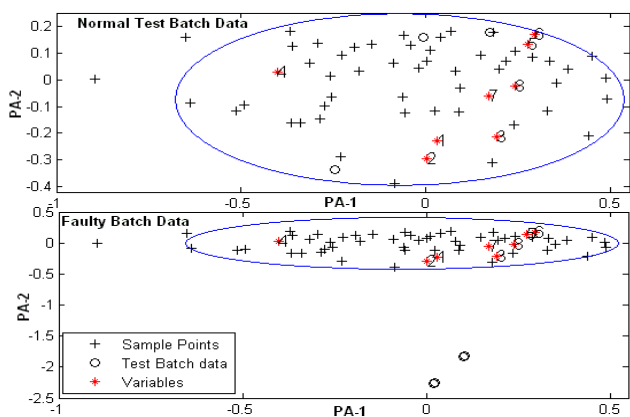


Fig. 5. Offline monitoring – Biplot - Approach B

VI. CONCLUSION

A new batch process monitoring scheme based on CA is proposed for fault detection and diagnosis. The ability of CA to capture time varying correlations makes it attractive for use in batch process monitoring. Both off-line and online monitoring schemes have been proposed and validated using simulations involving fed-batch fermentation. Analysis related to the interpretations of the bi-plot and comparison with PCA is currently in progress. The use of the multiway CA algorithm for batch quality predictions is a topic of future study.

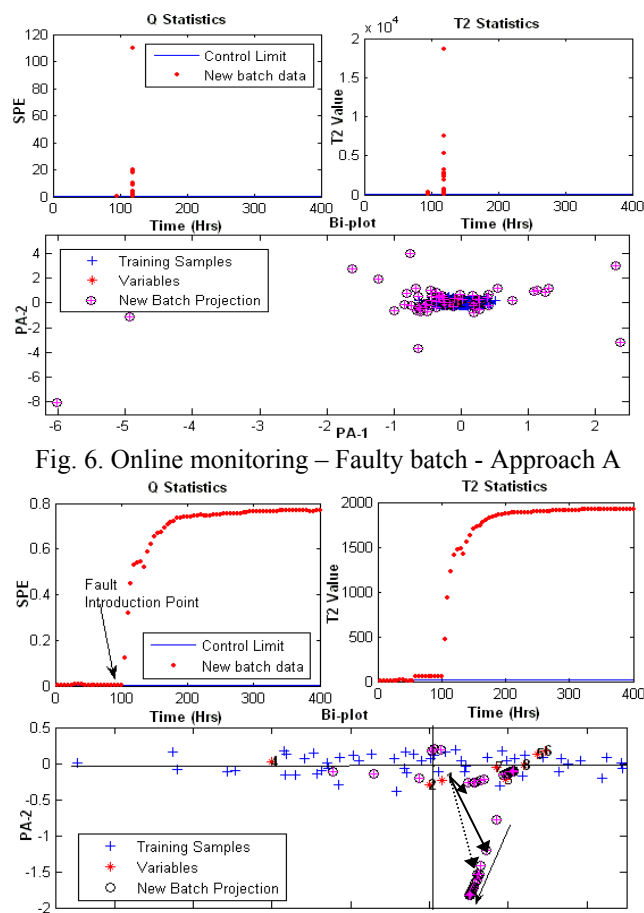


Fig. 6. Online monitoring – Faulty batch - Approach A

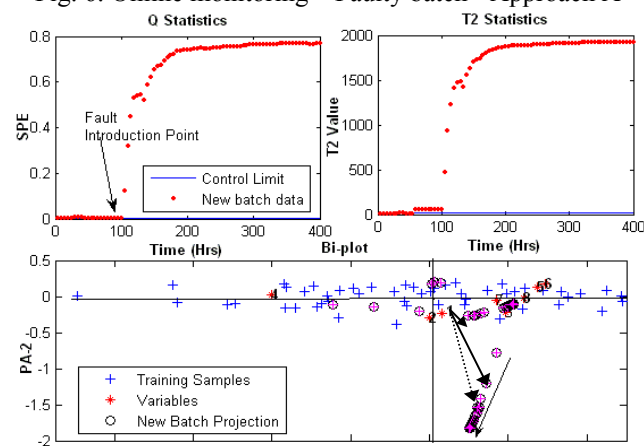


Fig. 7. Online monitoring – Faulty batch - Approach-B

REFERENCES

- [1] Birol, G., Ündey, C. and Çinar, Ali, "A modular simulation package for fed-batch fermentation: penicillin production", *Comput. Chem. Eng.*, 26, 1553-1565, 2002
- [2] Chen, J. and J. Liu, "Derivation of Function Space Analysis Based PCA Control Charts for Batch Process Monitoring", *Chemical Engineering Science*, 56, 3289-3304, 2001
- [3] Detroja K. P., R. D. Gudi, S. C. Patwardhan and K. Roy, "Fault Detection and Isolation Using Correspondence Analysis," *Ind. Eng. Chem. Res.*, 45, 223-235, 2005.
- [4] Greenacre, M. J., "Theory and Applications of correspondence Analysis," *Academic Press Inc.*, Orlando, Florida, 1984.
- [5] Kassidas, A., J. F. Macgregor and P. A. Taylor, "Synchronization of Batch Trajectories using Dynamic Time Warping," *AIChE Journal*, 44(4), 864-875, 1998
- [6] Lee, J. M., C. K. Yoo, and I. B. Lee, "Enhanced process monitoring of fed-batch penicillin cultivation using time-varying and multivariate statistical analysis" *Journal of Biotechnology*, 44(4), 864-875, 1998
- [7] Nomikos, P. and J. F. MacGregor, "Monitoring Batch Processes using Multiway Principal Component Analysis," *AIChE Journal*, 110(2), 119-136, 2004.
- [8] Tomasi G., Frans van den Berg and C. Anderson, "Correlation optimized warping and dynamic time warping as preprocessing methods for chromatographic data," *Journal of Chemometrics*, 26, 231-241, 2004
- [9] Wold, S., N. Kettaneh, H. Friden and A. Holmberg, "Modelling and Diagnostics of Batch Processes and Analogous Kinetic Experiments", *Chemometrics Intell. Lab. Syst.*, 44, 331-340, 1998