

# Canonical Variable Analysis for Fault Detection, System Identification and Performance Estimation

Xiaochuan Li <sup>1</sup>✉

Email [lix29@lsbu.ac.uk](mailto:lix29@lsbu.ac.uk)

Fang Duan <sup>1</sup>

Email [duanf@lsbu.ac.uk](mailto:duanf@lsbu.ac.uk)

Tariq Sattar <sup>1</sup>

Email [sattartp@lsbu.ac.uk](mailto:sattartp@lsbu.ac.uk)

Ian Bennett <sup>2</sup>

Email [Ian.bennett@shell.com](mailto:Ian.bennett@shell.com)

David Mba <sup>1</sup>

Email [mbad@lsbu.ac.uk](mailto:mbad@lsbu.ac.uk)

<sup>1</sup> School of Engineering, London South Bank University, London, UK

<sup>2</sup> Shell Global Solutions International B.V., Rijswijk, Netherlands

## Abstract

Condition monitoring of industrial processes can minimize downtime and maintenance costs while enhancing the safety of operation of plants and increasing the quality of products. Multivariate statistical methods are widely used for condition monitoring in industrial plants due to the rapid growth and advancement in data acquisition technology. However, the effectiveness of these methodologies in real industrial processes has not been fully investigated. This paper proposes a CVA-based approach for process fault identification, system modeling and performance estimation. The effectiveness of the proposed method was tested using data acquired from an operational industrial centrifugal compressor. The results indicate that CVA can be effectively used to identify abnormal operating conditions and predict performance degradation after the appearance of faults.

AQ1

## Keywords

Condition monitoring  
Canonical variable analysis  
Fault detection  
System identification  
Performance estimation

## 1. Introduction

Modern industrial facilities such as natural-gas processing plants are becoming increasingly complex and large-scale due to the use of machines of different nature. The complexity of large-scale industrial facilities makes it difficult to build first-principle dynamic models for condition monitoring (Russell et al. 2000). Thus, existing condition monitoring approaches for industrial processes are typically derived from routinely monitored system operating data. Due to the advancement in instrument and automation technology, long-term and high-frequency measurements can be taken with the different sensors mounted on the machinery

systems. The monitored data are easily stored and explored to extract important process condition information. Many methodologies have been developed to combine the multivariate process data for analysis, such as state-space-based models (Negiz and Çinar 1997a), time series analysis (Negiz and Çinar 1997b), and dimensionality reduction techniques (Chiang et al. 2000; Yang et al. 2012; Komulainen et al. 2004; Ku et al. 1995).

The main advantage of dimensionality reduction techniques over traditional approaches is that they can take into consideration the possible correlation between the different measured variables, hence facilitating fault detection and system identification (Ruiz-Carcel et al. 2016). Two early examples of process monitoring techniques are principal component analysis (PCA) (Ku et al. 1995) and partial least-squares (PLS) analysis (Muradore and Fiorini 2012). Both PCA and PLS assume that the monitored variables are time-independent (i.e., the observations at one time instant are not correlated with those in the past time instants). This assumption might not hold true for real industrial processes (especially chemical and petrochemical processes) because measurements driven by noises and disturbances often show strong correlation between the past and future time instances (Odiwei and Cao 2010). Dynamic extensions of PCA and PLS, so-called dynamic PCA and dynamic PLS, have been proposed to solve this problem, making them more suitable for dynamic processes monitoring. Although DPCA and DPLS have been successfully applied to dynamic systems, they have been reported not to be able to fully capture some important dynamic behaviors of the system working under varying operating conditions (Jiang et al. 2015a; Ruiz-Carcel et al. 2016).

Aside from approaches derived from PCA and PLS, the canonical variable analysis (CVA) is also a multivariate monitoring tool. CVA is a state-space-based method which takes both serial correlations and relationship between correlated variables into account, hence is more suitable for dynamic process modeling (Odiwei and Cao 2010). The performance of CVA has been tested by several researchers using computer-simulated data (Jiang et al. 2015b; Huang et al. 2015) and data obtained from small-scale test rigs (Ruiz Cárceles and Mba 2014). However, the effectiveness of CVA in real complex industrial processes has not been fully studied. In this investigation, we propose a CVA-based method for abnormal behavior detection, system identification, and performance estimation of petrochemical process. To prove the validity of the method, it was tested using process data acquired from an industrial centrifugal compressor operating in the real world. The results indicate that it is possible to perform fault detection and prognosis using real-life data.

## 2. Methodology

### 2.1. CVA for Fault Detection

CVA is a dimension reduction technique to monitor the process by converting the multidimensional observed data into a health indicator. Process data acquired from the system operating under normal operating conditions are used to determine the threshold for normal operating limits. The process faults can be identified when the value of the health indicator exceeds the threshold.

The objective of CVA is to maximize the correlation between two sets of variables (Russel et al. 2000). For this purpose, the measurement vector  $y_k \in \mathcal{R}^m$  (measurement at each time point containing  $m$  variables) is expanded at each time point  $k$  by considering  $p$  past measurements and  $f$  future measurements to give the past and future vectors  $y_{p,k} \in \mathcal{R}^{mp}$  and  $y_{f,k} \in \mathcal{R}^{mf}$ .

$$y_{p,k} = [y_{k-1}^T y_{k-2}^T \cdots y_{k-p}^T]^T \in \mathcal{R}^{mp} \quad 1$$

$$y_{f,k} = [y_k^T y_{k+1}^T \cdots y_{k+f-1}^T]^T \in \mathcal{R}^{mf} \quad 2$$

To avoid the domination of variables with larger absolute values, the past and future vectors are normalized to zero means to get  $\widehat{y}_{p,k}$  and  $\widehat{y}_{f,k}$ . The vectors calculated at different time points are arranged in columns to produce past and future truncated Hankel matrices  $\widehat{Y}_p$  and  $\widehat{Y}_f$ :

$$Y_p = [\widehat{y_{p,p+1}}, \widehat{y_{p,p+2}}, \dots, \widehat{y_{p,p+M}}] \in \mathfrak{R}^{mp \times M} \quad 3$$

$$Y_f = [\widehat{y_{f,p+1}}, \widehat{y_{f,p+2}}, \dots, \widehat{y_{f,p+M}}] \in \mathfrak{R}^{mf \times M} \quad 4$$

where  $M = n - f - p + 1$ ,  $n$  represents the total number of observations for  $y_k$ .

To find the linear combinations that maximize the correlation between the two sets of variables, the Hankel matrix  $H$  can be decomposed using Singular Value Decomposition (SVD):

$$H = \sum_{ff}^{-1/2} \sum_{pf}^{-1/2} \sum_{pp}^{-1/2} = U \sum V^T \quad 5$$

where  $\Sigma_{pp}$ ,  $\Sigma_{ff}$  and  $\Sigma_{pf}$  represent the sample-based covariance and cross-covariance matrix of matrices  $Y_p$  and  $Y_f$ . The  $mp$ -dimensional past vector  $Y_p$  can be converted into the  $r$ -dimensional canonical variates  $z$  by:

$$z = J \cdot Y_p \quad 6$$

where  $J$  represents the transformation matrix, and  $J = V_r^T \Sigma_{pp}^{-1/2}$ . The truncated matrix  $V_r \in \mathfrak{R}^{r \times M}$  can be obtained by selecting the first  $r$  columns of  $V$  having the highest pairwise correlation with those of  $U$  (Samuel and Cao 2015). Then, the Hotelling health indicator can be calculated as:

$$T_k^2 = \sum_{i=1}^r z_{k,i}^2 \quad 7$$

Since the Gaussian distribution does not hold true for nonlinear processes, the normal operating limits are derived from the actual probability density function of the indicator using Kernel Density Estimation (KDE) (Odiowei and Cao 2010). Faults will be considered every time the health indicator exceeds the threshold.

## 2.2. CVA for System Identification and Performance Estimation

CVA can be used to build a state-space model which describes the dynamic behavior of the system using process data. Given the past of the measured system inputs  $u_k$  and measured outputs  $y_k$ , the following state-space model can be built:

$$x_{k+1} = Ax_k + Bu_k + w_k \quad 8$$

$$y_k = Cx_k + Du_k + Ew_k + v_k \quad 9$$

where  $x_k$  is a  $r$ -order state vector,  $w_k$  and  $v_k$  are independent white noise, and  $A, B, C, D$  and  $E$  are coefficient matrices. According to (Larimore 1990), when the order of the system  $r$  is equal or greater than the actual order of the system, the state estimate  $z = J \cdot Y_p$  can be used ~~in-replace-of~~ to replace  $x_k$ . The unknown coefficient matrices  $A, B, C$  and  $D$  can then be estimated via multivariate regression (Larimore 1990):

$$\begin{bmatrix} \widehat{A} & \widehat{B} \\ \widehat{C} & \widehat{D} \end{bmatrix} = \text{Cov} \left[ \begin{pmatrix} z_{k+1} \\ y_k \end{pmatrix}, \begin{pmatrix} z_k \\ u_k \end{pmatrix} \right] \cdot \text{Cov}^{-1} \left[ \begin{pmatrix} z_k \\ u_k \end{pmatrix}, \begin{pmatrix} z_k \\ u_k \end{pmatrix} \right] \quad 10$$

The procedure of system identification and performance estimation using the model described above can be summarized as follows:

1. Determine the input (manipulated) and output (measured performance variables) variables;
2. Collect data for past system input  $u_{p,k}$  and output  $y_{p,k}$ ;
3. Estimate future system input  $u_{f,k}$  by looking at production plan or estimating from the past input values;
4. Determine the number of states and number of past and future lags for the collected data;
5. Estimate the parameters in the state-space model using Eq. (10);
6. Predict future output  $y_{f,k}$  as per Eqs. (8) and (9);
7. Validate the proposed model by looking at the average prediction error for each one of the measured variables.

This procedure will allow operators to access how the system will behave for the specified system inputs.

### 3. Application to Centrifugal Compressor Data

CVA has been successfully used to perform condition monitoring using computer-simulated data (Lee and Lee 2008; Juricek et al. 2001) and data acquired from small test rigs (Ruiz-carcel et al. 2016). In this investigation, the capabilities of CVA for fault detection and system identification were tested using data captured from an operational industrial centrifugal compressor.

AQ3

Centrifugal compressors are widely used in oil and gas industry for gas transport, gas lift, and gas injection. They are typically operating under high pressure and high load conditions, and are therefore subject to performance degradation. The compressor used in this study is automated using a condition monitoring system, where the signals from different sensors can be visualized. A total of 50 variables including three process inputs (i.e., rotational speed, inlet temperature, and inlet pressure) and 47 performance variables were recorded, sampling at a 1-hour interval. The recorded data consist of 25,900 observations (i.e., the total monitoring time is more than 3 years in length).

In order to fully capture the dynamic characteristics of the compressor under various operating conditions, nine periods of data were used to train the CVA algorithm to obtain the normal operating limits of  $T_k^2$ , and the training data sets were intentionally chosen to cover various operating speeds and ambient temperatures. The number of time lags ( $p$  and  $f$ ) is determined by computing the autocorrelation function of the summed squares of all measurements (Odiwei and Cao 2010). The autocorrelation function indicates how long the signal is correlated with itself, and thus can be used to determine the maximum number of significant lags. In this investigation, the number of  $p$  and  $f$  was set to 15. The optimal number of dimensions retained  $r$  is determined by considering the dominant singular values in the matrix  $D$  (Negiz and Çinar 1998). After several tests,  $r = 25$  was finally adopted to represent the order of the system for the purpose of fault detection. The 99% confidence interval was employed in this study to minimize the false alarm rate.

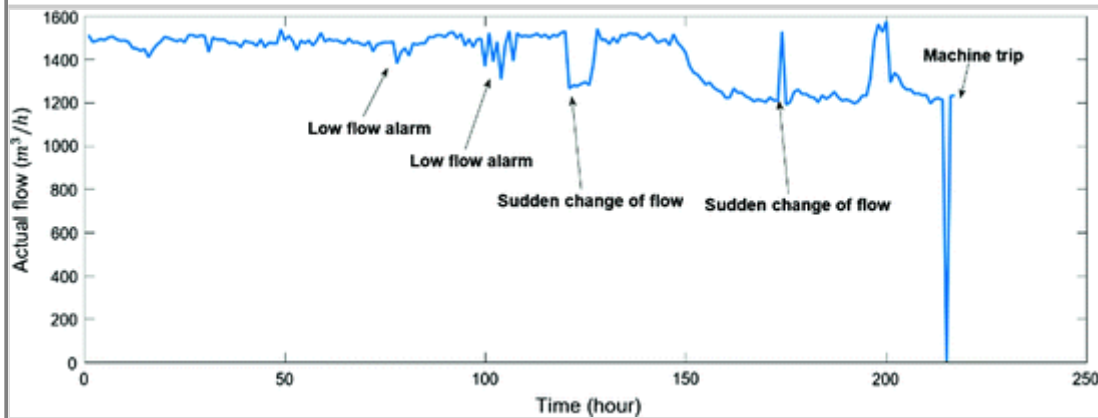
#### 3.1. Results Obtained for Fault Detection

As mentioned previously, nine data sets were used to train the CVA algorithm for fault detection. In addition to the training data sets, a period of data (hereafter referred to as data set C1) was obtained from the machine operating under faulty conditions for testing the trained algorithm. The fault evolution of the testing data can be seen in Fig. 1. According to the event logs provided by the machine operator, the monitoring system gave in total five warnings (as shown in Fig. 1) during this period of time. The first warning happened in the 78th

sample, and the machine continued to operate until the 215th sample and was then forced to completely shut down at the 217th sample. Figure 2 demonstrates the results in terms of fault detection. The fault was detected by the  $T^2$  indicator at sample 140 after several short alarms.

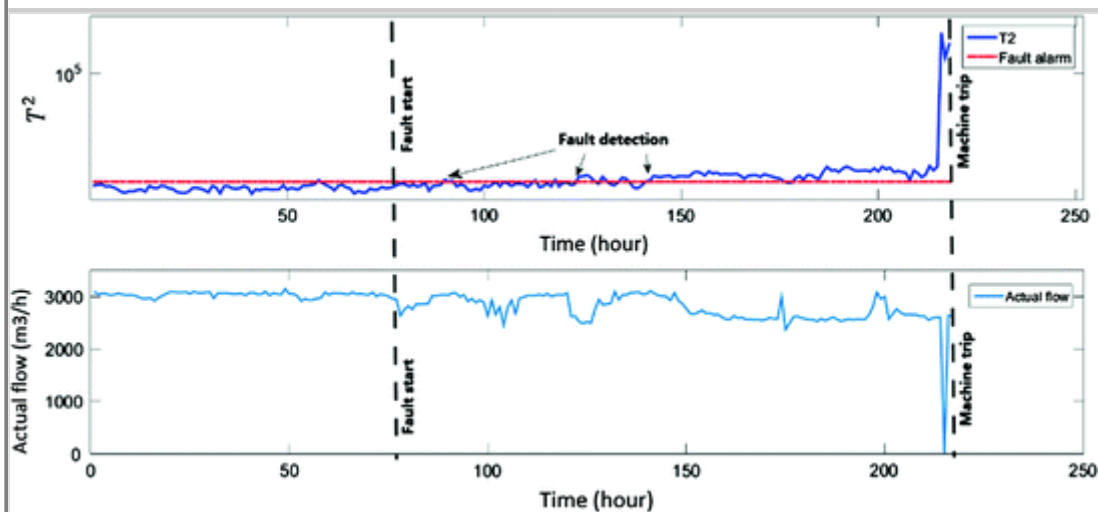
**Fig. 1**

Actual flow for data set C1



**Fig. 2**

$T^2$  for data set C1



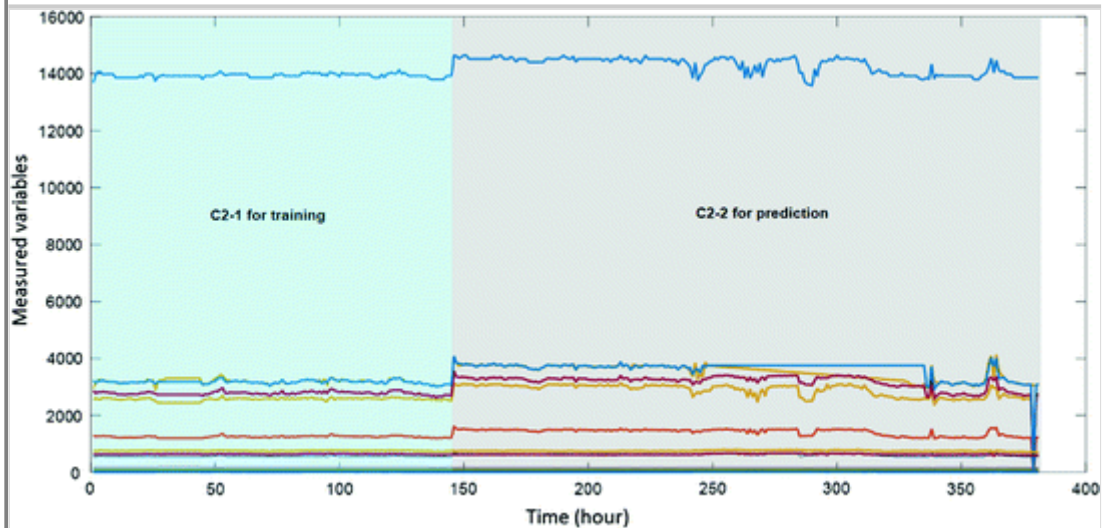
### 3.2. Results Obtained for System Identification and Performance Estimation

In addition to fault detection, plant operators may be more interested in how the system will behave given the future system input conditions and how the estimated behavior will affect the plant operation. The future system inputs can be obtained by looking at the production schedule or estimating from the past inputs. The processes described below are for prediction of future system behavior for the specified inputs.

In order to build a dynamic model as described in Eqs. (8) and (9), it is necessary to first determine the canonical variates  $z$ . Similar to the procedure described in Sect. 3, nine training data sets were used to train the CVA model to obtain  $z$ . Second, the past system inputs  $x_k$  and outputs  $y_k$  were obtained from the first 145 samples of the data set C2 (hereafter referred to as C2-1), as shown in Fig. 3. They were then used to construct a dynamic state-space model using the procedures described in Sect. 2.2.

**Fig. 3**

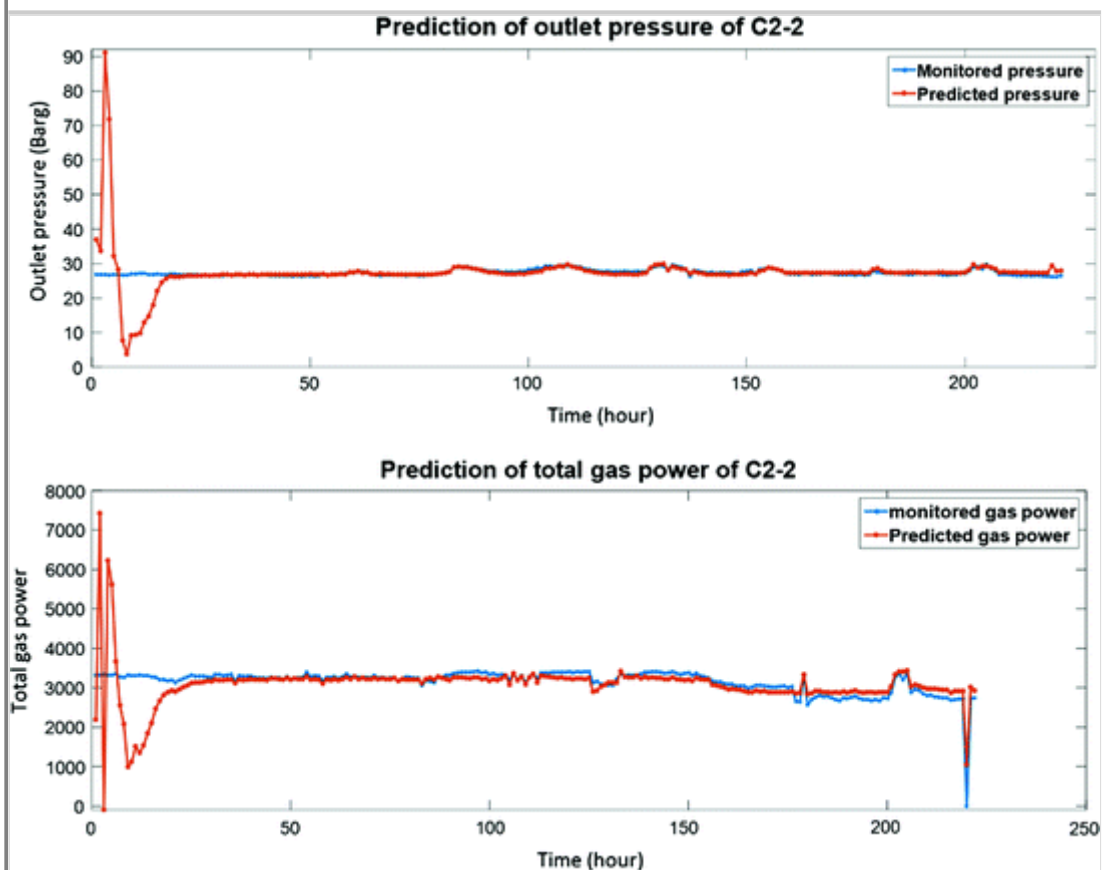
## All measured variables in data set C2



In order to maximize the accuracy of the built model, the past inputs  $x_k$  obtained from C2-1 was first used to predict the response of the system and the results were compared with the past outputs  $y_k$  to determine the prediction error. After several analyses testing, different values for  $r$ ,  $r = 32$  was finally adopted to give the minimum prediction error. The validated model was then used to make estimations of the process variables  $y_k$  of C2-2. Figure 4 shows the prediction results of the most significant variables of C2-2. The results show that the model causes large oscillations at the very beginning, but the oscillations dissipated quickly and the steady-state estimations are close to the actual measurements (Fig. 5). Fig 5 shows all the measured variables in dataset C3.

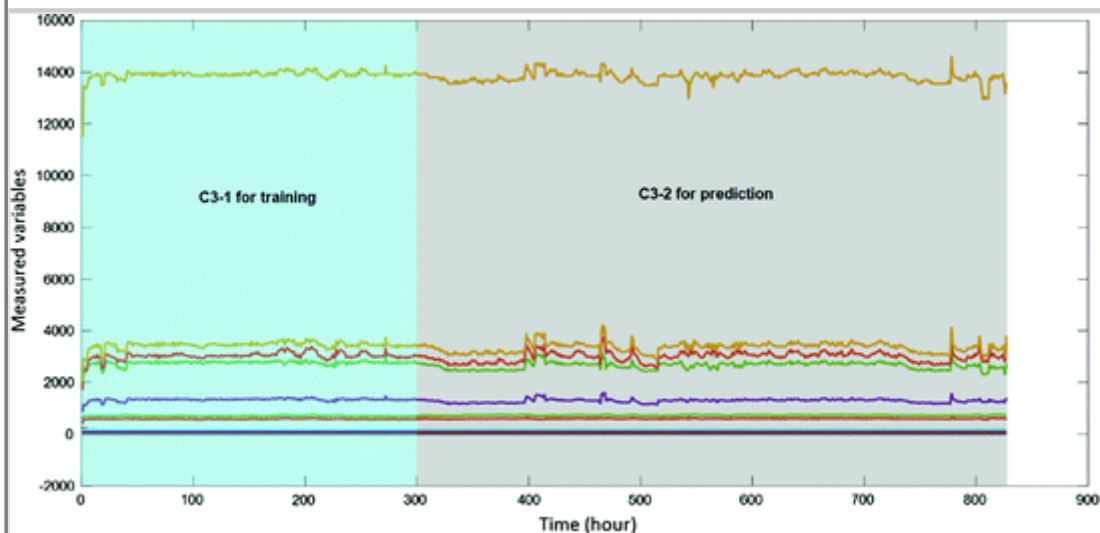
**Fig. 4**

Prediction results of the selected most significant variables of C2-2



**Fig. 5**

All measured variables in data set C3



AQ4

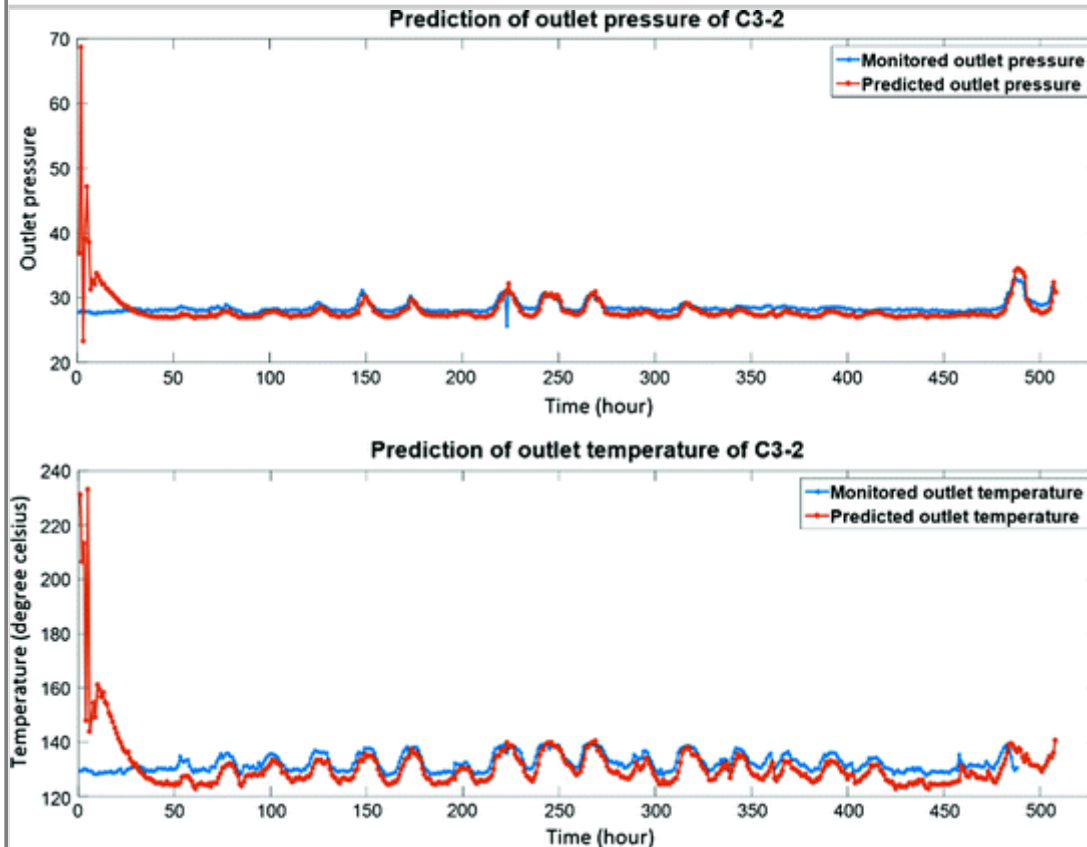
Similar to the procedure described above, the methodology was applied to predict the system outputs of another data set C3. Figure 6 shows the prediction results of the most significant variables of C3-2. The estimation is able to accurately represent the system behavior. Table 1 shows the average prediction error for the most significant variables of data set C2-2 and C3-2. The prediction error for each one of the significant variables was calculated by computing the mean of the difference between the predicted and measured signal:

$$e_i = \frac{1}{T} \sum_{t=1}^T \left| \frac{y_{t,i} - \hat{y}_{t,i}}{y_{t,i}} \right|$$

where  $y_{t,i}$  denotes the measured value of the  $i$ th variable at time  $t$ , and  $\hat{y}_{t,i}$  represents the estimated value of  $y_i$  and time  $t$ , and  $T$  denotes the total number of observations of the testing data set.

**Fig. 6**

Prediction results of the selected most significant variables of C3-2



**Table 1**

Average prediction error for significant variables

Data set	Outlet pressure (%)	Outlet temperature (%)	Ploy head (%)	Ploy head efficiency (%)	Bearing temperature V24 (%)	Total gas power (%)
C2-2	5.96	11.91	6.21	7.53	10.33	7.55
C3-2	3.63	3.67	6.14	2.75	15.33	8.37

## 4. Conclusion

Process data acquired from an operational industrial centrifugal compressor have been used to test the capabilities of CVA for fault detection, system identification, and performance estimation. The faults in data set C1 were successfully detected by  $T^2$  health indicator within a short detection time. CVA was also employed to build a state-space model for system identification. In order to fully capture the dynamics of the compressor, nine training data sets were selected from different operating conditions to train the CVA



algorithm. The trained CVA model was then used to predict the future system outputs for the specified system inputs. Although very fast and large oscillations were observed in the initial estimations, the average prediction error was low, proving that the model is able to represent the system dynamics under different operating conditions.

Although the results of this study clearly show the superior performance of the CVA algorithm for dynamic process monitoring, some things require further investigation. First, if CVA is employed to detect faults for a system operating under variable working conditions, it may produce high false alarm rates because the sudden changes in working conditions can be mistaken for performance degradations. Second, CVA-based performance estimation is based on the premise that the future system inputs are obtained from the production plan or from forecasts based on historic inputs, but the actual future inputs may be different from the forecasts due to process uncertainties, leading to inaccurate performance estimations. In addition, the proposed method is unable to predict the system behavior without knowing the future information. Therefore, more work should be conducted to apply the CVA for prognosis without future inputs.

## References

- Chiang LH, Russell EL, Braatz RD (2000) *Fault detection and diagnosis in industrial systems*. Springer, New York, London
- Huang LZ, Cao YP, Tian XM, Deng XG (2015) A Nonlinear quality-relevant process monitoring method with kernel input-output Canonical Variate analysis. *IFAC-PapersOnLine* 48:611–616. <https://doi.org/10.1016/j.ifacol.2015.09.035>
- Jiang B, Zhu X, Huang D, Braatz RD (2015a) Canonical variate analysis-based monitoring of process correlation structure using causal feature representation. *J Process Control* 32:109–116. <https://doi.org/10.1016/j.jprocont.#>
- Jiang B, Huang D, Zhu X, Yang F, Braatz RD (2015b) Canonical variate analysis-based contributions for fault identification. *J Process Control* 26:17–25. <https://doi.org/10.1016/j.jprocont.2014.12.001>
- Juricek BC, Seborg DE, Larimore WE (2001) Identification of the Tennessee Eastman challenge process with subspace methods. *Control Eng Pract* 9:1337–1351. [https://doi.org/10.1016/s0967-0661\(01\)00124-1](https://doi.org/10.1016/s0967-0661(01)00124-1)
- Komulainen T, Sourander M, Jämsä-Jounela SL (2004) An online application of dynamic PLS to a dearomatization process. *Comput Chem Eng* 28:2611–2619. <https://doi.org/10.1016/j.compchemeng>
- Ku W, Storer RH, Georgakis C (1995) Disturbance detection and isolation by dynamic principal component analysis. *Chemometr Intell Lab Syst* 30:179–196. [https://doi.org/10.1016/0169-7439\(95\)00076-3](https://doi.org/10.1016/0169-7439(95)00076-3)
- Larimore WE (1990) Canonical variate analysis in identification, filtering, and adaptive control. In: *Proceedings of the 29th IEEE conference on decision and control*, pp 596–604
- Lee C, Lee IB (2008) Adaptive monitoring statistics with state space model updating based on canonical variate analysis. *Korean J Chem Eng* 25:203–208. <https://doi.org/10.1007/s11814-008-0037-y>
- Muradore R, Fiorini P (2012) A PLS-Based statistical approach for fault detection and isolation of robotic manipulators. *IEEE Trans Ind Electron* 59:3167–3175. <https://doi.org/10.1109/tie.2011.2167110>
- Negiz A, Çınar A (1997) Statistical monitoring of multivariable dynamic processes with state-space models. *AIChE J* 43:2002–2020. <https://doi.org/10.1002/aic.690430810>
- Negiz A, Çınar A (1997) PLS, balanced, and canonical variate realization techniques for identifying VARMA models in state space. *Chemometr Intell Lab Syst* 38:209–221. [https://doi.org/10.1016/s0169-7439\(97\)00035-x](https://doi.org/10.1016/s0169-7439(97)00035-x)

- Negiz A, Çinar A (1998) Monitoring of multivariable dynamic processes and sensor aditing. *J Process Control* 8:375–380. [https://doi.org/10.1016/s0959-1524\(98\)00006-7](https://doi.org/10.1016/s0959-1524(98)00006-7)
- Odiowei PP, Cao Y (2010) Nonlinear dynamic process monitoring using canonical variate analysis and kernel density estimations. *IEEE Trans Ind Inf* 6:36–45. <https://doi.org/10.1109/tii.2009.2032654>
- Ruiz-Cárcel C, Lao L, Cao Y, Mba D (2016) Canonical variate analysis for performance degradation under faulty conditions. *Control Eng Pract* 54:70–80. <https://doi.org/10.1016/j.conengprac>
- Ruiz Cárcel C, Mba D (2014) A benchmark of canonical variate analysis for fault detection and diagnosis—IEEE Xplore document. <http://ieeexplore.ieee.org/document/6915178/?part=1>
- Russell EL, Chiang LH, Braatz RD (2000) Fault detection in industrial processes using canonical variate analysis and dynamic principal component analysis. *Chemometr Intell Lab Syst* 51:81–93
- Samuel R, Cao Y (2015) Kernel canonical variate analysis for nonlinear dynamic process monitoring. In: Cao Y (ed) 9th international symposium on advanced control of chemical processes. [http://ac.els-cdn.com/S2405896315011155/1-s2.0-S2405896315011155-main.pdf?\\_tid=8a1aa7e2-d649-11e6-a59a-00000aacb35d&acdnat=1483952382\\_e2740e26016619a37b573e7cc86ccc74](http://ac.els-cdn.com/S2405896315011155/1-s2.0-S2405896315011155-main.pdf?_tid=8a1aa7e2-d649-11e6-a59a-00000aacb35d&acdnat=1483952382_e2740e26016619a37b573e7cc86ccc74)
- Yang Y, Chen Y, Chen X, Liu X (2012) Multivariate industrial process monitoring based on the integration method of canonical variate analysis and independent component analysis. *Chemometr Intell Lab Syst* 116:94–101. <https://doi.org/10.1016/j.chemolab.2012.04.013>