# APPLICATIONS OF MULTIVARIATE ANALYSIS TECHNIQUES FOR FAULT DETECTION, DIAGNOSIS AND ISOLATION

PREM KRISHNAN

NATIONAL UNIVERSITY OF SINGAPORE

2011

# TABLE OF CONTENTS

# Summary

In this study, powerful multivariate tools such as Principal Component Analysis (PCA), Partial Least Squares (PLS) and Correspondence Analysis (CA) are applied to the problem of fault detection, diagnosis and identification and their efficacies are compared. Specifically, CA which has been recently adapted and studied for FDD applications is tested for its robustness when compared to other conventional and familiar methods like PCA and PLS on simulated datasets from three industry-based, high-fidelity simulation models. This study demonstrates that CA can negotiate time varying dynamics in process systems as compared to the other methods. This ability to handle dynamics is also responsible for providing robustness to CA based FDD scheme. The results also confirm previous claims that CA is a good tool for early detection and concrete diagnosis of process faults.

In, the second portion of this work, a new integrated CA and Weighted Pairwise Scatter Linear Discriminant Analysis method is proposed for fault isolation and identification. This tool tries to exploit the discriminative ability of CA to clearly distinguish between faults in the discriminant space and also predict if an abnormal event presently occurring in a plant is related to any previous faults that were recorded. The proposed method was found to give positive results when applied to simulated data containing faults that are either a combination of previously recorded failures or at intensities which are different from those previously recorded.

# LIST OF TABLES

# LIST OF FIGURES

# NOMENCLATURE

*A*                     The selected number of components/axes in PCA/PLS/CA

**A, B, C, D**          Parameter matrices in the state space model

*Aa*                    Principal axes (loadings) of the columns

*Bb*                    Principal axes (loadings) of the rows

*BB*                    The regression co-efficient matrix in PLS

*c*                     The vector of column sums in CA

*c*                     space of points of the class space in FDD system

*CC*                    The weight matrix of the output vector in PLS

*CM*                    The correspondence matrix in CA

*d*                     space of points of the decision space in FDD system

$D_\mu$                 Diagonal matrix containing the singular values for CA

$D_c$                   Diagonal matrix containing the values of the column sums from *c*

$D_r$                   Diagonal matrix containing the values of the row sums from *r*

*E*                     The residual matrix of the input in PLS

*EM*                    The expected matrix in CA

| | |
|---|---|
| *F* | The residual matrix of the output in PLS |
| *FF* | Scores of the row cloud in CA |
| *ff* | the score for the current sample |
| *g* | The scaling factor for chi-squared distribution in PLS model |
| *GG* | Scores of the column cloud in CA |
| *gg* | The grand sum of all elements in the input matrix in CA |
| $\boldsymbol{H}(z), \boldsymbol{G}(z)$ | Polynomial matrices in the input-output model |
| *I* | The number of rows in the input matrix in CA |
| *J* | The number of columns in the input matrix for CA |
| *K* | Number of decision variables in decision space in FDD system |
| *M* | Number of failure classes in class space in FDD system |
| *mc* | The number of columns (variables) in dataset *X* |
| *MO* | The number of columns in the output matrix in PLS |
| *mo* | The number of rows in the output matrix in PLS |
| *n* | Number of dimension in measurement space in FDD system |
| *NI* | The number of columns (variables) in the input matrix in PLS |
| *ni* | The number of rows in the input matrix in PLS |

| | |
|---|---|
| *nr* | The number of rows (samples) in dataset $X$ |
| *P* | The loadings (eigenvectors) of the Covariance Matrix in PCA |
| $P_A$ | The loadings only with the first $A$ columns included |
| *PP* | The matrix of loadings of the input in PLS |
| *q* | The new $Q$ statistic for the new sample $x$ |
| *QQ* | The matrix of the loadings of the output in PLS |
| $Q_\alpha$ | The $Q$ limit for the PCA/CA/PLS model at the $\alpha$ level of significance |
| *r* | The vector of row sums in CA |
| *res* | The residual vector formed for the new sample $x$ or $xx$ in PCA/CA |
| $r_{sample}$ | the row sum for the new sample |
| *S* | The variance-covariance matrix in PCA |
| *SM* | The chi squared matrix in CA |
| *t* | New score vector for a new sample $x$ |
| *T* | The scores (latent) variables obtained in PCA |
| $t^2$ | The $T^2$ statistic for the new sample $x$ |
| $T^2$ | The $T^2$ statistic used for the historical dataset |
| $T^2_\alpha$ | The $T^2$ limit for the PCA/CA/PLS model at the $\alpha$ level of significance |

| | |
|---|---|
| $T_A$ | The scores calculated for the first $A$ PCs alone in PCA |
| $t_{\text{new}}$ | The new score vector for input sample $x_{input-new}$ for PLS |
| $TT$ | The latent vector of the input variables in PLS |
| $U$ | The latent vector of the output variables in PLS |
| $\mathbf{u}(t)$ | Input signals for the state space model |
| $V$ | The eigenvectors (loadings) of the covariance matrix in PCA |
| $W$ | The weight matrix of the input vector in PLS |
| $X$ | The dataset matrix on which PCA will be applied |
| $x$ | Vector representation of the measurement space or new sample |
| $X_{input}$ | The input matrix for PLS calculations |
| $x_{input-new}$ | The new input sample for PLS |
| $xx$ | The new sample for CA |
| $\dot{x}_{new}$ | The predicted values of the new sample by the PLS model |
| $x'_{new}$ | The residual vector obtained for new sample in PLS |
| $Y$ | The output matrix for PLS calculations |
| $y$ | space of points of the feature space in FDD system |
| $\mathbf{y}(t)$ | Output signal for the state space model |

| XX | The input matrix in CA |
|---|---|

## Greek Letters

| $\Lambda$ | The diagonal matrix containing the eigenvalues in PCA |
|---|---|
| $\alpha$ | The level of significance for confidence intervals |
| $\Lambda_A$ | The diagonal matrix with eigenvalues equal to the chosen $A$ components |

## Abbreviations

| CA | Correspondence Analysis |
|---|---|
| CPV | Cumulative Percentage Variance |
| CUSUM | Cumulative Sum |
| CV | Cross Validation |
| DPCA | Dynamic Principal Component Analysis |
| EWMA | Exponentially Weighted Moving Average |
| FDA | Fisher Discriminant Analysis |
| FDD | Fault Detection and Diagnosis |
| KPCA | Kernel Principal Component Analysis |
| LDA | Linear Discriminant Analysis |
| MPCA | Multi-way Principal Component Analysis |

NLPCA            Non-Linear Principal Component Analysis

PCA              Principal Component Analysis

PLS              Partial Least Squares

WPSLDA           Weighted Pairwise Scatter Linear Discriminant Analysis

# 1. INTRODUCTION

## 1.1 Fault Detection and Diagnosis

It is well known that the field of process control has achieved considerable success in the past 40 years. Such a level of advancement can be attributed primarily to the computerized control of processes, which has led to the automation of low-level yet important control actions. Regular interventions like the opening and closing of valves, performed earlier by plant operators, have thus been completely automated. Another important reason for the improvement in control technology can be seen in the progress of distributed control and model predictive systems. However, there still remains the vital task of managing abnormal events that could possibly occur in a process plant. This task which is still undertaken by plant personnel involves the following steps

1) The timely detection of the abnormal event

2) Diagnosing the origin(s) of the problem

3) Taking appropriate control steps to bring the process back to normal condition

These three steps have come to be collectively called Fault Detection, Diagnosis and Isolation. Fault Detection and Diagnosis (FDD), being an activity which is dependent on the human operator, has always been a cause for concern due to the possibility of erroneous judgment and actions during the occurrence of the abnormal event. This is mainly due to the broad spectrum of possible abnormal occurrences such as parameter drifts, process failure or degradation, the size and complexity of the plant posing a need to monitor a large number of process variables and the insufficiency/non-reliability of process measurements due to causes like sensor biases and failures (Venkatasubramaniam *et al.,* 2003a).

## 1.2 The desirable characteristics of a FDD system

It is essential for any FDD system to have a desired set of traits to be acknowledged as an efficient methodology. Although there are several characteristics that are expected in a good FDD system, only some are extremely necessary for the running of today's industrial plants. Such characteristics include the quick detection of an abnormal event. The term 'quick' does not just refer to the earliness of the detection but also the correctness of the same, as FDD systems under the influence of process noise are known to lead to false alarms during normal operation. Multiple fault identifiability is another trait where the system is able to flag multiple faults despite their interacting nature in a process. In a general nonlinear system, the interactions would usually be synergistic and hence a diagnostic system may not be able to use the individual fault patterns to model the combined effect of the faults (Venkatasubramaniam *et al.,* 2003a). The success of multiple fault identifiability can also lead to the achievement of novel identifiability by which a fault occurring may be distinguished as being a known (previously occurred) or an unknown (new) one.

## 1.3 The transformations in a FDD system

It is essential to identify the various transformations that process measurements go through before the final diagnostic decisions could be made.

1) Measurement space: This is the initial status of information available from the process. Usually, there is no prior knowledge about the relationship between the variables in the process. It can literally be called as the plant or process data being recorded at regular intervals and can be represented as $x_1, x_2, \ldots, x_3$ where '$n$' refers to the number of variables.

2) Feature space: This is the space where the features are obtained from the data utilizing some form of prior knowledge to understand process behavior. This representation could be obtained by two means, namely feature selection and feature extraction. Feature selection simply deals with the selection of certain key variables from the measurement space. Feature extraction is the process of understanding the relationship between the variables in the measurement space using prior knowledge. This relationship between the variables is then represented in the form of a fewer parameters thus reducing the size of the information obtained. Another main advantage is that the features cluster well to aid in classification and discrimination for the remaining stages. The space can be seen as $y = [y_1, y_2, ..., y_3]$ where $y_i$ is the $i^{th}$ feature obtained.

3) Decision Space: This space is obtained by subjecting the feature space to meet an objective function which could be some kind of discriminant or simple threshold function. It is shown as $d = [d_1, d_2, ..., d_K]$ where '$K$' is the number of decision variables obtained.

4) Class Space: This space is a set of integers which can be presented as $c = [c_1, ..., c_M]$ that are a reference to '$M$' number of failure classes and normal class of data to any of which a given measurement pattern may belong.

## 1.4 Classification of FDD Algorithms

The classification of FDD classifier algorithms is usually based on the kind of search strategy employed by the method. The kind of search approach used to aid diagnosis is dependent on the way in which the process information scheme is presented which in turn is largely influenced by the type of prior knowledge provided. Therefore, the type of prior knowledge would provide the basis for the broadest classification of FDD algorithms. This *a priori* knowledge is supposed to

3

give the set of failures and the relationship between the observations and failures in an implicit or explicit manner. The two types of FDD methodologies under this basis include model-based methods and process history-based methods. The former refers to methods where fundamental understanding of the physics and chemistry (first principles) of the process is used to represent process knowledge while, in the latter, data based on past operation of the process is used to represent the normal/abnormal behavior of the process. Model based methods can, once again, be broadly classified into quantitative and qualitative models.

An important point to be noted here is that while it is indeed true that any type of model would require data finally to obtain its parameter values, and that all FDD methods need to create some kind of a model to aid their task. Therefore, the actual significance behind the use of the term model based methods is that the physical understanding of the process has already provided assumptions for the model framework and the form of prior knowledge. Meanwhile, process history methods are equipped with only large heaps of data from where the model is itself created from the same in such a form so to have extracted features from the data.

### 1.4.1 Quantitative and Qualitative models

Quantitative models portray the relationships between the inputs and outputs in the form of mathematical functions whereas qualitative models represent the same association in the form of causal models.

The work with quantitative models began as early as the late 1970's with attempts to apply first principles model directly (Himmelblau, 1978) but this was often associated with computational complexity rendering the models of questionable utility in real time applications. Therefore, the main kind of models usually employed were the ones relating the inputs to the outputs (input-

output models) or those related with the identification of the input output link via internal system states (State Space models).

Let us consider a system based on *'m'* inputs to the system and *'k'* outputs. Let, $\mathbf{u}(t) = [u_1(t)\, u_2(t)\, ...\, u_m(t)]'$ be the input signals and $\mathbf{y}(t) = [y_1(t)\, y_2(t)\, ...\, y_k(t)]'$ be the output signals, then the basic system model in the state space form is,

$$\mathbf{x}(t+1) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \tag{1.1}$$

$$\mathbf{y}(t+1) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t) \tag{1.2}$$

where **A, B, C** and **D** are parameter matrices with appropriate dimensions and $x(t)$ refers to the state vector.

The input - output form is given by,

$$\mathbf{H}(z)\mathbf{y}(t) = \mathbf{G}(z)\mathbf{u}(t) \tag{1.3}$$

where $\mathbf{H}(z)$ and $\mathbf{G}(z)$ are polynomial matrices.

When the fault does occur, the model will generate inconsistencies between the actual and expected value of the measurements. This indicates deviation from normal behavior and such inconsistencies are called residuals. The check for such inconsistencies requires redundancy. The main task, here, consists of the detection of faults in the processes using the dependencies between different measurable signals established through algebraic or temporal relationships. This form of redundancy is termed analytical redundancy (Chow & Willsky, 1984; Frank, 1990) and is more frequently used than hardware redundancy which involves using more sensors.

There are two kinds of faults that are modeled. On one hand, we have additive faults which refer to the offset of sensors and other disturbances such as actuator malfunctioning or a leakages in pipelines. On the other hand, we have multiplicative faults which represent parameter changes in the process model. These changes are known to have an important impact on the dynamics of the model. Problems caused by fouling, contamination usually come under this category (Huang *et al.*, 2007). Incorporation of terms for both these faults in both state space and input–output models can be found in control literature (Gertler, 1991, 1992). As mentioned earlier, residuals generated are required to perform FDI actions in quantitative models; this is done on the basis of analytical redundancy in both static and dynamic systems. For static systems, the residual generator will also be static i.e. a rearranged form of the input-output models (Potter & Suman, 1977) or material balance equations (Romagnoli & Stephanopoulus, 1981). In dynamic systems, residual generations is developed using techniques such as diagnostic observers, Kalman filters, parity relations, least squares and several others. Since process faults are known to either affect the state variables (additive faults) or the process parameters, it is possible to estimate the state of the system using Kalman filters (Frank & Wunnenberg, 1989). Dynamic observers are algorithms that estimate the states based on the process model's observed inputs and outputs. Their aim is to develop a set of robust residuals which will help to detect and uniquely identify different faults such that their decision making is not affected by unknown inputs or noise. The least squares method is more concerned with the estimation of model parameters (Isermann, 1989). Parity equations, a transformed version of the state space and input output models have also been used for generation of residuals to aid in diagnosis (Gertler, 1991, 1998). Li & Shah (2000) developed a novel structured residual based technique for the detection and isolation of sensor faults in dynamic systems which was more sensitive as compared to the scalar based

counterparts developed by Gertler (1991, 1998). The novel technique was able to provide a unified approach to the isolation of single and multiple sensor faults together. A novel FDI system for non-uniformly sampled multirate system was developed by Li & Shah (2004) by extending the Chow-Willsky scheme from single rate systems to multirate systems. This generates a primary residual vector (PRV) for fault detection and then by structuring the PRV to have different sensitivity/insensitivity to different faults, fault isolation is also performed.

As mentioned earlier, quantitative models express the relationship between the inputs and outputs in the form of mathematical functions. In contrast, qualitative models present these relationships in the form of qualitative functions. Qualitative models are usually classified based on the type of qualitative knowledge used to develop these qualitative functions; these include diagraphs, fault trees and qualitative physics.

Cause-effect relations or models can be represented in the form of signed digraphs (SDG). A digraph is a graph with directed arcs between the nodes and SDG is a graph in which the directed arcs have a positive or negative sign attached to them. The directed arcs lead from the 'cause' nodes to the 'effect' nodes. SDGs provide a very efficient way of representing qualitative models graphically and have been the most widely used form of causal knowledge for process fault diagnosis (Iri *et al.*, 1979; Umeda *et al.*, 1980; Shiozaki *et al.*, 1985; Oyeleye and Kramer, 1988; Chang and Yu, 1990). Fault trees models are used in analyzing system reliability and safety. Fault tree analysis was originally developed at Bell Telephone Laboratories in 1961. Fault tree is a logic tree that propagates primary events or faults to the top level event or a hazard. The tree usually has layers of nodes. At each node different logic operations like AND and OR are performed for propagation. Fault-trees have been used in a variety of risk assessment and reliability analysis studies (Fussell, 1974; Lapp and Powers, 1977). Qualitative physics

knowledge in fault diagnosis has been represented in mainly two ways. The first approach is to derive qualitative equations from the differential equations termed as confluence equations. Considerable work has been done in this area of qualitative modeling of systems and representation of causal knowledge (Simon, 1977; Iwasaki and Simon, 1986; de Kleer and Brown, 1986). The other approach in qualitative physics is the derivation of qualitative behavior from the ordinary differential equations (ODEs). These qualitative behaviors for different failures can be used as a knowledge source (Kuipers, 1986; Sacks, 1988).

### 1.4.2 Process history based models

Process history based models are concerned with the transformation of large amounts of historical data into a particular form of prior knowledge which will enable proper detection and diagnosis of abnormalities. This transformation is called feature extraction, which can be performed qualitatively or quantitatively.

Qualitative feature extraction is mostly developed in the form of expert systems or trend modeling procedures. Expert Systems may be regarded as a set of if-else rules set on analysis and inferential reasoning of details in the data provided. Initial work in this field has been attempted by Kumamato *et al.* (1984), Niida *et al.* (1986), Rich *et al.* (1989). Trend modeling procedures tend to capture the trends in the data samples at different timescales using slope (Cheung & Stephanopoulos, 1990), finite difference (Janusz & Venkatasubramanian, 1991) calculations and other methods after initially removing the noise in the data using noise-filters (Gertler, 1989). This kind of analysis facilitates better understanding of the process and hence diagnosis.

Quantitative procedures are more prompted towards the classification of data samples into separate classes. Statistical methods like Principal Component Analysis (PCA) or PLS perform this classification on the basis of prior knowledge in class distributions, while non-statistical methods like Artificial Neural Networks use functions to provide decisions on the classifiers.

## 1.5  Motivation

In present day industries, plant engineers are on the lookout for tools and methods that tend to be more robust in nature i.e. those that indicate less number of false alarms even at the compromise of mild delays in detection or relatively less detection rates. The reason for this is that, repeated occurrences of false alarms events would leave plant personnel in a state of ambiguity and lacking faith in the tool. Another major problem in the industry is multiple fault identifiability when some of the faults follow a similar trend and cannot be distinguished clearly leading to improper diagnosis.  The part that multiple fault identifiability plays in providing a clear picture of the nature of faults in a process will eventually lead to the proper identification of future fault i.e. novel fault identifiability. The solution and handling of these three problems are important in better running of industrial plants and will eventually lead to greater profits. In this regard, statistical tools are found to be the most successful in application to industrial plants. This can be attributed to their low requirements in modeling efforts and less a priori knowledge of the system involved (Venkatasubramaniam *et al.,* 2003c). The main motivation for this work would be to identify a statistical tool which would satisfy the above mentioned traits at an optimum level. This is determined by comparing the FDD application of contemporary popular statistical tools alongside recent ones on certain examples.

Table 2.1: Comparison of Various Diagnostic methods

| | Observer | Diagraphs | Abstraction hierarchy | Expert Systems | QTA | PCA | Neural networks |
|---|---|---|---|---|---|---|---|
| Quick detection and diagnosis | ✓ | ? | ? | ✓ | ✓ | ✓ | ✓ |
| Isolability | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ |
| Robustness | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Novel Identifiability | ? | ✓ | ✓ | ✗ | ? | ✓ | ✓ |
| Classification Error | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Adaptability | ✗ | ✓ | ✓ | ✗ | ? | ✗ | ✗ |
| Explanation Facility | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ |
| Modeling Requirement | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Storage and Computation | ✓ | ? | ? | ✓ | ✓ | ✓ | ✓ |
| Multiple fault Identifiability | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |

Source: Venkatasubramaniam *et al.* (2003c).

Table 1.1 shows the comparison between several methods on the basis of certain traits that are expected in FDD tools. It is quite clear from Table 1.1 that statistical tool PCA is almost on par with other methods and also seems to satisfy two of the three essential qualities required in the industry. PCA, being a linear technique, is prone to only satisfy these qualities as long as the data comes from a linear or mildly non-linear system.

In this regard, the objective of this thesis is to compare a few statistical methods and determine which are most effective in FDD operations. The tools involved would include well known and

implemented methods such as PCA and PLS alongside Correspondence Analysis (CA) which is a recent addition to the FDD area. CA has been highlighted as having the ability to effectively handle time-varying dynamics of the process because it simultaneously analyzes the rows and columns of datasets. This work will show results which will compare robustness, extent of early detection and diagnosis of all the considered techniques. In addition to that, it will be demonstrated that an integrated technique featuring CA and Weighted Pairwise Scatter Linear Discriminant Analysis (CA-WPSLDA) will provide better multiple fault identifiability and novel identifiability as compared to PCA, FDA and WPSLDA.

**1.6 Organization of the thesis**

This thesis is divided into five chapters. Chapter 2 comprises of the literature survey and algorithms of the basic conventional methods such as PCA, PLS and CA. A comparison between PCA and CA is also made based on previous literature. Chapter 3 will feature results which will prove the robustness of CA as a fault detection tool based on the simulated datasets obtained from three systems, a Quadruple tank system, the Tennessee Eastman Challenge Process (TEP) and a Depropanizer process. Chapter 4 will provide a brief introduction and literature survey to feature extraction by FDA and its current role in FDD. This will be followed by a comparison of the FDA and CA techniques and the explanation of the integrated CA-WPSLDA technique for fault identification. The chapter will end with the application of these techniques to the quadruple tank system and Depropanizer process. The final chapter (Chapter 5) will contain the conclusions of the study and the prospects for future works.

# 2. LITERATURE REVIEW

This chapter will focus on the work that had been done in the field of fault detection, diagnosis (FDD) and with regard to the multivariate statistical techniques PCA, PLS and CA. The initial stages of this chapter will first explain the origins of PCA and PLS as FDD tools followed by an explanation of their algorithms and monitoring strategies based on them. This will be succeeded by the advances and modifications that have taken place with respect to these methods. A similar explanation of CA will then be provided involving its origin and algorithm followed by its comparison to PCA and PLS. The chapter will finally conclude stating the advantages of CA as compared to the other two methods.

## 2.1 Statistical Process Control

Statistical Process Control (SPC) may be referred to as one of the earliest versions of FDD based on statistics. SPC is a statistical procedure which determines if a process is in a state of control by discriminating between what is called common cause variation and assignable cause variation (Baldassarre *et al.*, 2007). Common cause variation refers to the variations that are inherent in the process and cannot be removed without changing the process. In contrast, assignable cause variation refers to the unusual disruptions and abnormalities in the process. In this context, a process is said to be "in statistical control" if the probability distribution representing the quality characteristic is constant over time (Woodall, 2000). Thus, one could check if the process adheres to the distribution by setting the parameter values that include the Central Line (CL) or tangent, the Upper Control Limit (UCL) and the Lower Control Limit (LCL) for the process based on the properties of the distribution. The CL would be the best representation of quality while the UCL and LCL would encompass the region for common cause variation. If the data

monitored violates the UCL or LCL, one can come to the conclusion that there is the strong possibility of an abnormal event in progress. The first control chart to be developed was the Shewhart chart (Shewhart, 1931) Chart. The Shewhart chart is the simplest example of a control chart based on the Gaussian distribution. The CL in this chart would be the average of all the samples which appear to be in the normal region, the LCL is three times the standard deviation of the dataset subtracted from the average while the UCL is three times the standard deviation of the dataset added to the average. Thus, in accordance with the properties of normal distribution, the limits are set such that only 1% of the data points are expected to fall outside the limits "by chance". SPC gained more prominence with the use of other univariate control charts such as Cumulative Sum (CUSUM) (Woodward and Goldsmith, 1964), Exponentially Weighted Moving Average (EWMA) (Roberts, 1959; Hunter, 1986) to monitor important quality measurements of the final product. The problem with analyzing one variable at a time is that not all the quality variables are independent of each other making detection and diagnosis difficult (MacGregor and Kourti, 1995). This led to the need to treat all the variables simultaneously, thus creating the need for multivariate methods. This problem was at first solved using multivariate versions of all the previously mentioned control charts (Sparks, 1992). These methods were the first to use the $T^2$ statistic (Hotelling, 1931), a multivariate form of the Student's t-statistic which would set the control limits for the multivariate control charts.

The main problem encountered then was the fact that a large number of quality and process variables were being monitored in process plants due to being measured in process plants owing to improvements in instruments as well as their lowered costs. This rendered the application of multivariate control charts to be impractical for such high dimensional systems that exhibited significant collinearities between variables (Bersimis *et al.,* 2006). There was, therefore, a need

for methods that can reduce the dimensions in the dataset and utilize the high correlations existing amongst the process as well as quality variables. Such a need led to the use of PCA and PLS for FDD tasks.

**2.2 PCA and PLS**

**2.2.1 PCA – the algorithm**

PCA is a multivariate dimensional reduction technique that has been applied in the field of process monitoring and FDD for the past two decades. PCA transforms a number of possibly correlated variables in a dataset into a smaller number of uncorrelated pseudo or latent variables. This is done by a bilinear decomposition of the variance-covariance matrix of the dataset. The uncorrelated (orthogonal) variables obtained are called the principal components and they represent the axes obtained by rotation of the original co-ordinate system along the direction of maximum variance. The main assumptions in this method are that the data follows a Gaussian distribution and that all the samples are independent of one another.

The steps involved in the formulation of the PCA model for FDD operations are as follows:

Consider a dataset organized in the form of a matrix **X**, with $nr$ rows (samples) and $mc$ columns (variables). This matrix is initially pre-processed and normalized to give $X_0$. Normalization is necessary when the variable of the dataset will belong to different units and doing so will bring all the variables down to a mean value of zero and unit variance. This will ensure that all the variables have an equal opportunity to participate in the development of the model and subsequent analysis (Bro and Smilde, 2003). $X_0$ will then be decomposed to provide scores (latent variables) and loadings based on the NIPALS algorithm (Wold et al., 1987) or by Singular Value Decomposition (SVD) or Eigenvalue decomposition. The SVD or Eigenvalue

decomposition method (EVD) is preferred due to its advantages over NIPALS in PCA. These include fewer uncertainties associated with the eigenvalues and less round-off errors in the calculation (Seasholtz *et al.,* 1990).

Step 1: The sample covariance matrix is given by

$$S = \frac{1}{(nr-1)} X_0{}^T X_0 \tag{2.1}$$

Step 2: This covariance matrix S is then subjected to eigenvalue decomposition.

$$S = V\Lambda V^T \tag{2.2}$$

where matrix $\Lambda$ is the diagonal matrix containing the non-negative eigenvalues arranged in decreasing order ($\lambda_1 > \lambda_2 > \lambda_3 ... > \lambda_{mc}$). Matrix V contains the eigenvectors corresponding to the eigenvalues in $\Lambda$.

Step 3: Formulation of loadings and scores

$$P = V \tag{2.3}$$

$$T = X_0 P \tag{2.4}$$

The loadings P are the eigenvectors in the matrix V corresponding to the eigenvalues. The eigenvectors with the largest eigenvalues correspond to the dimensions that have the strongest correlation in the data set. The PCA scores T may be defined as transformed variables obtained as a linear combination of the original variables based on the maximum amount of variance captured. They are the observed values of the Principal Components for each of the nr original sample vectors.

Step 4: Monitoring and Detection

In the first step to monitoring, it is essential to choose the number of PCs required to capture the dominant information about the process (i.e. the signal space). The selection of A principal components could be done through the cross validation (CV) technique (Jackson, 1991) or the Cumulative Percentage Variance (CPV) technique. CV involves the splitting of the dataset into two (training and testing sets) or more parts a specified number of times. This is followed by the calculation and construction of a Predictive Residual Sum of Squares Plot (PRESS) in descending order and looks for the "knee" or "elbow" in the curve. The numbers of selected components is the one that is at the "knee" or "elbow" of the process plot.

The CPV is given by,

$$CPV = \frac{\sum_{i=1}^{A} \lambda_i}{\sum_{j=1}^{mc} \lambda_j} \times 100 \qquad (2.5)$$

When the CPV is found to be greater than a value (usually fixed at 80% or 85%), then A is fixed as the required number of components. This is then followed by the use of the $T^2$ and Q statistic for monitoring purposes.

The calculation of the $T^2$ statistic for the historical dataset is given by

$$T^2 = T_A^T \Lambda_A^{-1} T_A \qquad (2.6)$$

where, $T_A$ represents the scores calculated for the first A PCs and $\Lambda_A$ represent the diagonal matrix containing the first A eigenvalues. The $T^2$ statistic is a representation of the correlation within the dataset over several dimensions. It is the measurement of the statistical distance of the score values from the centre of the A-dimensional PC space (Mason and Young, 2002).

Monitoring of this statistic for any new m dimensional sample x is done by first normalizing it to give $x_0$. The new score vector t for the sample is given by,

$$t = x_0 P_A \tag{2.7}$$

where, $P_A$ represents the first A columns of the loadings matrix

$$t^2 = t^T \Lambda_A^{-1} t \tag{2.8}$$

Thus, the $T^2$ statistic value of any new sample can be calculated

The limit for this statistic for monitoring purposes can be obtained using the F-distribution as follows.

$$T_\propto^2 = \frac{A(nr^2-1)}{nr(nr-A)} F_\propto(A, nr - A) \tag{2.9}$$

The above mentioned equation expresses the fact that the limit is the value of the F-distribution with A and nr-A degrees of freedom at α level of significance (the level of alpha is mostly 90, 95 or 99 %). Any deviation from normality is indicated when $t^2 > T_\propto^2$.

The limitation of the $T^2$ statistic is that it will only detect an event if the variation in the latent variables is greater than the variation explained by common causes. This led to the development of the Q-statistic which is the sum of the squares of the residuals of the model and is a measure of the variance not captured by the model.

$$r = \left(I - P_A P_A^T\right) x_0 \tag{2.10}$$

where r is the residual vector and,

$$q = r^T r \tag{2.11}$$

17

The upper limit for the Q-statistic is given by,

$$Q_\alpha = [\frac{(h_0 c_\alpha \sqrt{2\theta_2})}{\theta_1} + 1 + \frac{(\theta_2 h_0 (h_0 - 1))}{\theta_1^2}]^{\frac{1}{h_0}} \qquad (2.12)$$

with,

$$\theta_i = \sum_{j=A+1}^{nr} \lambda_j^i \qquad (2.13)$$

$$h_0 = 1 - (\frac{2\theta_1 \theta_3}{3\theta_2^2}) \qquad (2.14)$$

Abnormalities which affect the correlation between the variables can be detected using the Q statistic when $q > Q_\alpha$.

Another use of the residual vector r is in the generation of contribution plots where each of the residual values is divided by the sum of all elements in it and presented in the form of bar plots to identify the variables that is most likely associated with the fault. Contribution plots are still being used as effective diagnostic tools.

PCA was initially used for SPC alone (application to quality variables) but was later applied to process variables as well, thus enabling it to act as a tool for Statistical Process Monitoring (SPM). Kresta *et al.* (1991) were the first to apply PCA to both process as well as quality variables. The main advantages of doing so was the  improved diagnosis and understanding of faults through the changes in process variables and the identification of drifts in process variables which cannot usually be noticed in quality variables for the same operating condition (Qin, 2003). It also enabled the application of the tool to processes where the quality variables are not recorded in the historical datasets (Bersimis *et al.*, 2007).

## 2.2.2 PLS – the algorithm

Partial Least Squares (PLS) is a dimensional reduction as well as a regression technique that finds a new set of latent variables which maximize the covariance between the input data matrix $X_{inp}$ and the output data matrix Y. The main objective here is to approximate $X_{inp}$ and Y into reduced dimensional forms as well as model a linear relationship between them. The application of PLS to systems for FDD is mostly done such that, the process variables are assigned to the data matrix $X_{inp}$ and the quality variables are assigned to the output matrix Y. PLS is performed mainly using two algorithms, namely the NIPALS algorithm (Geladi and Kowalski, 1986) and the SIMPLS algorithm (de Jong, 1993).

The input and output matrices are first normalized as in PCA. This is done by mean centering and dividing the values by the corresponding variance to give $X_{inp0}$ and $Y_0$. This brings all the variables in both matrices down to having a zero mean and unit variance and can hence be treated equally during the analysis. The NIPALS algorithm is applied to the PLS regression in order to sequentially extract the latent vectors TT and U and the weight vectors W and CC from the $X_{inp0}$ and $Y_0$ matrices in a decreasing order of their corresponding singular values of the cross-covariance matrix $X_{inp0}{}' Y_0$. As a result, PLS decomposes $X_{inp0}(ni \times NI)$ and $Y_0(mo \times MO)$ matrices into the form.

$$X_{inp0} = (TT)(PP)' + E \tag{2.15}$$

$$Y_0 = U(QQ)' + F \tag{2.16}$$

where TT and QQ are $(n \times A_{PLS})$ matrices of the extracted A score vectors, $PP(NI \times A_{PLS})$ and $QQ(MO \times A_{PLS})$ are matrices of loadings, and $E(ni \times NI)$ and $F(ni \times MO)$ represent matrices of residuals. The $A_{PLS}$ vectors are extracted using cross validation (CV).

The PLS regression model can be expressed with regression coefficient BB and residual matrix R as follows:

$$Y_0 = X_{inp0}BB + R \tag{2.17}$$

$$BB = W((PP)'W)^{-1}CC' \tag{2.18}$$

Rannar *et al.* (1994) derived the following equalities:

$$W = X_{inp0}'U \tag{2.19}$$

$$(PP) = X_{inp0}'(TT)((TT)'(TT))^{-1} \tag{2.20}$$

$$CC = Y_0' (TT)((TT)'(TT))^{-1} \tag{2.21}$$

Substituting the Equations $(2.19 - 2.21)$ into Equation $(2.18)$ using the orthogonality of the matrix TT columns, the matrix B can be written in the following form:

$$BB = X_{inp0}'U((TT)'X_{inp0}X_{inp0}'U)^{-1}(TT)'Y_0 \tag{2.22}$$

This will be used to make predictions in PLS regression i.e. compared with principal component regression, PLS considers the amount of input information and also accounts for the contribution of the input latent variables to the output.

The monitoring scheme for PLS with a new sample $x_{input-new}$ of the process variables is as follows:

$$t_{new} = R'x_{input-new} \tag{2.23}$$

where $t_{new}$ is the new score vector for the X-subspace.

$$\dot{x}_{new} = ((PP)R')x_{input-new} \tag{2.24}$$

$$r_{new} = (I - ((PP)R'))x_{input-new} \tag{2.25}$$

where $\dot{x}_{new}$ is the value predicted by the model and $r_{new}$ is the residual attached to the $X_{input}$ subspace. The $T^2$ and Q statistics are given by:

$$t^2 = t_{new}\Lambda_{PLS}^{-1}t_{new} \tag{2.26}$$

where, $\Lambda_{PLS} = \frac{1}{(n-1)}TT'TT$

$$q = \|r_{new}\|^2 \tag{2.27}$$

The calculation of the statistic limits remains the same for $T^2$ but varies for the Q statistic which is given by $g\chi_h^2$. Where, g is the scaling factor for the Chi-squared distribution with h degrees of freedom.

It must be noted that PLS which attempts to understand the covariance between $X_{inp}$ and Y does not provide the components in $X_{inp}$ in a descending order of its variance as some of them may be orthogonal to Y and therefore be useless in its prediction. Thus there is a possibility for large variability in the residual space after the selection of $A_{PLS}$ components leaving the Q statistic unsuitable for monitoring purposes (Zhou *et al.*, 2010).

## 2.2.3 The evolution of PCA and PLS in FDI

Some of the earliest works in PCA and PLS for SPC/SPM were done by Denney *et al*. (1985) and Wise *et al*. (1991). Finally, MacGregor and Kourti (1995) had successfully established that both PCA and PLS can be applied to several industrial processes such as sulphur recovery unit, low-density polyethylene process or fluidic bed catalytic cracking with the largest system containing a total of 300 process variables and 11 quality variables.

Nomikos and MacGregor (1994) extended PCA to batch processes by employing the Multi-way PCA (MPCA) approach where they proposed estimating the missing data on trajectory deviations from the current time until the end of the batch. Rannar *et al*. (1998) proposed the use of hierarchical PCA for adaptive batch monitoring to overcome the problem of estimating missing data. Since the simple PCA technique is based on the development of linear relationships among variables and their subsequent representation of industrial processes which are non-linear in nature, there was a need to develop techniques which were more effective in representing the non-linearity in the system, this necessity led to the first work on Non-Linear PCA (NLPCA) developed by Kramer (1991) who used neural networks to achieve the required non-linear dimensional reduction and representation. Dong and McAvoy (1996) improved the NLPCA method by employing Principal Component Curves but the methods were still difficult to use owing to the need for non-linear optimization and estimation of number of components prior to training of the network. The problem of non-linear optimization in NLPCA was handled by the use of Kernel PCA (KPCA) where the nonlinear input is transformed to a hidden high dimensional space where features are extracted using a Kernel function. The earliest attempts at KPCA were by Scholkopf *et al.* (1998). Some variants of the KPCA include the Dynamic KPCA by Choi and Lee (2004) using a time lagged matrix. Application of Multi-way KPCA to batch

processes was demonstrated by Lee *et al.* (2004).  One important problem involved  in KPCA were increase the size of the dataset to higher dimensions leading to computational difficulties (Jemwa & Aldrich, 2006) but this was taken care of by representing the calculations in the feature space in the form of dot products. Another important problem present in PCA is that it is time invariant while most of the processes are time varying and dynamic in nature. This led to the development of recursive PCA developed by Li *et al.* (2000). Dynamic PCA (DPCA) was seen as another tool to handle this problem; it was developed by incorporating time as an additional column in the dataset using time series models such as the ARX model (Russell *et al.*, 2000).

The use and development of PLS in the field of process monitoring was also widespread especially owing to its ability to identify relationships between the process and quality variables in the system. MacGregor and Kourti (1995) were the first to suggest the use of multi-block PLS as an efficient tool for diagnosis when there are a large number of process variables to be handled. As PLS too being a linear technique like PCA had limitations dealing with non-linearities, Qin and McAvoy (1992) developed the first neural network PLS method which employed feedforward networks to tackle this problem. The problem of time-invariance in PLS led to the development of the first dynamic PLS algorithm by Kaspar and Ray (1993) to be used in the modeling and control of processes. Lakshminarayanan *et al.* (1997) later used a dynamic PLS algorithm towards the simultaneous identification and control of chemical processes and also provided a design for feed forward controllers in multivariate processes using the PLS framework. A Recursive PLS algorithm was developed by Qin (1998) to handle the same issue. Vijaysai *et al.* (2003) later extended this algorithm to provide a blockwise recursive PLS

technique based on the segregation of old and new data for dynamic model identification under closed loop conditions.

## 2.3 Correspondence Analysis

### 2.3.1 The method and algorithm

Correspondence analysis (CA) is a multivariate exploratory analysis tool that aims to understand the relationship between the rows and columns of a dataset. It has come a long way in the 30 years since the publication of Benzécri's seminal work, Analyse des Données (Benzécri *et al.,*1973) and, shortly thereafter, in Hill's paper on applied statistics, (Hill, 1974). This work was further explained by Greenacre (1987 and 1988) and made popular in various applications including social sciences, medical data analysis and several other areas (Greenacre, 1984 and 1992). CA can be defined as a two way analysis tool which seeks to understand the relationship between the rows and columns of a contingency table (cross tabulation calculations which are clearly explained by Simpson (1951)).

In this approach, let us assume that we have a matrix XX with I rows and J columns. Initial scaling of the data is necessary as, only a single form (common unit/mode of measurement) of data could be fit into several categories; it would not make much sense to analyze different scales of data in the form of relative frequencies (Greenacre, 1993). The form of scaling adopted is to bring all the values in the matrix within the scale of 0 to 1 as CA being a categorical variable method cannot handle negative values (Detroja *et al.,* 2006).

Step 1:  Calculation of the Correspondence Matrix CM.

$$CM = \left(\frac{1}{gg}\right)XX \tag{2.28}$$

where, CM is the correspondence matrix and $gg$ is the grand sum (sum of all elements in the matrix). The main objective here is to convert all values along rows and columns to the form of relative frequencies.

Step 2: In this step, the row sums and column sums of CM are calculated, they are given by,

$$r_i = \sum_{j=1}^{J} cm_{ij} \tag{2.29}$$

$$c_j = \sum_{i=1}^{I} cm_{ij} \tag{2.30}$$

where, $r$ and $c$ are vectors containing the row (I values) and column sums (J values).

Step 3: In this step, the null hypothesis of independence is assumed by which no row or column is associated to one another. According to this assumption, the actual values of the correspondence matrix CM should be such that each element is given by the product of the corresponding row and column sum of the matrix. These expected values are stored in what is called the Expected Matrix EM, where,

$$em_{ij} = r_i c_j \tag{2.31}$$

The centering would involve calculating the difference between the observed and expected difference between the expected and observed relative frequencies, which is then normalized by dividing the difference of each value by the square root of the corresponding expected value,

$$sm_{ij} = \frac{(cm_{ij} - em_{ij})}{\sqrt{em_{ij}}} \tag{2.32}$$

This equation can also be written as,

$$sm_{ij} = \frac{(cm_{ij} - r_i c_j)}{\sqrt{r_i c_j}} \qquad (2.33)$$

In matrix form, SM can be written as :

$$SM = D_r^{-\frac{1}{2}}(CM - rc^T)D_c^{-\frac{1}{2}} \qquad (2.34)$$

This matrix is similar to the Chi-squared matrix which represents the weighted departure of the original dataset from total independence. It may also be treated as the measure of weighted distance from the centroid in terms of rows and columns.

Step 4: The Chi-squared matrix is then subjected to singular value decomposition.

$$SM = LD_\mu M^T \qquad (2.35)$$

The SVD signifies an optimization problem where the orientations of the axes are obtained at the most reduced weighted distance from the cloud of row points and column points simultaneously. The sum of the squared values along the diagonal of $D_\mu$ represents the inertia of the cloud. The inertia is a term derived from the 'moment of inertia' and may be considered as the total mass of the weighted distance for the row or column cloud from the centroid. The calculation of the inertia along each principal axis (direction) is given by,

$$In(j^{th} \text{ axis}) = \frac{\mu_j^2}{\sum_{j=1}^{J} \mu_j^2} \qquad (2.36)$$

$$Aa = D_r^{\frac{1}{2}}L \qquad (2.37)$$

$$Bb = D_c^{\frac{1}{2}}M \qquad (2.38)$$

where, Aa and Bb represent the principal axes (loadings) of the columns and rows.

26

Step 5: Choice of number of components:

The number of components is usually chosen when the cumulative inertia values are found to exceed 80% in the same manner as the CPV calculations in equation (2.4) where the eigenvalues are replaced by the squares of the singular values from the diagonal of $D_\mu$ . Thus, in this manner, A components are chosen.

Step 6: Calculation of row and column scores.

The coordinates (scores) of the row cloud and column cloud for the new principal axis can be computed by projection on the first A columns of Aa and Bb.

$$FF = D_r^{-1}(Aa)D_\mu \tag{2.39}$$

$$GG = D_c^{-1}(Bb)D_\mu \tag{2.40}$$

where, FF and GG are the scores of the row cloud and the column cloud.

It must be noted that as both rows and column profiles have been considered in the SVD of the problem, the principal axes is used to show both the row cloud and column cloud on the same plot, hence these graphs are called bi-plots. These bi-plots are known to reveal useful information on the dependencies in the row, column and joint row-column space (Detroja *et al.,* 2006).

Step 7: Monitoring scheme for CA.

The monitoring scheme for Correspondence Analysis in FDD was developed by Detroja *et al.* (2007). In this procedure, a new sample xx i.e. $xx = [xx_1 \ xx_2 \ xx_3 \ ... \ xx_J \,]'$, can have its score calculated as:

$$r_{sample} = \sum_{j=1}^{J} xx_j \tag{2.41}$$

$$ff = [\frac{1}{r_{sample}} xx^T (GG) D_\mu^{-1}]' \tag{2.42}$$

where, $r_{sample}$ is the row sum for the current sample and ff is the score for the current sample.

The limits for the $T^2$ and Q statistics are calculated in the same was as in equations (2.6) and (2.12) except for the replacement of the eigenvalues by the square of the singular values in CA. The $T^2$ and Q statistics for CA are calculated as follows:

$$t^2 = (ff)' D_\mu^{-1} (ff) \tag{2.43}$$

$$res = (Bb)ff - (\frac{1}{r_{sample}} xx - c) \tag{2.44}$$

$$q = res' res \tag{2.45}$$

where, res is the residual vector for the sample.

## 2.3.2 Advances in CA

CA was applied quite recently in the field of FDD by Detroja *et al*. (2006). However, much before this, the method had been identified as a powerful multivariate tool in the field of categorical data analysis due to its abilities such as simultaneous analysis, graphical representation and flexibility in requirements. It has therefore been quickly adopted into several fields of study such as archeology (Baxter, 1994; Clouse, 1999), marketing research (Carroll *et al.,* 1989), ecology (ter Braak, 1987) and the social sciences (Clausen, 1998). An extension of simple correspondence analysis is Multiple Correspondence Analysis which refers to more than a couple of categorical variables.

Over the past few decades, CA has also been deeply analyzed by several researchers - many have tried to modify the method so that it can be adapted to interdisciplinary problems that have come about. Hill & Gauch Jr (1980) developed Detrended Correspondence Analysis (DCA). In this method, CA is performed as usual to obtain the principal axes but then, the first axis is divided into segments, and each segment is rescaled to have mean value of zero on the 2nd axis. This was found to be effective in removing a horse shoe curve where the first axes distort the second. Another method called Canonical Correspondence Analysis (CCA) was developed by ter Braak (1986) which conducts correspondence analysis by inducing the additional step of selecting the linear combination of row variables that maximizes the variation of the column scores. Greenacre developed what was called Joint Correspondence Analysis (JCA) which is considered a multiple correspondence analysis adjustment which can also be used for the analysis of two way contingency tables thus simplifying calculations. It was later improved by Boik (1996).

In the field of FDD, Detroja, *et al*. (2006 and 2007), had successfully applied CA to the quadruple tank system. Pushpa, *et al*. (2009) developed a polar classification procedure in which several faults are clustered after applying CA to a simulated dataset of a non-linear distillation column and experimental data from a quadruple tank system setup. Patel and Gudi (2009) have recently proposed a scheme to apply CA to penicillin fed batch fermentation process.

## 2.4 A Comparison between PCA and CA

CA has often been regarded as a form of PCA simultaneously performed for rows and columns (Jolliffe, 2002). It is known that PCA decomposes the covariance matrix to obtain a new set of axes. In geometrical terms, the covariance matrix is the Euclidean distance measure of n samples over an m-dimensional space. The same concept can also be noticed in CA where, the chi-

squared distance may be treated as a form of weighted Euclidean distance measure of the row and column cloud from a weighted centroid, where the weights correspond to the inverse of the row and column frequency sums for the respective row and column profiles (Detroja *et al*., 2006). Therefore, it is indicated that CA attempts to decompose a form of distance measure for both rows and columns of a dataset while PCA performs a similar type of decomposition for the columns of the dataset alone.

According to Detroja *et al*. (2007), CA has the advantage of analyzing dynamic data to a much better extent as compared to conventional and dynamic PCA. This can be seen in the fact that CA attempts to establish a relationship between rows and columns and in doing so can capture serial correlations in the dataset. PCA has the disadvantage of assuming independence of samples in its dataset while dynamic PCA has the need to create a data matrix of larger size to accommodate the same level of statistical significance. In Detroja *et al.* (2007), the authors applied CA to the Tennessee Eastman Challenge Process (Downs and Vogel, 1993) and successfully proved that CA possesses better detection and diagnosis capabilities as compared to both PCA and DPCA. This included superior features such as lower dimensional representation, higher detection rates and better diagnosis based on contribution plots. In consistency with the previous statements, CA can also be considered as a better tool than PLS which is again aimed to establishing a linear relationship between the inputs and outputs of the process yet again assuming independence of the samples. Thus, it can be concluded that CA is a superior multivariate tool which can be used for the fault detection and diagnosis in industrial processes where process dynamics is known to play a key role.

# 3. APPLICATION OF MULTIVARIATE TECHNIQUES TO SIMULATED CASE STUDIES

The following chapter will compare results regarding the fault detection and diagnosis of three systems, namely the quadruple tank system, Tennessee Eastman Challenge Process and the Depropanizer process. The first three sections will each begin with a description of the process followed by the tabulation and graphical representation of results. The results will contain the outcomes of using PCA, PLS and CA as detection and diagnosis tools. Detection is acknowledged by those data samples that exceed the 99% confidence limit of the $T^2$ or 95% confidence limit of Q statistics before and after the fault is introduced. The Q statistic is not employed while applying PLS as it is considered unsuitable for monitoring purposes as mentioned in Chapter 2. Diagnosis is performed with the aid of contribution plots for the various faults studied. . The contributions are calculated by first obtaining the aggregate for consecutive sets of six abnormal points detected. These aggregates are later used to obtain an overall contribution vector for the complete run. The last section will have an overall discussion on all the results arrived at earlier.

## 3.1     Quadruple Tank System

### 3.1.1   Process description

The quadruple tank process, as shown in Figure 3.1 is a multivariate process which is extensively used in the field of process control and monitoring as a test problem.  It was originally developed by Johansson (2000). This system consists of four interconnected water tanks, two pumps and

associated valves. The inputs to the system are the voltages supplied to the pumps $v_1$ and $v_2$ and the outputs are the water levels in the tanks $h_1$, $h_2$, $h_3$ and $h_4$. The flow to each tank is fixed using the associated valves $\gamma_1$ and $\gamma_2$ (range varies between 0 and 1), before each experiment.
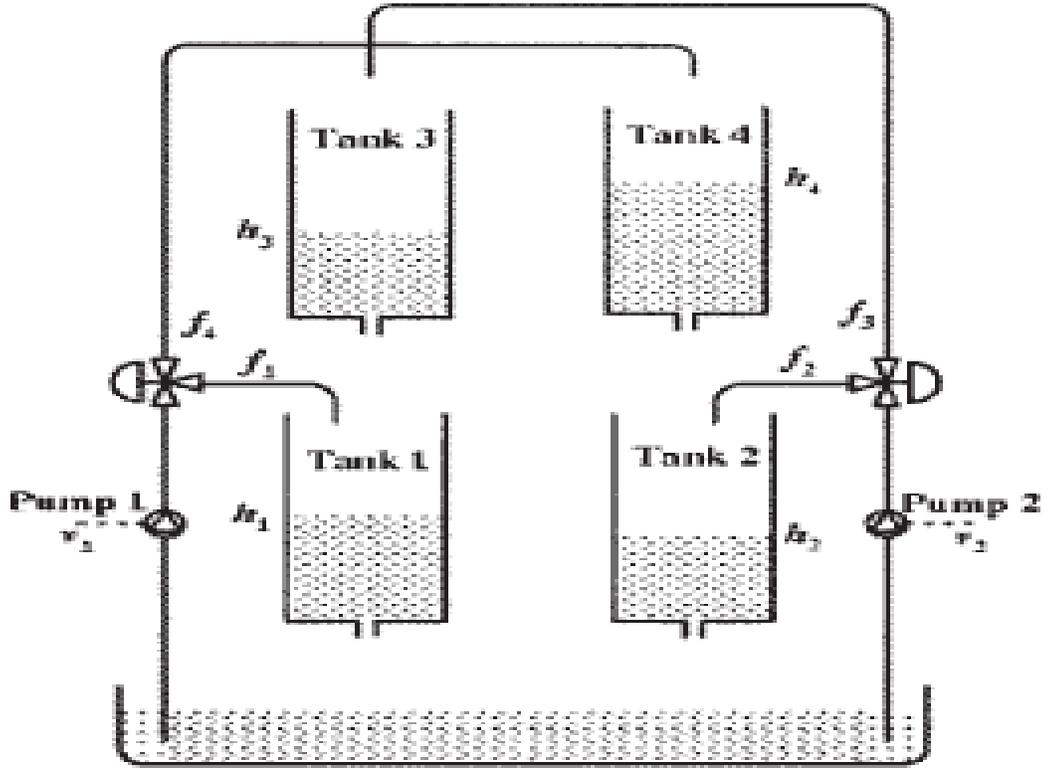


Figure 3.1: Quadruple Tank System

The equations of the non-linear model based on mass balances and Bernoulli's law are given as follows:

$$\frac{dh_1}{dt} = -\frac{a_1}{A_1}\sqrt{2gh_1} + \frac{a_3}{A_1}\sqrt{2gh_3} + \frac{\gamma_1 k_1}{A_1} v_1 \tag{3.1}$$

$$\frac{dh_2}{dt} = -\frac{a_2}{A_2}\sqrt{2gh_2} + \frac{a_4}{A_2}\sqrt{2gh_4} + \frac{\gamma_2 k_2}{A_2} v_2 \tag{3.2}$$

$$\frac{dh_3}{dt} = -\frac{a_3}{A_3}\sqrt{2gh_3} + \frac{(1-Y_2)k_2}{A_3} v_2 \tag{3.3}$$

$$\frac{dh_4}{dt} = -\frac{a_4}{A_4}\sqrt{2gh_4} + \frac{(1-Y_1)k_1}{A_4} v_1 \tag{3.4}$$

For each tank $i$, the Area is given by $A_i$. The cross section of the outer hole of each tank is $a_i$ and the voltages applied to each pump are given by $v_1$ and $v_2$ corresponding to the valves $\gamma_1$ and $\gamma_2$. The acceleration due to gravity is denoted by $g$. The flow rate to tank 1 is given by

$$f_1 = Y_1 k_1 v_1 \tag{3.5}$$

Similarly, flowrate to tank 2 is given by

$$f_2 = Y_2 k_2 v_2 \tag{3.6}$$

Then the flowrates to tanks 3 and 4 are,

$$f_3 = (1 - Y_2)k_2 v_2 \tag{3.7}$$

$$f_4 = (1 - Y_1)k_1 v_1 \tag{3.8}$$

The model of the quadruple tank system is simulated using SIMULINK in MATLAB. The level of the four tanks is controlled using two PID controllers which regulate the voltage values in the system. The set points for the two controllers are with respect to the heights of tank 1 and tank 2. The set points are referred to by variables h1_set and h2_set. A total of eight variables comprising the flow rate and heights of the four tanks are collected as data from the system. Gaussian white noise having a mean of zero and a standard deviation of 0.05 are added to the voltage values of $v_1$ and $v_2$ during the simulation thus corrupting the data generated with noise. The parameter values for the simulation are listed below in Table 3.1.

Table 3.1: Simulation parameters for the quadruple tank system

| Parameter | Unit | Value |
|---|---|---|
| $A_1$ ; $A_3$ | $cm^2$ | 28 |
| $A_2$ ; $A_4$ | $cm^2$ | 32 |
| $a_1$ ; $a_3$ | $cm^2$ | 0.071 |
| $a_2$ ; $a_4$ | $cm^2$ | 0.057 |
| $k_1$ ; $k_2$ | $cm^3 V^{-1} s^{-1}$ | 3.33 |
| $g$ | $cm\ s^{-2}$ | 981 |

The two major kinds of faults introduced in the system include sensor biasing and leakage of tanks. These faults have been introduced at different intensities and combinations to the system. Faults related to the sensor biasing of tanks is created by adding or deducting a fixed value from certain variables in the system. The leakage of tanks 1 & 2 is simulated by assuming that there are small holes at the bottom of each tank with areas $a_{leak1}$ and $a_{leak2}$. The equations 3.1 and 3.2 are replaced by the following equations in order to simulate the leakage.

$$\frac{dh_1}{dt} = -\frac{a_1}{A_1}\sqrt{2gh_1} + \frac{a_3}{A_1}\sqrt{2gh_3} + \frac{\gamma_1 k_1}{A_1} v_1 - \frac{a_{leak1}}{A_1}\sqrt{2gh_1} \tag{3.9}$$

$$\frac{dh_2}{dt} = -\frac{a_2}{A_2}\sqrt{2gh_2} + \frac{a_4}{A_2}\sqrt{2gh_4} + \frac{\gamma_2 k_2}{A_2} v_2 - \frac{a_{leak2}}{A_2}\sqrt{2gh_2} \tag{3.10}$$

The total number of variables used are 8 which are arranged in such a way that, one sample of the simulation would be given by $[f_1\ f_2\ f_3\ f_4\ h_1\ h_2\ h_3\ h_4]$. The normal operating condition is simulated for 350 samples with a sampling period of 5 seconds. The set points for the controllers

during operation are set at h1_set = 12.4 and h2_set = 12.7. The faults are simulated by introducing the fault after the 50$^{th}$ sample till a total of 400 data samples. The list of faults simulated along with their description is provided in Table 3.2. Fault 3 and fault 8 were simulated at slightly different operating conditions where, the set point h1_set was changed from 12.4 to 12.5. This was done to study the effect that such a change would have on the detection ability of the methods as such would be the case in an actual plant.

Table 3.2: Description of faults simulated for the Quadruple tank system

| Fault no. | Description | Important values |
|---|---|---|
| 1 | Leakage in tank 1 alone | $a_{leak1} = 0.005$ |
| 2 | Leakage in tank 2 alone | $a_{leak2} = 0.005$ |
| 3 | Negative sensor bias in height of tank 1 | $h1_{set} = 12.5, \Delta h_1 = 0.4$ |
| 4 | Negative sensor bias in height of tank 2 | $\Delta h_2 = 0.4$ |
| 5 | Simultaneous leakage in tank 1 & 2 | $a_{leak1} = a_{leak2} = 0.025$ |
| 6 | Leakage in tank 1 alone at low $a_{leak1}$ value | $a_{leak1} = 0.002$ |
| 7 | Positive sensor bias in height of tank 1 | $\Delta h_1 = 0.4$ |
| 8 | Positive sensor bias in height of tank 2 | $h1_{set} = 12.5, \Delta h_2 = 0.4$ |

The data generated for the normal operating condition and faults is then subjected to detection tests using PCA, CA and PLS. In PLS testing, the four flow rates are treated as the inputs and the heights of the respective tanks are taken as the outputs.

### 3.1.2 Results

The models obtained using PCA, PLS and CA are shown in Figures 3.2 to 3.7. The results for specific faults are shown from Figures 3.8 onwards. Table 3.3 displays all the values for the detection rates (DR) and false alarm rates (FAR) for all the faults based datasets. Detection delays involved in using each of the methods are shown in Table 3.4.
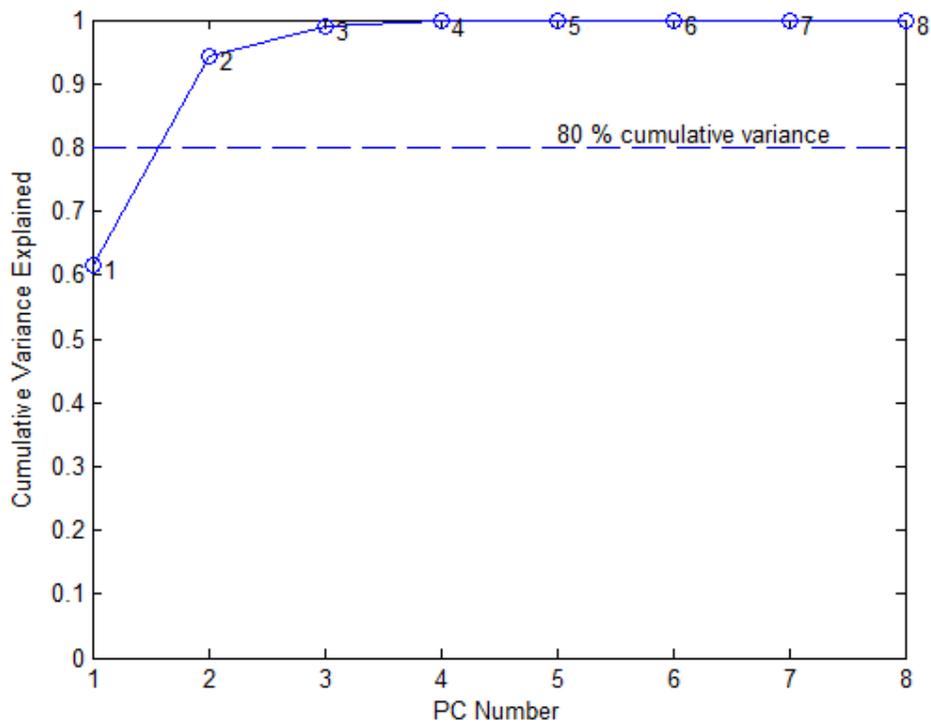


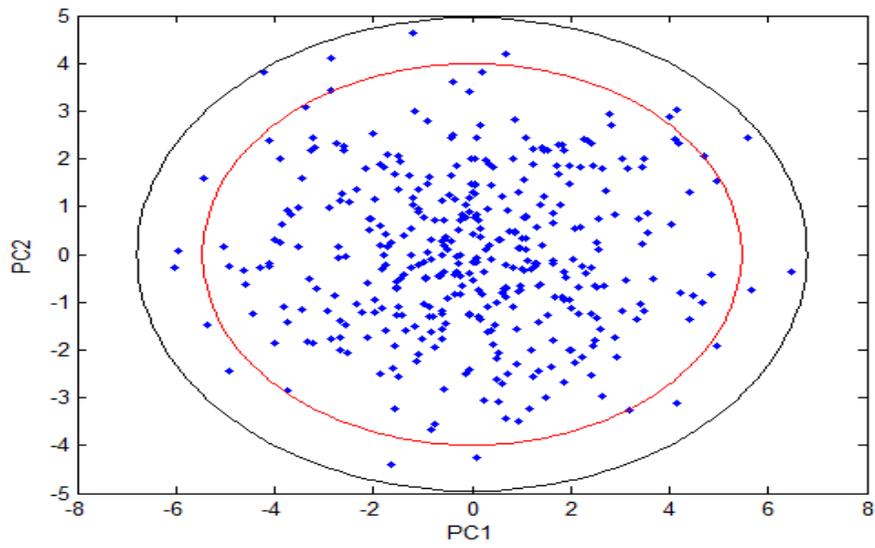Figure 3.2: Cumulative variance explained in the PCA model - Quadruple Tank system

Figure 3.3: PCA scores plot for first two PCs - Quadruple Tank system
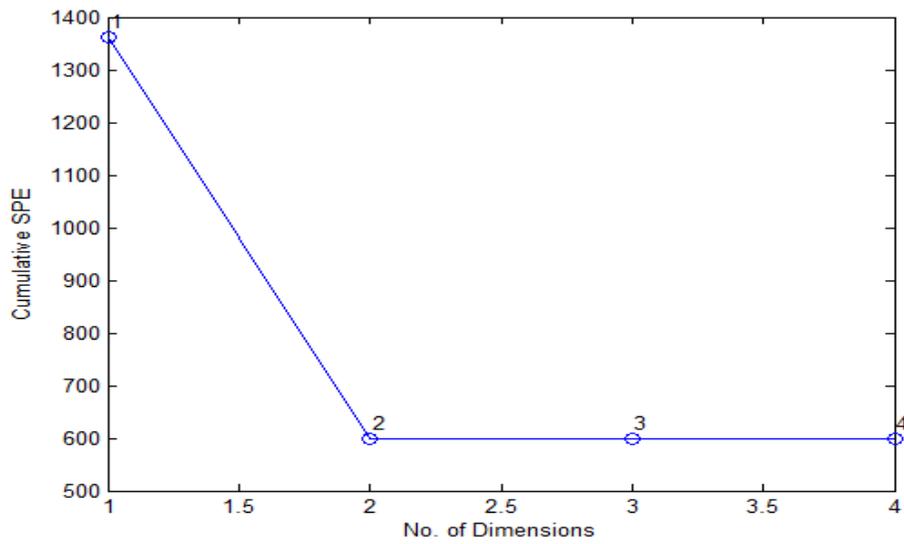


Figure 3.4: PLS cross validation to choose the number of PCs - Quadruple Tank system
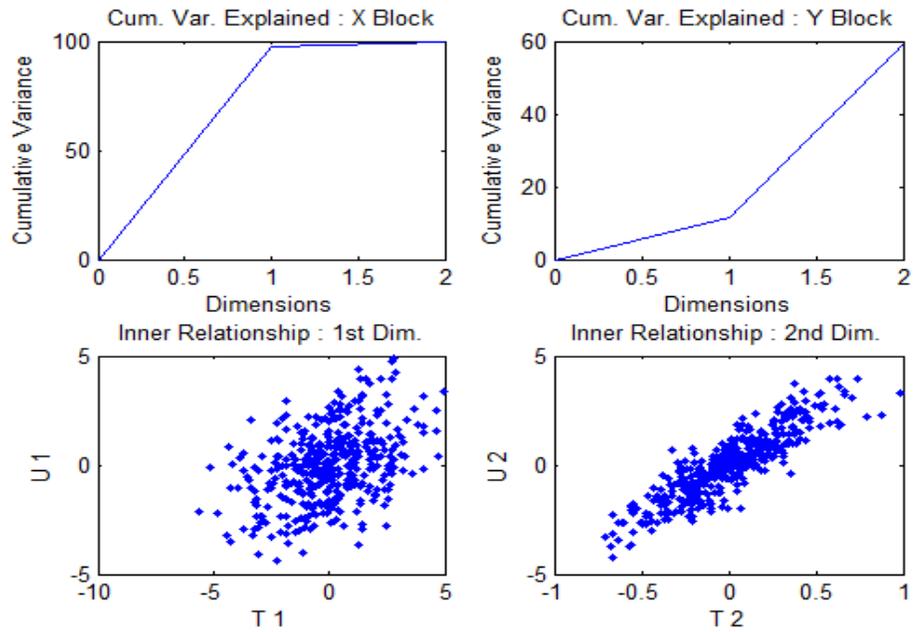
Figure 3.5: PLS Cumulative input-output relationships for first two PCs- Quadruple Tank system
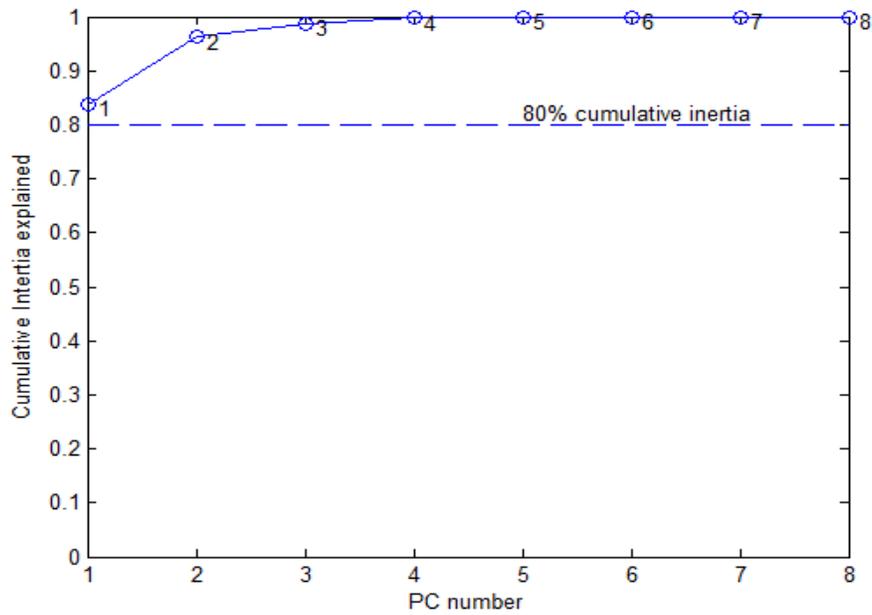


Figure 3.6: Cumulative Inertia explained by each PC in the CA model- Quadruple Tank system
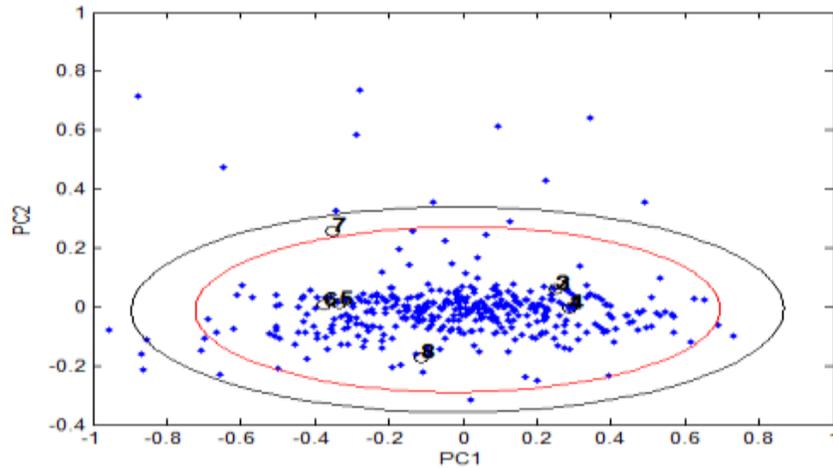
Figure 3.7: CA row and column scores bi- plot for first two PCs- Quadruple Tank system

In PCA, it is clear from Figure 3.2 that the first two PCs which explain about 95% of the variance are good enough to develop a model of the system. PLS uses the leave one out cross validation technique to choose the number of dimensions and according to Figure 3.4, the number of PCs required is 2. It is also clear from Figure 3.5 that the first two components alone account for 100 % of the variance in the input matrix X explaining about 60% of the variance in the output matrix Y. Therefore, it would not be possible for the model to use the Q statistic for the inputs in the analysis due to the extremely negligible amount of variance involved in the residual space. In Figure 3.6, the first two PCs for the CA model account for 97% of the inertia in the system. Although one cannot draw a clear comparison between inertia and variance, it is proper to state that both PCA and CA capture most of the information in the system with their first two PCs. Figure 3.7 shows the bi-plot developed by CA, where the blue dots denote the row scores and the black squares are the column scores. The bi-plot will be useful in graphically understanding the relationship between the rows and columns. But, for the sake of monitoring purposes, one can only use the row scores to develop a confidence region. Both Figures 3.3 and
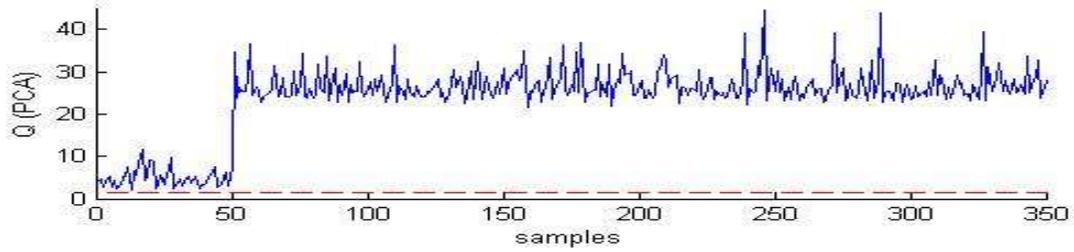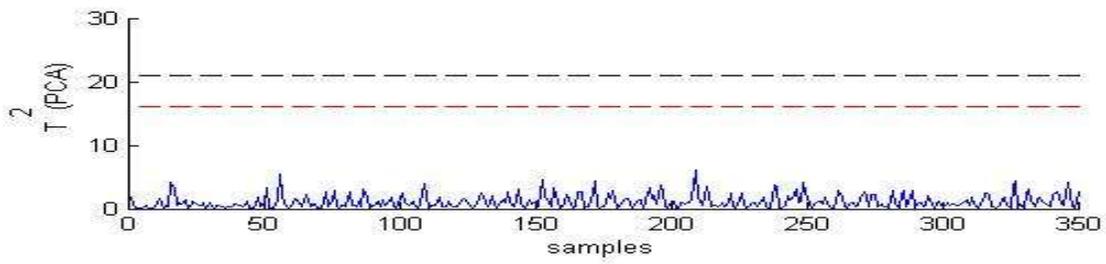
3.7 have confidence ellipses to isolate the zone of normal operation where the red ellipse refers to the area with a 95% confidence limit and the black ellipse refers to the area with a 99% confidence limit.

Table 3.3: Detection rates and false alarm rates – Quadruple tank system

| Faults | DR | | | | | FAR | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PCA | | PLS | CA | | PCA | | PLS | CA | |
| | $T^2$ | $Q$ | $T^2$ | $T^2$ | $Q$ | $T^2$ | $Q$ | $T^2$ | $T^2$ | $Q$ |
| 1 | 0.0033 | 0.9967 | 0.98 | 0.9734 | 0.0100 | 0 | 0.2040 | 0 | 0 | 0 |
| 2 | 0.0033 | 0.9967 | 0.9867 | 0.9867 | 0.9867 | 0 | 0.2040 | 0 | 0 | 0 |
| 3 | 0 | 1 | 0 | 0.0598 | 0.9967 | 0 | 0.8775 | 0 | 0 | 0 |
| 4 | 0 | 1 | 0 | 0.1628 | 0.9967 | 0 | 0.2040 | 0 | 0 | 0 |
| 5 | 0.9900 | 0.9967 | 0.9933 | 0.9502 | 0.9934 | 0 | 0.2040 | 0 | 0 | 0 |
| 6 | 0 | 0.9967 | 0.3567 | 0.3389 | 0.0033 | 0 | 0.2040 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 0 | 0.9967 | 0 | 0.2040 | 0 | 0 | 0 |
| 8 | 0 | 1 | 0 | 0 | 0.9967 | 0 | 0.8775 | 0 | 0 | 0 |

Table 3.4: Detection delays (in seconds) – Quadruple tank system

| Faults | PCA | PLS | CA |
|---|---|---|---|
| 1 | 5 | 0 | 10 |
| 2 | 5 | 0 | 10 |
| 3 | 0 | 0 | 5 |
| 4 | 0 | 0 | 5 |
| 5 | 5 | 0 | 10 |
| 6 | 5 | 0 | 10 |
| 7 | 0 | 0 | 5 |
| 8 | 0 | 0 | 5 |

a) PCA analysis results



b) CA analysis results



c) PLS analysis results

Figure 3.8: Fault 3 results – Quadruple tank system

a) PCA analysis results



b) CA analysis results



c) PLS analysis results

Figure 3.9: Fault 6 results – Quadruple tank system

a)  PCA analysis results



b)  CA analysis results



c)  PLS analysis results

Figure 3.10: Fault 8 results – Quadruple tank system

Table 3.5: Contribution plots with PCA and CA analysis – Quadruple tank system

| Faults | PCA | CA |
|--------|-----|-----|
| Fault 1 | | |
| Fault 2 | | |
| Fault 3 | | |
| Fault 4 | | |
| Fault 5 | | |
| Fault 6 | | |
| Fault 7 | | |
| Fault 8 | | |



Figures 3.8, 3.9 and 3.10 show the fault detection results for faults 3,6 and 8 while Table 3.5 contains the contribution plots for all faults where, the variable 1, 2, 3, 4, 5, 6, 7 and 8 correspond to $f_1$, $f_2$, $f_3$, $f_4$, $h_1$, $h_2$, $h_3$, $h_4$.

In the results, faults 1, 2 and 5 which were related to the leakage in tanks 1 and 2 or both together were mostly detected by the $T^2$ statistic in the case of CA and PLS while it was more properly

detected by the $Q$ statistic in the case of PCA. This shows that the CA model structure was able to understand the relationship between the variables to a much better extent due to its visualization in a weighted space while the right choice of predictor and response variables in PLS helped establish a proper regression model. PCA has to depend on the residual statistic to understand the anomaly. Faults 3, 4, 7 and 8 were all related to sensor bias in tanks 1 and 2 and were well detected by both CA and PLS with very mild differences in detection rates. But, the use of slightly different operating conditions in faults 3 and 8 immediately displayed the fact that the PCA model is quite rigid and time-invariant. One can notice that both these faults recorded the value of 0.87 as false alarm rates shown clearly in Figures 3.8 and 3.10 and is therefore incapable of proper detection, whereas CA is found to not record any such false alarms at all thus displaying its ability to understand the dynamics of the process and remain flexible. The only negative point in terms of fault detection was fault 6 as shown in Figure 3.9 where the leakage in tank 1 was too mild to detect for CA and PLS while PCA was able to do so effortlessly. This can be attributed to the fact that the Q statistic in PCA was able to pick up the slight modification to the model structure in its residuals while CA's Q statistic was influenced and distorted by the cross-tabulation interaction between the rows and columns of the model's original dataset. In regard to PLS, the same could be said where the statistic was not able to identify the mild change in the relationships between the input and the output. The only silver lining even in this fault's analysis is that once again the $T^2$ statistic of CA and PLS was able to perform much better than that of PCA.

With regard to the fault diagnosis capabilities of PCA and CA in terms of their contribution plots, one can see in Table 3.5 that both methods were able to provide accurate information on the major variables related to the sensor bias based faults i.e. faults, 3, 4, 7 and 8. When it came

to faults 1, 2, 5 and 6, which were based on leakage of tanks 1 and 2, the results were not accurate. According to equations 3.1 – 3.4, 3.9 and 3.10, leakages caused to tanks 1 or 2 would tend to change the voltage values of $v_1$ and $v_2$ to regulate control. The same voltage values will also be used to regulate the flow to tanks 3 and 4 changing their values in the process, hence in the case of the contribution plots for 1, 2, 5 and 6, all four values rise to different values and would thus exhibit some conflicting values in variables 5, 6, 7 and 8 in the bar plots. In the case of PCA, the variables 7 and 8 corresponding to $h_3$ and $h_4$ show significant values thus proving that these variables carry more weightage in the model as compared to others. In CA, although variable 8 corresponding to $h_4$ is found to have higher contribution as compared to other variables in faults 1,2 and 5, the issue of conflicting values can be confirmed by seeing the bar plots in Table 3.5. The only difference in diagnosis turned out to be fault 6 which was properly diagnosed by CA; this could be due to the fact that the few samples detected by CA (detection rate – 0.3389) could have understood the actual dynamics of the abnormality and provided an accurate estimate.

## 3.2    Tennessee Eastman Process (TEP)

### 3.2.1   Process description

The Tennessee Eastman Process (Downs and Vogel, 1993) is a popular benchmark problem used in the field of process control and fault detection. It is based on a real chemical process plant where the components, kinetics and operating conditions were modified for proprietary reasons. As shown in Figure 3.11, the process consists of five major unit operations: the reactor, a product condenser, a vapor-liquid separator, a recycle compressor and a product stripper. The process consists of 12 manipulated variables from the controller and 41 process measurements. Gaseous

reactants A, C, D, E and an inert B are fed to the reactor. They react to form the liquid products G and H and other byproducts while gas phase reactions in the same are catalyzed by a non-volatile catalyst dissolved in the liquid phase. The products



Figure 3.11: Tennessee Eastman Challenge Process

stream from the reactor then passes through the condenser for condensation of products and then flows to the vapor-liquid separator. Here, the non-condensed components recycle back through a centrifugal compressor to the reactor feed. Condensed components move to a product stripping column to remove remaining reactants by stripping with feed stream number 4. The required products G and H exit the stripper base and are collected separately.

Table 3.6: Process faults: Tennessee Eastman Process

| Fault | Description | Type |
|---|---|---|
| IDV(1) | A/C Feed ratio, B composition constant (Stream 4) | Step |
| IDV(2) | B Component, A/C ratio constant (Stream 4) | Step |
| IDV(3) | D Feed temperature (Stream 2) | Step |
| IDV(4) | Reactor cooling water inlet temperature | Step |
| IDV(5) | Condenser cooling water inlet temperature | Step |
| IDV(6) | A Feed loss (Stream 1) | Step |
| IDV(7) | C Header pressure loss–reduced availability (Stream 4) | Step |
| IDV(8) | A, B, C Feed component (Stream 4) | Random |
| IDV(9) | D Feed temperature (Stream 2) | Random |
| IDV(10) | C Feed temperature (Stream 4) | Random |
| IDV(11) | Reactor Cooling Water Inlet temperature | Random |
| IDV(12) | Condenser Cooling Water Inlet temperature | Random |
| IDV(13) | Reactor kinetics | Slow drift |
| IDV(14) | Reactor Cooling Water valve | Sticking |
| IDV(15) | Condenser Cooling Water valve | Sticking |
| IDV(16) | Unknown | - |
| IDV(17) | Unknown | - |
| IDV(18) | Unknown | - |
| IDV(19) | Unknown | - |
| IDV(20) | Unknown | - |
| IDV(21) | The valve for stream 4 was fixed at the steady state position | Constant position |

Source: Detroja *et al.* (2007).

The TEP simulation setup has a total of 21 pre-programmed process faults. From Table 3.6, it can be seen that faults IDV(1) – IDV(15) and IDV(21) are of a known nature and the rest are not. Of those faults, IDV(1) – IDV(7) are related to a step change in process variables while IDV(8) – IDV(12) are involved in the random variability of certain process variables. IDV(13) is influenced by a slow drift in the reaction kinetics and IDV(14), IDV(15) and IDV(21) are associated with sticking valves. The datasets for the system was obtained from the website http://brahms.scs.uiuc.edu (link is no longer functional). The datasets obtained were generated using the control structure recommended by Lyman and Georgakis (1995). The data comprised of testing and training datasets for the normal operating condition and the 21 faults. Each training dataset had 480 to 500 samples collected at an interval of three minutes each for 52 variables (the manipulated variable related to the speed of the stirrer in the reactor was not recorded) with the fault being introduced at the $20^{th}$ sample. The testing sets contained 960 samples with the fault being introduced at the $160^{th}$ sample. Only 34 (23 process and 11 manipulated variables) variables out of the total 53 (41 process and 12 manipulated variables) are used in the simulation runs. About 22 of the 23 process variables used along with the 11 manipulated variables are continuous process measurements such as temperatures, pressures, levels, flow rates, work rates and speeds which are usually available in a real plant. The remaining 19 process measurements are related to component analysers at various points in the process which are measured at discrete intervals of 6 to 15 min. Of these 19 measurements the analyser value related to component G in stream 9 alone is chosen to act as the quality variable. The main reason for choosing such a combination of variables is to mimic the pragmatic nature of plants where continuous measurements would be available easily. Faults IDV(3), IDV(9) and IDV(15) will be neglected in the final results as they were found to show very low or negligent detection rates.

This was confirmed by Russell *et al*. (2000) when all 52 variables were used to obtain the results for PCA. The authors had stated that no observable change in the mean or the variance could be detected by visually comparing the plots of each associated observation variable in these faults.

### 3.2.2   Results

PCA, PLS and CA models obtained using the training datasets are shown in Figures 3.12 - 3.17. The detection rates, detection delays and diagnosis results are tabulated in Tables 3.7 to 3.9. In the case of this system, main contribution variables alone will be mentioned in Table 3.10 as there are a large number of variables and faults to provide a detailed explanation for all of them. The main contribution variables are chosen as those that exceed a value of greater than or equal to 5%. The contribution variables obtained in PCA and CA will then be compared for analysis.
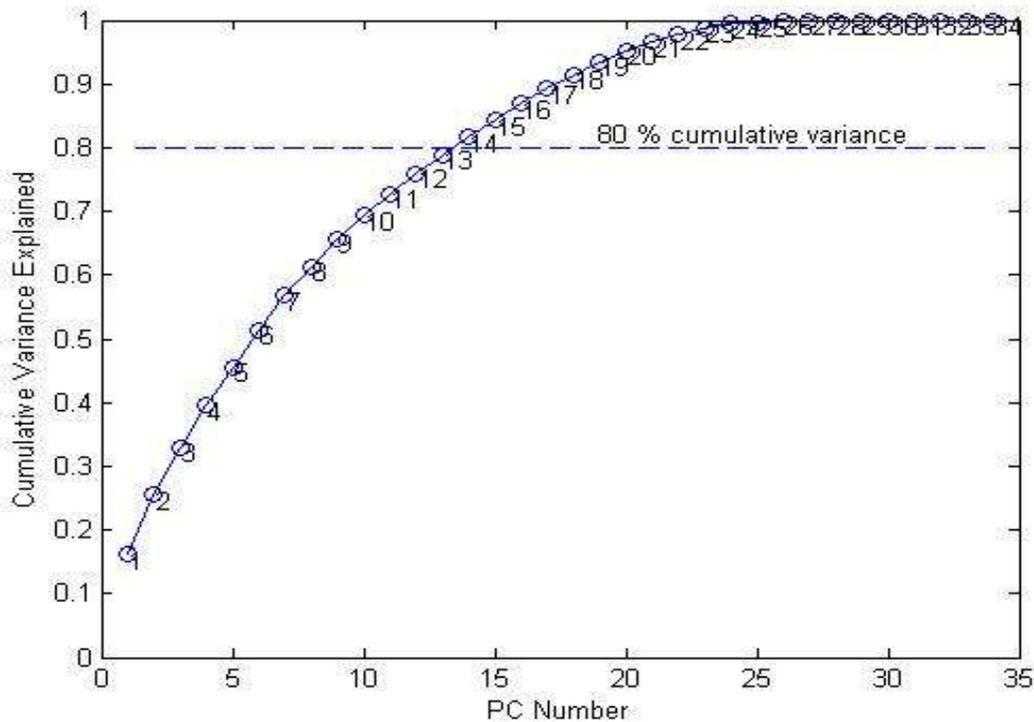


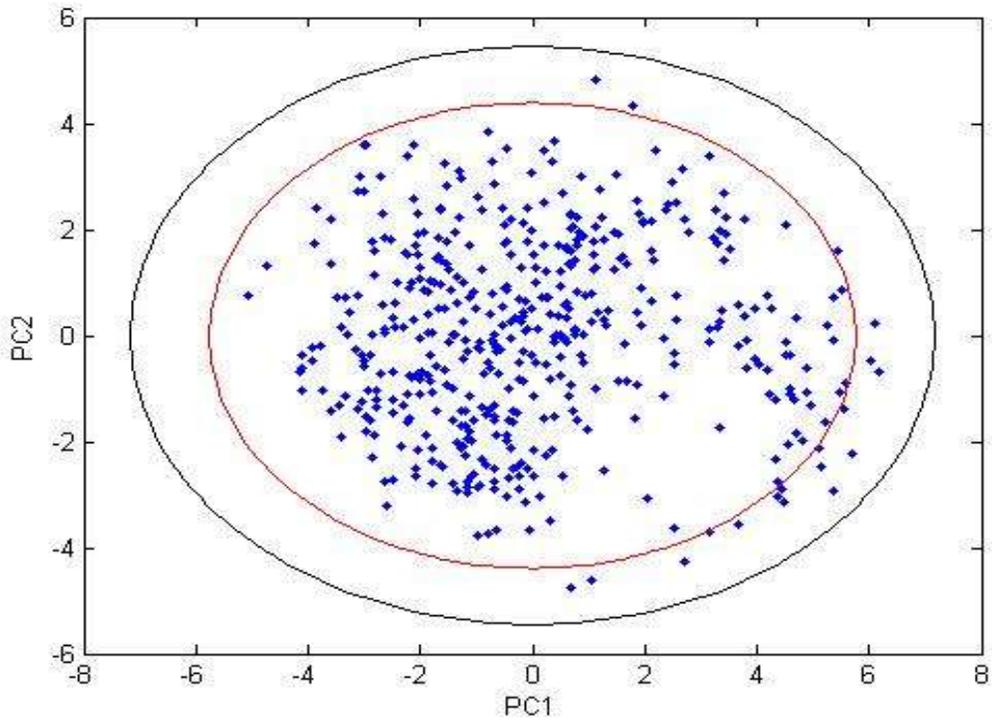Figure 3.12: Cumulative variance explained in the PCA model - TEP

Figure 3.13: PCA scores plot for first two PCs - TEP



Figure 3.14: PLS cross validation to choose the number of PCs - TEP

Figure 3.15: PLS Cumulative input-output relationships for first 12 PCs- TEP



Figure 3.16: Cumulative inertia explained in the CA model - TEP

Figure 3.17: CA scores bi-plot for first two PCs - TEP

From Figure 3.12, it is clear that about 14 components are required to represent a cumulative variance in excess of 80 % for the PCA model while in Figure 3.16 only about 6 components were required to obtain a cumulative inertia in excess of 80 %. In order to avoid having to compare the physical significance of variance with that of inertia, we will be choosing a total of 15 components for both the PCA and CA models. Fifteen components in the PCA model were found to account for 84.53 % of the variance in the system while the same number of components was found to represent 98.86% of the inertia in the system. In the case of PLS, 12 components were chosen for detection purposes based on the cross validation diagram given in Figure 3.14.

Table 3.7: Detection rates and false alarm rates – Tennessee Eastman Process

| Faults | DR | | | | | FAR | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PCA | | PLS | CA | | PCA | | PLS | CA | |
| | $T^2$ | $Q$ | $T^2$ | $T^2$ | $Q$ | $T^2$ | $Q$ | $T^2$ | $T^2$ | $Q$ |
| IDV(1) | 0.9850 | 1 | 0.9950 | 0.9850 | 0.9513 | 0 | 0.3563 | 0 | 0 | 0.0063 |
| IDV(2) | 0.9725 | 0.9950 | 0.9787 | 0.9775 | 0.9850 | 0 | 0.2500 | 0 | 0 | 0 |
| IDV(4) | 0.0013 | 1 | 0.4150 | 0.2900 | 0.9463 | 0 | 0.4125 | 0 | 0 | 0.0063 |
| IDV(5) | 0.1513 | 0.8363 | 0.2225 | 0.2125 | 0.9988 | 0 | 0.4125 | 0 | 0 | 0.0063 |
| IDV(6) | 0.9800 | 1 | 0.9900 | 0.9875 | 1 | 0 | 0.3500 | 0 | 0 | 0.0063 |
| IDV(7) | 0.9800 | 1 | 1.000 | 1 | 0.5800 | 0 | 0.3500 | 0 | 0 | 0.0063 |
| IDV(8) | 0.8488 | 1 | 0.9675 | 0.9538 | 0.9125 | 0 | 0.5563 | 0 | 0 | 0 |
| IDV(10) | 0 | 0.9388 | 0.7362 | 0.1538 | 0.5133 | 0 | 0.4438 | 0 | 0 | 0.0063 |
| IDV(11) | 0.1275 | 0.9650 | 0.4050 | 0.2963 | 0.5663 | 0 | 0.4688 | 0 | 0 | 0.0125 |
| IDV(12) | 0.8050 | 1 | 0.9775 | 0.9638 | 0.9425 | 0 | 0.5438 | 0 | 0 | 0 |
| IDV(13) | 0.8450 | 0.9938 | 0.9412 | 0.9225 | 0.9525 | 0 | 0.2625 | 0 | 0 | 0 |
| IDV(14) | 0.7888 | 1 | 0.9987 | 0.7438 | 1 | 0 | 0.4125 | 0 | 0 | 0.0188 |
| IDV(16) | 0 | 0.9600 | 0.5375 | 0.0525 | 0.7638 | 0 | 0.7250 | 0 | 0 | 0.0125 |
| IDV(17) | 0.5350 | 0.9825 | 0.7850 | 0.4775 | 0.7650 | 0 | 0.6125 | 0 | 0 | 0.0188 |
| IDV(18) | 0.8813 | 0.9688 | 0.8912 | 0.8825 | 0.9038 | 0 | 0.4375 | 0 | 0 | 0.0188 |
| IDV(19) | 0 | 0.8475 | 0.0562 | 0 | 0.4400 | 0 | 0.4000 | 0 | 0 | 0.0063 |
| IDV(20) | 0.0588 | 0.9325 | 0.3287 | 0.2450 | 0.5188 | 0 | 0.3625 | 0 | 0 | 0 |
| IDV(21) | 0.0813 | 0.8125 | 0.3400 | 0.2388 | 0.5650 | 0 | 0.6438 | 0 | 0 | 0.0188 |

Table 3.8: Detection delays (in minutes) – Tennessee Eastman Process

| Faults | PCA | PLS | CA |
|--------|-----|-----|-----|
| IDV(1) | 0 | 12 | 3 |
| IDV(2) | 12 | 45 | 36 |
| IDV(4) | 0 | 3 | 0 |
| IDV(5) | 0 | 6 | 3 |
| IDV(6) | 0 | 18 | 0 |
| IDV(7) | 0 | 3 | 0 |
| IDV(8) | 0 | 63 | 39 |
| IDV(10) | 0 | 75 | 81 |
| IDV(11) | 0 | 21 | 15 |
| IDV(12) | 0 | 9 | 6 |
| IDV(13) | 0 | 129 | 111 |
| IDV(14) | 0 | 6 | 0 |
| IDV(16) | 0 | 39 | 27 |
| IDV(17) | 6 | 78 | 60 |
| IDV(18) | 0 | 261 | 45 |
| IDV(19) | 0 | 33 | 3 |
| IDV(20) | 0 | 258 | 255 |
| IDV(21) | 0 | 792 | 276 |

a) PCA analysis results



b) CA analysis results



c) PLS analysis results

Figure 3.18: IDV(16) results – TEP

From the results provided in Tables 3.7 and 3.8, about 11 faults in the TEP were detected with detection rate that is greater than 0.9 while the same was achieved by 15 faults in PCA and 9 faults in PLS. All three methods were able to detect most of the faults created by step input in the variables while CA was unable to detect the faults IDV(10) and IDV(11) as compared to PCA (which still had a detection rate greater than 0.9 in these cases) and PLS (which fared better than CA in the case of IDV(10) alone). The false alarms rates were recorded for the $T^2$ and $Q$ statistics of all three methods. The PCA method recorded high false alarm rates for all the faults and they were found to lie within a range of 0.25 to 0.72 while CA and PLS recorded negligible values. The detection delays recorded in terms of minutes indicated that both CA and PLS have high values of detection delays as compared to PCA. Only IDV(4), IDV(6), IDV(7) and IDV(14) were found to give zero time delays for CA and were comparable to PCA. IDV(10), IDV(13), IDV(17) and IDV(21) were found to have an excess time delay greater than 50 minutes as compared to PCA while IDV(18) along with the previous mentioned faults was found to have similar excessive time delays in PLS as compared to PCA. IDV(13) related to the slow changing kinetics of the process, IDV(20) which is of an unknown nature and IDV(21) related to the constant position of valves were found to have the highest time delay values going into three digit Figures. Comparison between CA and PLS in terms of detection delays indicated that CA fared to a slightly better extent than PLS in most of the cases.

Table 3.9: Tennessee Eastman Process

| Variable number | Variable reference in TEP |
| --- | --- |
| 1 | XMEAS(1) |
| 2 | XMEAS(2) |
| 3 | XMEAS(3) |
| 4 | XMEAS(4) |
| 5 | XMEAS(5) |
| 6 | XMEAS(6) |
| 7 | XMEAS(7) |
| 8 | XMEAS(8) |
| 9 | XMEAS(9) |
| 10 | XMEAS(10) |
| 11 | XMEAS(11) |
| 12 | XMEAS(12) |
| 13 | XMEAS(13) |
| 14 | XMEAS(14) |
| 15 | XMEAS(15) |
| 16 | XMEAS(16) |
| 17 | XMEAS(17) |
| 18 | XMEAS(18) |
| 19 | XMEAS(19) |
| 20 | XMEAS(20) |
| 21 | XMEAS(21) |
| 22 | XMEAS(22) |
| 23 | XMEAS(35) |
| 24 | XMV(1) |
| 25 | XMV(2) |
| 26 | XMV(3) |
| 27 | XMV(4) |
| 28 | XMV(5) |
| 29 | XMV(6) |
| 30 | XMV(7) |
| 31 | XMV(8) |
| 32 | XMV(9) |
| 33 | XMV(10) |
| 34 | XMV(11) |

Table 3.10: High fault contribution variables - Tennessee Eastman Process

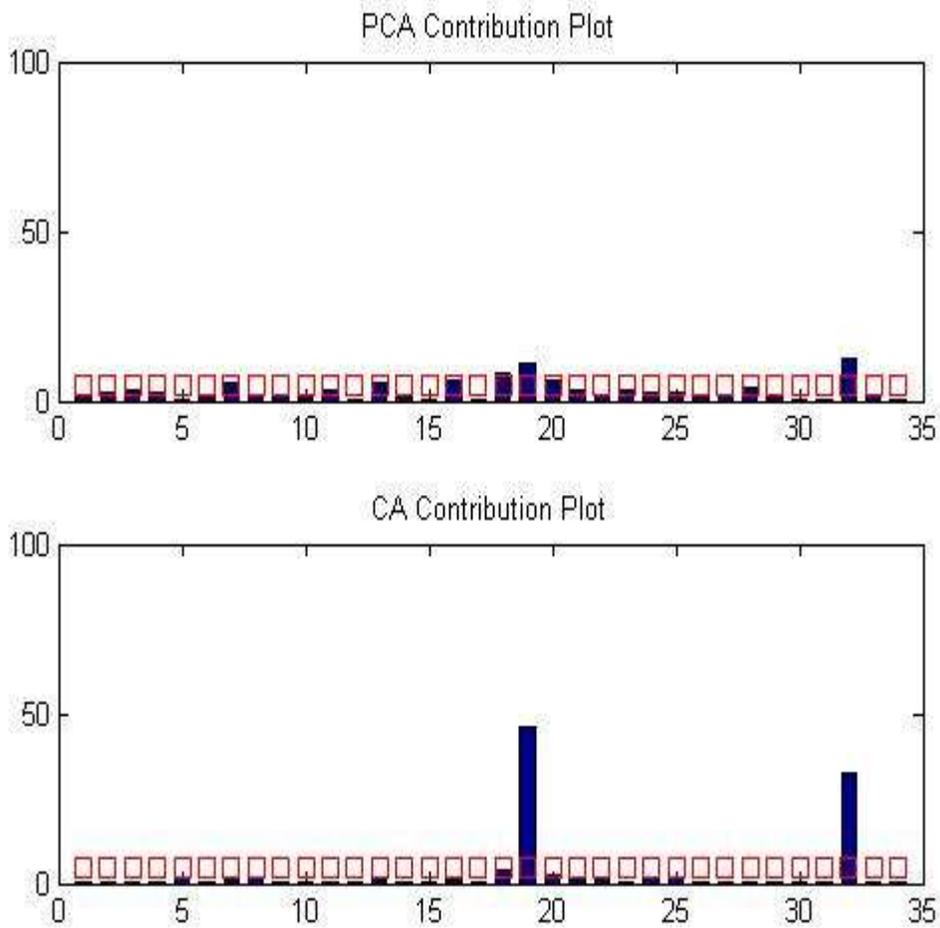| Faults | PCA | CA |
|---|---|---|
| IDV(1) | 1, 3, 4, **18**, **21**, 26 | 8, **18**, 19, 20, **21**, 32 |
| IDV(2) | 10, **11**, 16, **19**, 29 | 7, **11**, **19**, 20, 21, 25 |
| IDV(4) | 9, **33** | 8, 21, 24, **33** |
| IDV(5) | 7, 13, 16, 20, **34** | 8, 17, **34** |
| IDV(6) | 7, 16, **20**, 28, 33 | 1, 17, **20**, 26 |
| IDV(7) | 4, 27 | 5, 7, 8, 11, 20, 22, 25 |
| IDV(8) | **7**, 11, **13**, **16**, **20**, 28 | **7**, 8, **13**, **16**, **20** |
| IDV(10) | 7, 13, 16, **18**, **19**, **20**, 32 | **18**, **19**, **20** |
| IDV(11) | **8**, 9, **33** | **8**, 21, 24, **33** |
| IDV(12) | **7**, **11**, **13**, **16**, **18**, **20**, 21 | **7**, 8, **11**, **13**, **16**, **18**, 19, **20** |
| IDV(13) | **7**, **13**, **16**, 18, **19**, **20**, **32** | **7**, **13**, **16**, **19**, **20**, **32** |
| IDV(14) | **8**, 9, **21**, **33** | **8**, **21**, 24, **33** |
| IDV(16) | 7, 13, 16, 18, **19**, 20, **32** | **19**, **32** |
| IDV(17) | 9, **21** | 8, **21** |
| IDV(18) | 7, 13, 16, 24, **25**, **28**, 33 | 8, 17, 19, 21, 22, **25**, **28**, 32 |
| IDV(19) | 3, 16, 21, **28** | 5, 8, 13, 19, 20, 24, **28** |
| IDV(20) | **7**, 13, **16**, 20, **28** | **7**, 11, **16**, 22, **28** |
| IDV(21) | **7**, **8**, 11, **13**, **16**, 19 | **7**, **8**, **13**, **16**, 20 |

Figure 3.19: IDV(16) results – contribution plots - TEP

In fault diagnosis using contribution plots, there were 11 instances in which CA was found to be on par or show less number of main contribution variables as compared to PCA. Of the 11 faults where detection was greater than 0.9 in both PCA and CA, two faults (IDV(1) and IDV(14)) were found to show the same number of contribution variables which IDV(5), IDV(6), IDV(8) and IDV(13) were found to show less number of contribution variables. IDV(16) and IDV(17) with average detection rates exceeding 0.7 were also found to show more concrete diagnosis with CA. A good example of diagnosis by CA would be IDV(16) where the fault is of an unknown nature. CA indicates variables variables 19 and 32 which are XMEAS(19) and

XMV(9) which are both related to the stripper steam flow and indicate the problem to be there as compared to 7 main variables indicated in PCA which also show variables 19 and 32.

## 3.2 Depropanizer Process

### 3.2.1 Process description

The depropanizer unit consists of a fractionating column which is used to separate a mixture of $C_3$ and $C_4$ hydrocarbons so that the top product would yield the lighter $C_3s$ based hydrocarbons while the bottom of the unit will give the $C_4s$ and heavier hydrocarbons. The unit described here comprises of a 40 tray fractionation tower, a condenser, reflux drum and the reboiler. This process has a total of 36 variables that are monitored. The initial stage of the process involves the input containing the above mentioned mixture being fed to the middle of the bubble tray fractionating tower C11. The flow of this input is directly controlled by a flow controller FC11. During fractionation, the extent of separation is controlled by the tower temperature controller TC11. The tower bottom level is controlled by the tower bottom level controller LC11. LC11 controls the level by adjusting the bottom product draw.

After fractionation, the overhead vapors obtained are condensed in a shell and tube condenser E12 with cooling water, a condenser bypass valve is also present for regulating the pressure. The condensed liquid is then fed to the bottom of a horizontal vessel called the reflux drum while, vapours passing through the condenser bypass valve are directed to the top of the same vessel. The pressure in the tower is regulated by a pressure controller PC11. PC11 regulates the pressure of the tower by controlling the condenser bypass valve and an off-line gas valve. The off-line gas valve is part of the off-gas line connected to the top of the reflux drum. The condenser bypass valve is opened when the pressure is too low and closed when it's too high. If a situation arises

where the pressure cannot be maintained by closing the condenser bypass valve, the off-gas line valve is opened to let out vapors. The reflux and pumped by pumps P11A and P11B. Only one of the pumps is used during operation while the other is on stand-by, the discharge from this pump is then separated into two streams, the reflux stream and the top product stream. The flow of the reflux stream is controlled using a reflux flow controller FC12. The purpose of this controller is to maintain an optimum value of reflux flow to sustain the desired extent of separation and hence preserve product quality. The top product's flow is regulated by reflux drum level controller LC12. The top product is now collected separately or sent to the next stage in a wider process. The product from the bottom of the tower is vaporized using hot oil that is fed to the shell side of reboiler E11; the flow of this oil is regulated by the tower temperature controller TC11. This control of flow helps in maintaining the bottom temperature of the tower. The bottom product is pumped out with one of two pumps P12A or P12B and collected separately.

The simulation data is collected over a period of three hours for each of the normal and fault conditions and the samples are recorded at a regular interval of 12 seconds. A total of 15 faults were generated as shown in Table 3.11 with 901 samples were collected for the normal operating region as well as the faults. In the fault induced datasets, the fault has been introduced at the $51^{st}$ sample. The normal operating region was divided into testing and training datasets such that first 60% of the samples formed the training dataset and the remaining 40% formed the testing dataset.
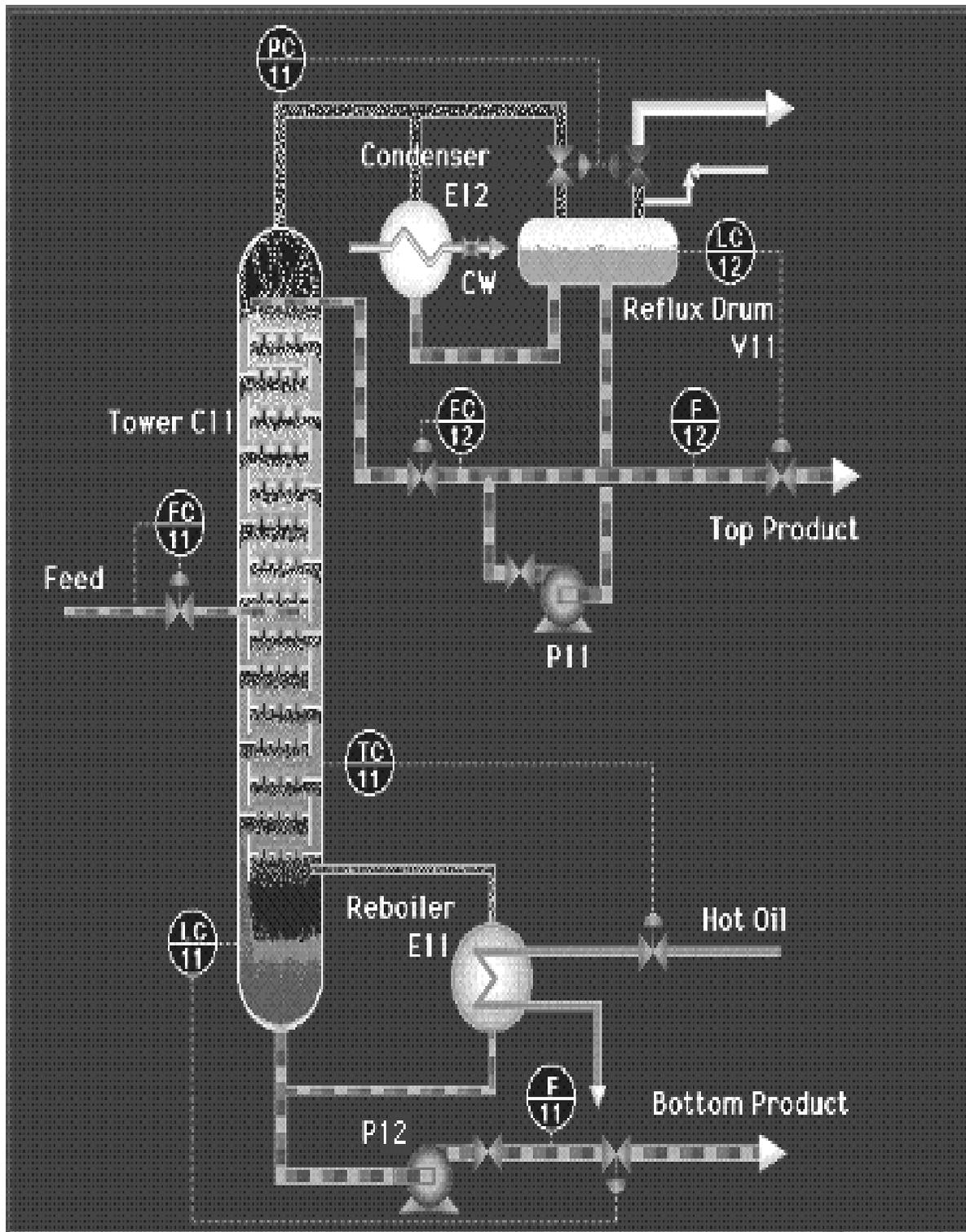
Figure 3.20: Depropanizer Process

Table 3.11: Process faults: Depropanizer Process

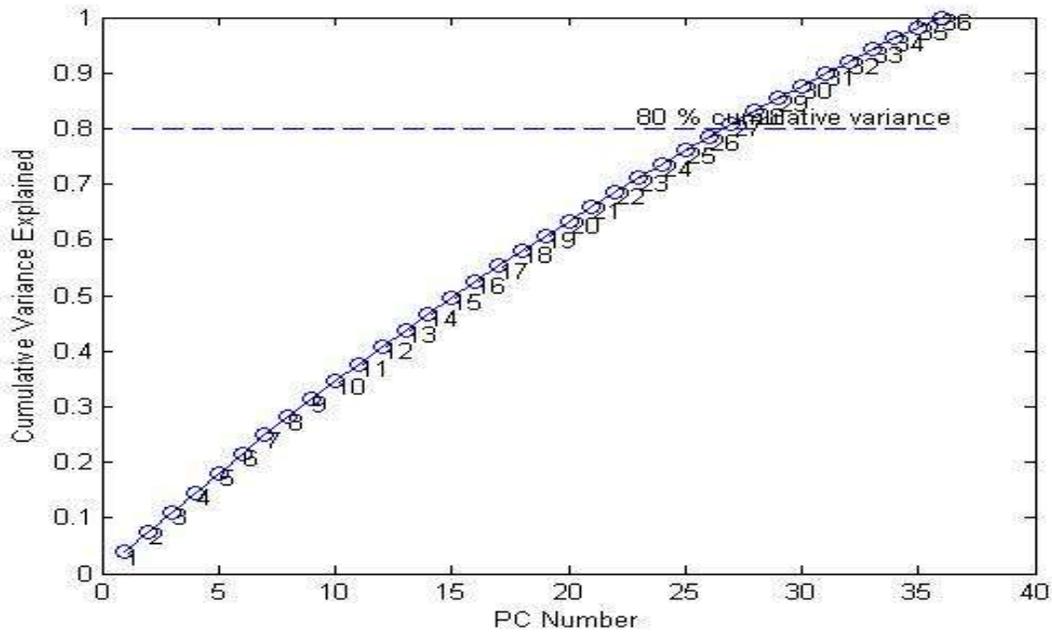| Fault | Description | Additional details |
|---|---|---|
| F1 | Complete leakage in tower C11 bottom | - |
| F 2 | Tower Feed Flow Control Valve, FV11 Fails Closed | - |
| F 3 | Tower Bottom Level Control Valve, LV11 Fails Closed | - |
| F 4 | Reflux Pump – P11A Degradation | - |
| F 5 | Loss of Feed | - |
| F 6 | Reflux Drum Level Control Valve, LV12 Fails Closed | - |
| F 7 | Tower Pressure Control Valve, PV11A Fails Closed | - |
| F 8 | Tower Reboiler - E11 Fouling – variable intensity | severity - 25% at 10 min and 50% after 60 min |
| F 9 | Tower Bottom Level Transmitter, LT11 Drifts | severity – 50% at 10 min and 75% after 60 min |
| F10 | Fault 1 and fault 2 occur simultaneously | - |
| F 11 | Fault 4 and fault 5 occur simultaneously | - |
| F12 | Fault 2 and fault 6 occur in staggered manner | Fault 2 at 10 min, fault 2 and fault 6 occur at 60 min, both deactivated at 120 min |
| F13 | Fault 1 and fault 2 occur in staggered manner | Fault 1 at 10 min, fault 2 and fault 6 occur at 60 min, both deactivated at 120 min |
| F14 | Fault 8 occurs | Deactivated after 120 min |
| F15 | Fault 9 – full intensity | severity – 100% at 10 min, deactivated at 120 min |

### 3.2.2 Results



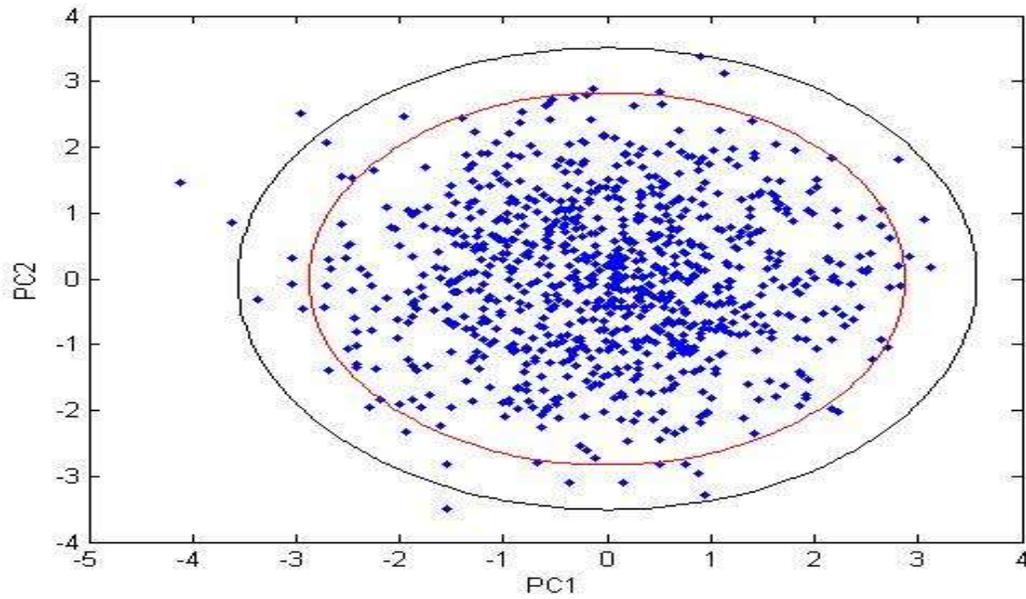Figure 3.21: Cumulative variance explained in the PCA model - DPP



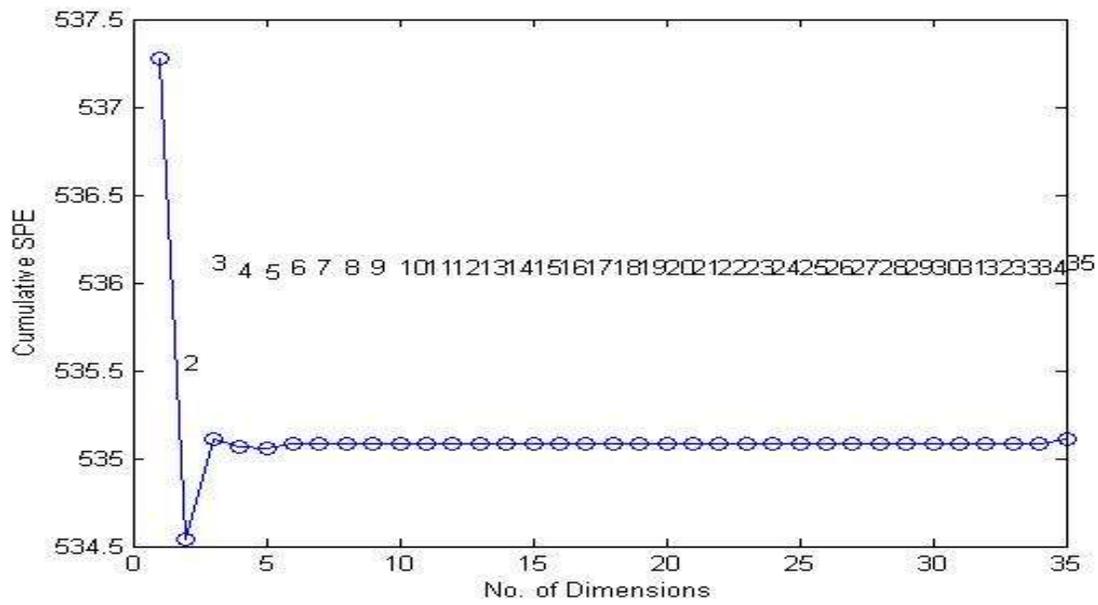Figure 3.22: PCA scores plot for first two PCs - DPP

Figure 3.23: PLS cross validation to choose the number of PCs - TEP
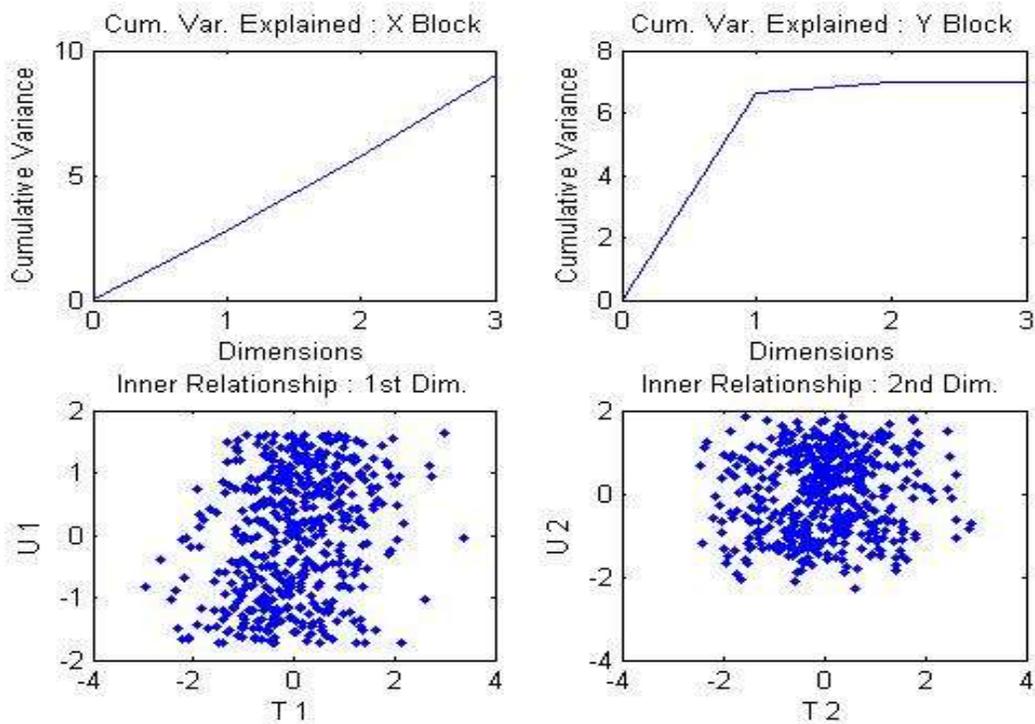


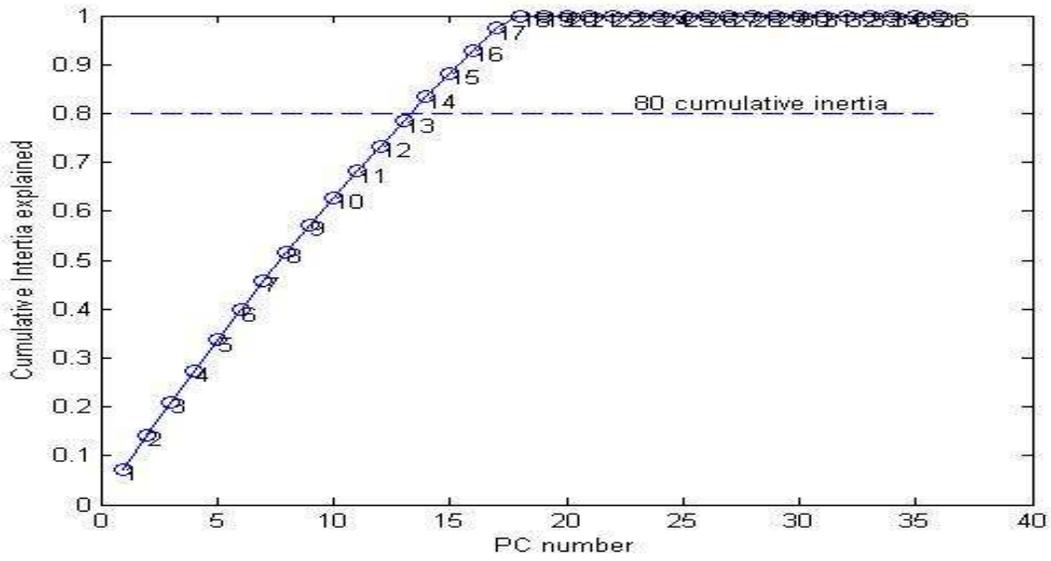Figure 3.24: PLS input-output relationships for 3 PCs - DPP

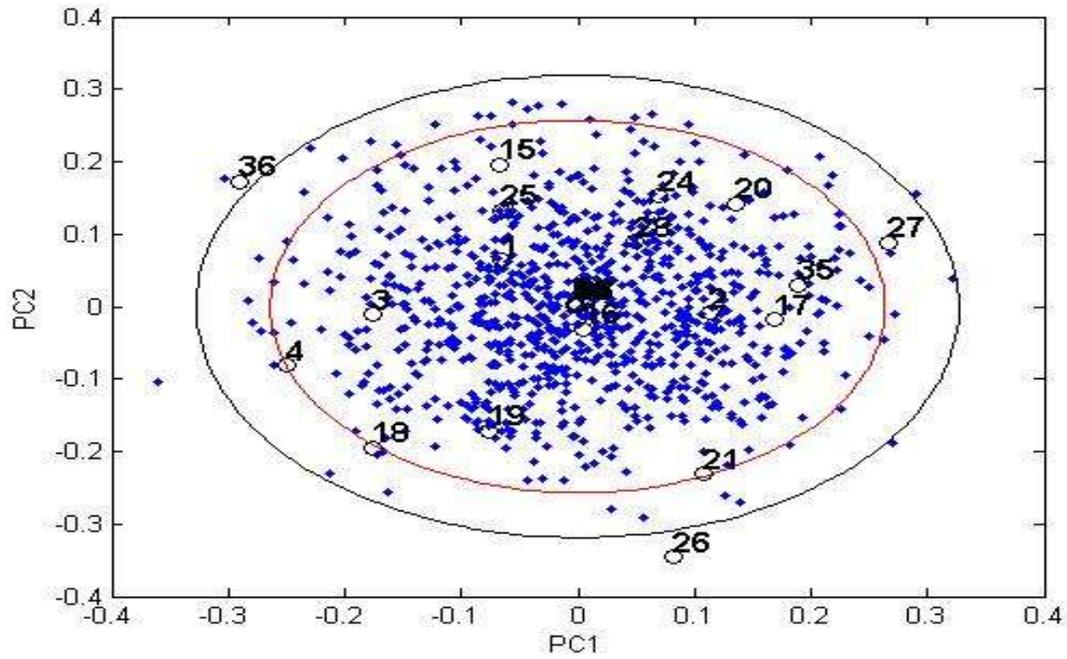Figure 3.25: Cumulative inertia explained in the CA model - DPP



Figure 3.26: CA scores bi- plot for first two PCs - DPP

Table 3.12: Detection rates – Depropanizer Process

| Faults | DR | | | | | FAR | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PCA | | PLS | CA | | PCA | | PLS | CA | |
| | $T^2$ | $Q$ | $T^2$ | $T^2$ | $Q$ | $T^2$ | $Q$ | $T^2$ | $T^2$ | $Q$ |
| F1 | 0.9918 | 0.9953 | 0.9882 | 0.9894 | 0.9800 | 0 | 0.2041 | 0 | 0 | 0 |
| F2 | 0.9977 | 0.9977 | 1 | 0.9988 | 0.9988 | 0 | 0.2653 | 0 | 0 | 0 |
| F3 | 0.9977 | 0.9977 | 1 | 0.9988 | 0.9988 | 0 | 0.1020 | 0 | 0 | 0 |
| F4 | 0.9965 | 0.9977 | 0.9941 | 0.8931 | 0.3314 | 0 | 0.1633 | 0 | 0 | 0 |
| F5 | 0.9965 | 0.9977 | 0.9952 | 0.9871 | 0.9847 | 0 | 0.0408 | 0 | 0 | 0 |
| F6 | 0.9977 | 0.9988 | 0.9917 | 0.9988 | 0.9988 | 0 | 0.2041 | 0 | 0 | 0 |
| F7 | 0.9918 | 0.9930 | 0.9729 | 0.9671 | 0.5781 | 0 | 0.2041 | 0 | 0 | 0 |
| F8 | 0.9883 | 0.9918 | 0.9823 | 0.9871 | 0.9730 | 0 | 0.1224 | 0 | 0 | 0 |
| F9 | 0.9977 | 0.9977 | 1 | 0.9988 | 0.9295 | 0 | 0.2041 | 0 | 0 | 0 |
| F10 | 0.9977 | 0.9988 | 1 | 0.9988 | 0.9988 | 0 | 0.1429 | 0 | 0 | 0 |
| F11 | 0.9977 | 0.9977 | 0.9976 | 0.9882 | 0.9847 | 0 | 0.1633 | 0 | 0 | 0 |
| F12 | 0.9977 | 0.9977 | 0.9788 | 0.9401 | 0.7779 | 0 | 0.1020 | 0 | 0 | 0 |
| F13 | 0.9918 | 0.9941 | 0.9894 | 0.9802 | 0.9812 | 0 | 0.1837 | 0 | 0 | 0 |
| F14 | 0.9883 | 0.9918 | 0.9658 | 0.9530 | 0.7814 | 0 | 0.1633 | 0 | 0 | 0 |
| F15 | 0.9977 | 0.9977 | 0.9835 | 0.9530 | 0.7485 | 0 | 0.2653 | 0 | 0 | 0 |

Table 3.13: Detection delays (in seconds) – Depropanizer Process

| Faults | PCA | PLS | CA |
|--------|-----|-----|-----|
| F1 | 12 | 132 | 108 |
| F2 | 24 | 12 | 12 |
| F3 | 24 | 12 | 12 |
| F4 | 24 | 72 | 180 |
| F5 | 24 | 60 | 132 |
| F6 | 12 | 12 | 12 |
| F7 | 72 | 168 | 240 |
| F8 | 0 | 192 | 132 |
| F9 | 24 | 12 | 12 |
| F10 | 0 | 12 | 12 |
| F11 | 24 | 36 | 120 |
| F12 | 24 | 12 | 12 |
| F13 | 12 | 120 | 120 |
| F14 | 12 | 168 | 156 |
| F15 | 24 | 12 | 12 |

Table 3.14: High contribution variables - Depropanizer Process

| Faults | PCA | CA |
|--------|-----|-----|
| F1 | 10, 19, **28**, 31, 35 | 16, 27, **28** |
| F2 | 9, 13, 23, **28**, 29, 31 | 2, 7, 14, 16, 26, **28** |
| F3 | 1, 14, 20, **21**, 22, 23, 26, **28** | 16, 19, **21**, 24, 27, **28** |
| F4 | **3**, 8, 10, 13, 20, **27**, 30 | **3**, 16, 19, 21, **27**, 28 |
| F5 | 9, 13, 23, **28**, 29, 31 | 2, 7, 14, 16, 26, **28** |
| F6 | 1, 13, 22, 23, 28, 29, 31 | 3, 15, 17, 18, 25 |
| F7 | 6, 8, 13, 23, **28**, 31, 35 | 16, 26, 27, **28** |
| F8 | 13, 20, **28** | 16, 27, **28** |
| F9 | 14, 20, 22, 26, **28**, 31 | 16, 19, 27, **28** |
| F10 | 7, 10, 14, **19**, 31 | 1, **19**, 24, 27, 28 |
| F11 | **2**, 13, 23, **28**, 29, 31 | **2**, 7, 14, 16, 26, **28** |
| F12 | 2, **3**, **28**, 29, 31, 32 | **3**, 16, 26, **28** |
| F13 | 10, 14, **19**, 31 | **19**, 24, 27, 28 |
| F14 | 13, **28** | 3, 16, 27, **28** |
| F15 | 1, 14, 20, 22, 26, **28**, 31 | 3, 16, 19, 27, **28** |

In the case of the Depropanizer process, all results were found to be quite consistent with all 15 faults showing a detection rate greater than 0.9. PCA was still found to exhibit false alarms but the only fell within the range of 0.10 to 0.26. In the case of diagnosis with contribution plots, CA was found to be on par or better than PCA in 14 of the 15 cases, the exception being fault 14 where PCA showed only two main contribution variables as compared to four shown by CA. In the 14 cases where CA indicated better diagnosis faults 2, 5, 8, 10, 11 and 13 were found to show the same number of contribution variables while faults 1, 3, 4, 6, 7, 9, 12 and 15 were found to show less number of main contribution variables in the case of CA diagnosis.

## 3.4    Discussion

From the results obtained for the three systems, it is clear that PCA is the most powerful of the three tools when it comes to detection, but it also has the biggest disadvantage of false alarm rates caused by its inability to understand non linearity and serial correlation dynamics. CA is noted to overcome these problems but its detection delays are found to be higher and hence its detection rates are found to be comparatively lower except in a few cases. The main problem was found to be with the Q statistic which was found to not be effective; this may be because the residual space is affected by the cross tabulation dual analysis which distorts the analysis. Therefore, there is need to find an improved or modified statistics which can help monitor the residual space to a better extent. PLS which performs its analysis between two sets of variables was also found to be quite effective but could not gain the edge over CA since it was still a linear technique. As far as diagnosis was concerned CA was found to be a more concrete tool in diagnosis in all three systems and could be relied upon under any circumstances.

# 4. FAULT ISOLATION AND IDENTIFICATION METHODOLOGY

The main aim of this chapter is to highlight the importance of Linear Discriminant Analysis (LDA) in the field of diagnosis. The chapter will explain the basis of LDA along with a literature survey on its application in fault detection and diagnosis. This will be followed by a comparison of diagnosis performance with CA and the formulation of the integrated CA-WPSLDA technique for fault isolation and identification. The formulation of this technique will also include an explanation on the superior discriminative abilities of CA as compared to PCA.

In the field of fault diagnosis, fault isolation involves isolating the specific fault that occurred. It also includes determining the kind of fault, the location of the fault, and the time of detection while fault identification deals with determining the size and time-variant behaviour of a fault. In this regard, the newly integrated algorithm will use all the information available from historical datasets to create a model which will try to isolate a new fault during the monitoring phase by identifying whether it is related to ones that have previously occurred and would then identify the intensity of the fault with respect to the ones in the model.

## 4.1 Linear Discriminant Analysis

### 4.1.1 LDA - Introduction

LDA or Fisher's Linear Discriminant (FLD) is an optimal dimensionality reduction technique in terms of maximizing the separability of these classes. It determines a set of projection vectors that maximize the inter-class scatter while minimize the intra-class scatter.

In fault diagnosis, data collected from the plant during specific faults is categorized into classes, where each class contains data representing a particular fault. Let $X \in \mathcal{R}^{nn \times mm}$ be a set of mm-dimensional samples containing all the data related to the various faults (classes) where the total number of classes is $M$. Then, $x \in \mathcal{R}^{mm}$ and the matrix $X_i$ is the subset which contains $nn_i$ rows corresponding to the samples from class i .

Then,

$$\bar{x} = \frac{1}{nn} \sum x \tag{4.1}$$

$$\bar{x}_\iota = \frac{1}{nn_i} \sum_{x \in X_i} x \tag{4.2}$$

Where $\bar{x}$ is the overall mean for all samples in $X$ and $\bar{x}_\iota$ is the mm–dimensional mean for the samples belonging to each class $i$. The within-class scatter matrix $S_w$ is calculated as a measure of the spread within a class of data.

$$S_i = \sum_{x \in X_i} (x - \bar{x}_\iota)(x - \bar{x}_\iota)' \tag{4.3}$$

$$S_w = \sum_{i=1}^{M} S_i \tag{4.4}$$

The inter-class matrix, which is a measure of the overall spread between the class is given by,

$$S_b = \sum_{i=1}^{M} nn_i \, (\bar{x}_\iota - \bar{x})(x - \bar{x}_\iota)' \tag{4.5}$$

$$S_t = S_w + S_b \tag{4.6}$$

$$S_t = \sum_{j=1}^{nn} (x_j - \bar{x})(x_j - \bar{x})' \tag{4.7}$$

Where, $S_t$ is called the total scatter matrix. The optimal Fisher direction is found by maximizing the following Fisher criterion $J(\varphi)$:

$$J(\varphi) = \frac{\varphi' S_b \varphi}{\varphi' S_w \varphi} \tag{4.8}$$

The maximizer $\varphi$ is the Fisher optimal discriminant direction that maximizes the ratio of the inter-class scatter to the intra-class scatter. The maximizer contains the discriminant vectors equal to the generalized eigenvectors of the eigenvalue problem.

$$S_b \varphi = \Lambda_{LDA} S_W \varphi \tag{4.9}$$

If, $S_W$ is non singular, the eigenvector could be further modified to give,

$$S_w{}^{-1} S_b \varphi = \Lambda_{LDA} \varphi \tag{4.10}$$

where the eigenvalues $\Lambda_{LDA}$ indicates the degree of overall separability among the classes. The score matrix $T_{LDA}$ is obtained by projecting the observations X onto the Fisher directions $\varphi$.

$$T_{LDA} = X \, \varphi \tag{4.11}$$

### 4.1.2 Literature Survey

The first attempt to use LDA for fault diagnosis was done by Raich and Çinar (1994) who developed a methodology to integrate PCA and LDA in order to determine out-of-control status of a continuous process and to diagnose the source causes for abnormal behaviour. Chiang *et al*. (2000) later applied LDA to most of the faults in the Tennessee Eastman process simulation to obtain one lower dimensional model which could be used for diagnosis as well as detection by including another class containing data from the normal operating condition. He *et al*. (2005) developed a fault diagnosis method based on fault direction using PCA and LDA, which they successfully applied to the quadruple tank system for sensor and leakage faults as well as to an industrial film polyester manufacturing process. Both Jiang *et al.* (2008) and He *et al.* (2009)

later used partial F-values and the Cumulative Percentage Variance (CPV) values along with FDA for the identification of key variables responsible for abnormalities and the development of a Variable weighted FDA (VW-FDA) technique for better discrimination.

## 4.2 The integrated CA-WPSLDA methodology

The integrated CA-WPSLDA methodology is a technique developed for the isolation and identification of faults detected during the monitoring stages of a system. It attempts to use the FDA space as a monitoring space instead of just diagnosis, and tries to provide a simple graphical plot which may be used by operators to understand the nature of a fault that they encounter in a plant

### 4.2.1 Motivation

The motivation for the WPSLDA algorithm was based on the fault diagnosis methodology by He *et al.* (2005). In this paper, the authors first developed an algorithm based on PCA and LDA which is used to detect and isolate fault related data in historical data sets for monitoring purposes. A PCA model of the normal operating data was first used to detect other faults in the historical dataset. These datasets would be combined and later subjected to PCA where certain clusters would be visible and K-means clustering could be used to roughly isolate normal and abnormal clusters. The final dataset after removing samples based on K-means clustering is then subjected to LDA for better visualization in much lower dimensions. Then, pairwise LDA was applied to the normal operating dataset and each class of fault alone to obtain a LDA vector which is treated as a contribution plot to understand the nature of the fault involved. This work provides the basis for a similar yet modified algorithm which could also be used for monitoring

as well as isolation and identification purposes. The several modifications and the reasons for doing so will be explained in the subsequent sections.

### 4.2.2 A combined CA plus LDA model

In the work by He *et al.* (2005), the authors had used the PCA for two reasons, primarily for fault detection. We wish to replace this method with CA as it had been proved earlier that CA is a more robust detection tool. This was verified in chapter 3 during the application of CA to the quadruple tank system where all the faults were detected at much lower false alarm rates and almost acceptable detection rates. It was also noticed in Table 3.2 that fault 3 and 8 which were simulated at slightly different operating conditions were detected properly by CA, while high false alarm rates were recorded in PCA due to its inability to account for dynamics of the system. This would be very useful, especially in historical datasets that are recorded for longer lengths of time and could therefore have been recorded under different operating conditions.

The other use of PCA in the original algorithm was for pre-analysis with K-means clustering to roughly identify the clusters. Later on, PCA did not play a role in pairwise LDA calculations as the direction vectors obtained are treated as contribution plots and require all the original variables from the system to understand the cause of the abnormality. In our case, we wish to use the tool for isolation purposes, but, not by means of any contribution plots; hence the need to retain the original variables for the final calculations. Therefore, CA was used to develop the final combined model for pre-analysis but its row scores would later be used for LDA and not the original dataset. This was done for two reasons, the first being that applying a technique like PCA or CA will lead to dimensional reduction and will not lead to much loss of information; this was proved by Yang *et al.* (2003) in the case of PCA. Since CA can store much more improved

information as compared to PCA, it is a better choice to be used. The other reason for using CA is that it has better discriminative properties than PCA (both PCA and K-means clustering tend to fail as the extent of non-linearity is found to increase). This property of CA is called the process of self-aggregation, where CA can provide better discriminative clusters and is attributed to the fact that generalized SVD is performed in CA. The process of self-aggregation was first explained by Ding *et al.* (2002) who explained that self-aggregation is governed by connectivity and occurs in a space obtained by a nonlinear scaling of PCA called Scaled Principal Component Analysis (SPCA). They had stated that nonlinear scaling in PCA can be performed by obtaining scaling factors in the form of a diagonal matrix, where each value along the diagonal is the sum of the corresponding row of the covariance/correlation matrix represented by $WW$.

Let,

$$WW = ww_{ij} \tag{4.12}$$

Then, et the scaling factor be

$$D = d_i \tag{4.13}$$

and,

$$d_i = \sum_j ww_{ij} \tag{4.14}$$

Thus, the new scaled matrix is,

$$\widehat{WW} = D^{-\frac{1}{2}} WW D^{-\frac{1}{2}} \tag{4.15}$$

which leads to,

$$WW = D \sum_k (q_k \Lambda_k q_k') D \tag{4.16}$$

where,

$$q_k = D^{-\frac{1}{2}} z_k \tag{4.17}$$

And, the final eigenvalue problem is defined as,

$$\widehat{WW}z = \Lambda z \qquad (4.18)$$

or,

$$WWq = \Lambda Dq \qquad (4.19)$$

In the above formulation, Ding *et al.* (2002) explained that when there are $K$ clusters and there are no overlaps between them in the regular Euclidian space, then the scaled K principal components $(q_1, q_2, \dots, q_k) = Q_K$ get the same maximum eigenvalue equal to 1. In the SPCA space spanned by $Q_K$, all the objects within same cluster self-aggregate into a single point. However, when overlaps between different clusters are present, samples within same cluster tend closer to each other in Scaled SPCA space than in Euclidean space. Khare *et al.* (2008) compared the SPCA algorithm to normal PCA and FDA where he stated that SPCA could be comparable to FDA as it is an unsupervised tool which also has the ability to greatly reduce intra-clustering distances enhancing segregation. Now, taking the case of SPCA and comparing it to CA, non-linear scaling is applied by the use of generalized SVD. Generalized SVD is usually applied when there is a need to impose constraints on the rows and columns of a matrix by using two positive definite matrices. In the formulation of CA in chapter 2, the term $(CM - rc^T)$ can be subjected to generalized SVD where,

$$(CM - rc^T) = AaD_\mu Bb \qquad (4.20)$$

subject to the constraints,

$$AaD_r^{-1}Aa = I_r \qquad (4.21)$$

and

$$BbD_c^{-1}Bb = I_c \qquad (4.22)$$

The above three expressions are the same as the SVD equation given in 2.32 where,

$$SM = LD_\mu M^T$$

Equations 2.34 and 2.35 show that,

$$Aa = D_r^{\frac{1}{2}}L$$

$$Bb = D_c^{\frac{1}{2}}M$$

One can notice that 2.34 and 2.35 are similar to equation 4.17. Thus, one may conclude that CA is slightly different from non-linear SPCA applied to the rows and columns of the dataset. This was in agreement with the statements provided by Detroja *et al.* (2006). From the above points, it can be concluded that a CA plus LDA formulation is preferred over the methodology used by He *et al.* (2005).

An example over the discriminative property of CA was applied by following the first two steps in the algorithm alone where data from the TEP process was taken from the website http://brahms.scs.uiuc.edu (link is no longer functional) as in chapter 3 but for a total of 52 variables for the normal operating condition, fault 4 and fault 11. Both fault 4 and 11 are associated with the same fault variables. But fault 4 is related to a step change in the reactor cooling water temperature while fault 11 is more related to the reactor cooling water inlet temperature and is subjected to random variation as compared to the step change in fault 4. The faults were first monitored using both PCA and CA separately and then subjected to a combined model using the respective algorithm in each case.

Table 4.1: Detection rates and false alarm rates – TEP with fault 4 and fault 11

| Datasets | Symbol | Detection Rates | |
|---|---|---|---|
| | | PCA | CA |
| Normal Condition | Green circle | - | - |
| IDV(4) | Red Circle | 1 | 1 |
| IDV(11) | Blue Circle | 0.2991 | 0.5663 |

The number of PCs obtained for the PCA combined model was found to be 25 while for CA the number of PCs were found to be 2 for a Cumulative percentage in variance and inertia of 80%. The scores of PCA and the row scores of CA were projected onto the first two dimensions in fig 4.2 and 4.4.



Figure 4.1: Cumulative variance shown in the combined PCA model for TEP example

Figure 4.2: Scores plot for first two components of the combined PCA model – TEP



Figure 4.3: Cumulative inertial change shown in combined CA model for TEP example

Figure 4.4: Row scores plot for first two components of combined CA model – TEP

Thus, it can be clearly seen from Figures 4.1, 4.2, 4.3, and 4.4 that CA can distinctly present the clusters for a normal operating condition and two other faults even when both the faults share a certain amount of similarity to one another. It was also proved that CA can provide this visualization at much lower representation.

### 4.2.3 A weighted LDA algorithm

Following the development of the combined CA model, the scores are subjected to LDA at two levels. The first application of LDA will be to the complete set of row scores corresponding to the selected number of components in CA. The main aim here is the visualization of the transformed dataset in the Fisher space. Visualization is usually preferred corresponding to the two largest eigenvalues (2-D space). The second analysis involves the use of pairwise LDA to a combination of the normal operating condition and each fault. Since there are only two classes used in pairwise LDA there will be only one significant non-zero eigenvector which will be used for projection purposes and to later on develop the monitoring scheme based on control charts.

The need for a weighted LDA algorithm arouse from the fact that the presence of overlapping cluster despite applying CA would tend to disrupt the algorithm. Variable weighted techniques had been applied earlier using partial F-values along with CPV in LDA (He *et al.,* 2009 ) but the procedure was found to be quite tedious and complex and did not provide any weights to the class of normal operating data. Therefore there was a need to identify a weighting technique which was simple and treated all the classes of data equally while providing improved discriminative visualization. The solution to this problem was seen in the form of a weighted pairwise scatter linear discriminant analysis (WPSLDA) algorithm which was suggested by Li *et al.* (2000). According to these authors, an implicit assumption in LDA is that each class may be equally confused with other classes. This can be explained by deriving the following equations. We know that the within-class scatter matrix is given by,

$$S_w = \sum_{i=1}^{M} S_i \tag{4.23}$$

This can be rewritten using equation 4.3:

$$S_w = \sum_{i=1}^{M} (x - \bar{x_i})(x - \bar{x_i})' \tag{4.24}$$

Then,

$$S_w = \sum_{i=1}^{M} nn_i \Sigma_i \tag{4.25}$$

where, $\Sigma_i$ is the covariance matrix for each class of data. The covariance matrix for the different classes of data as well as the whole dataset $X$ can be given by,

$$\Sigma_i = \frac{1}{nn_i} \sum_{x \in X_i} (x - \bar{x_i})(x - \bar{x_i})' \tag{4.26}$$

This can again be written as,

$$\Sigma_i \ = \ \frac{1}{nn_i}\Sigma_{x \in X_i}(xx' - x_i x_i')$$ (4.27)

Let, the covariance matrix for the whole dataset be given by,

$$\Sigma \ = \ \frac{1}{nn}\Sigma_{j=1}^{nn}(x_j - \bar{x})(x_j - \bar{x})'$$ (4.28)

This can again be written as,

$$\Sigma \ = \ \frac{1}{nn}\Sigma_{j=1}^{nn}(x_j x_j' - \bar{x}\bar{x}')$$ (4.29)

We know that, the total scatter matrix, $S_t$ which is the sum of the within-class scatter matrix and the between class scatter matrix can be written as:

$$S_t = \Sigma_{j=1}^{nn}(x_j - \bar{x})(x_j - \bar{x})'$$ (4.30)

Then

$$S_t = \Sigma_{j=1}^{nn}(x_j - \bar{x_\iota} + \bar{x_\iota} - \bar{x})(x_j - \bar{x_\iota} + \bar{x_\iota} - \bar{x})'$$ (4.31)

where, $\bar{x_\iota}$ is the class mean corresponding to the class of data that each sample $x_j$ belongs to, equation 4.31 can now be written as:

$$S_t = \Sigma_{j=1}^{nn}((x_j - \bar{x_\iota})(x_j - \bar{x_\iota})') + ((\bar{x_\iota} - \bar{x})(\bar{x_\iota} - \bar{x})')$$ (4.32)

The transformation from equation 4.31 to 4.32 is similar to the ones that take place between equations 4.26 and 4.27 as well as between 4.28 and 4.29. Equation 4.32 finally becomes:

$$S_t = \Sigma_{i=1}^{M} nn_i \Sigma_i \ + \Sigma_{i=1}^{M} nn_i (\bar{x_\iota} - \bar{x})(x - \bar{x_\iota})'$$ (4.33)

Then from equations 4.5, 4.25, and 4.33, we again arrive at the fact that,

$$S_t = S_w + S_b \qquad (4.34)$$

The focus of the previous formulations is inter-class scatter matrix $S_b$. According to Li *et al.* (2000), the inter-class scatter matrix in its regular form neglects any discriminatory information if the distance between certain classes are much closer to each other as compared to others. This was demonstrated with the following case where we have four clusters spanning over a two dimensional space each having the same number of samples and equal variance as shown in Figure 4.5 where the mean of each class is given as $\mu_1 = (1, \delta)$, $\mu_2 = (-1, \delta)$, $\mu_3 = (-1, -\delta)$, and $\mu_4 = (1, -\delta)$.



Figure 4.5: WPSLDA case study

The inter-class scatter matrix is given by,

$$\frac{1}{4}S_b = \begin{pmatrix} 1 & 0 \\ 0 & \delta^2 \end{pmatrix} \qquad (4.35)$$

Now as $\delta \to 0$, the matrix is of the form $\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$ where it is only possible to discriminate between the of class pairs of (1,4) and (2,3) whose covariance dominates the model. Although it is true that both these pairs are important in the model, it still does prove the fact that the

between class scatter matrix does not accurately represent the discriminatory information in the model. Therefore, the between class scatter matrix was redefined to be the sum of pairwise scatter matrices. This new version of the within-class scatter matrix $S_{bw}$ is given by:

$$S_{bw} = \frac{1}{2nn} \sum_{i,z=1}^{M} mw_{iz} nn_i \, nn_z (\bar{x}_i - \bar{x}_z)(\bar{x}_i - \bar{x}_z)' \tag{4.36}$$

This new form of between class scatter matrix is developed such that a certain set of weights in a matrix given by $MW$ where, $mw_{iz}$ is the value provided for a pair of classes referenced by 'i' and 'z' to improve the information via scatter so that the pair could be treated with a certain required amount of importance in the LDA model. The weightage value for a certain class is calculated as follows based on their mean values.

$$mw_{iz} = \frac{1}{(\bar{x}_i - \bar{x}_z)'(\bar{x}_i - \bar{x}_z)} \tag{4.37}$$

Equation 4.36 can be simplified to equation 4.5 when the weightage value is assumed to be 1 in all cases. This will mean that each pairwise scatter will contribute equally to the between class scatter matrix. Then, the equation is found to be:

$$S_{bw} = \frac{1}{2nn} \sum_{i,z=1}^{M} nn_i \, nn_z (\bar{x}_i - \bar{x}_z)(\bar{x}_i - \bar{x}_z)' \tag{4.38}$$

$$S_{bw} = \frac{1}{2nn} \sum_{i,z=1}^{M} nn_i \, nn_z (\bar{x}_i - \bar{x} + \bar{x} - \bar{x}_z)(\bar{x}_i - \bar{x} + \bar{x} - \bar{x}_z)' \tag{4.39}$$

$$S_{bw} = \sum_{i=1}^{M} nn_i \, (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})' \tag{4.40}$$

Equation 4.39 is the same as that of 4.5 for regular $S_b$ and thus one can say that $S_b$ is a special case of $S_{bw}$ when the weights are uniform.

## 4.2.4 Fault intensity calculations

After applying the WPSLDA algorithm for better visualization, pairwise FDA is performed between each of the fault classes to the normal data. Since there are only two classes involved in these pairwise calculations, the number of significant discriminant vectors is 1.
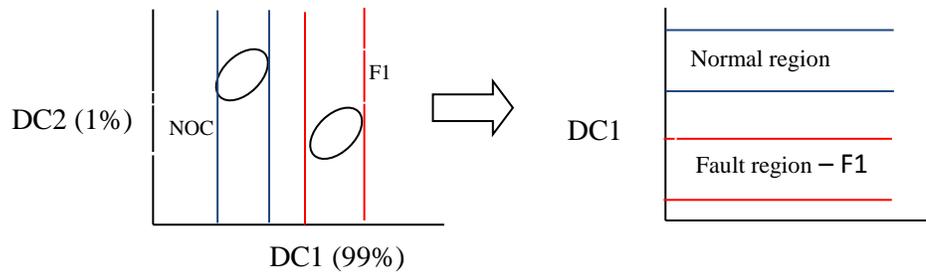


Figure 4.6: Control chart like monitoring scheme from pairwise LDA-1

Since the number of significant discriminant vectors is 1, the two-dimensional plot as shown in Figure 4.5 for pairwise FDA of two classes can be converted to another having just one dimension, i.e. the most significant discriminant direction. The bounds for the two regions can be chosen by selecting the maximum and minimum value along the same to provide bounds for the two regions, if the monitored data (after undergoing a series of transformations) is found to exceed the bound of the normal region and approach the fault region. This can be indicated by certain bar plots which conduct fault intensity calculations. The main aim of calculating this intensity value is to understand how strongly the samples are related to a certain fault in the simplest way possible as visualization of the sample in a multi-class LDA visualization may not provide a clear picture of the outcome. The fault intensity values are expressed in percentage between 0 -100 %. Calculations are carried out as follows according to the following set of rules:
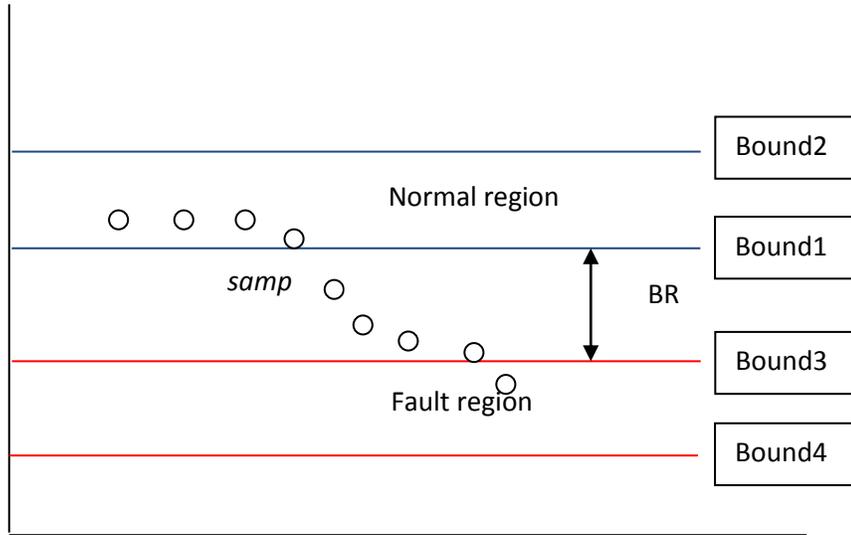
Figure 4.7: Control chart like monitoring scheme from pairwise LDA-2

When a sample *samp* is being monitored, it has to move from the normal region to the fault region, this transitional region is called the buffer region and its distance is termed *BR*. The two limits that are necessary for the calculation of intensity would first include the limit that the sample has to cross to leave the normal region and the limit it has to cross to enter the fault region. Each of these limits would be referred to as *Bound1* and *Bound2*. Thus, the intensity of the sample would be calculated as:

$$Intensity = \frac{Bound1 - samp}{BR} \tag{4.41}$$

The bound values would be interchanged if the fault region were to lie above the normal region. In this case, the equation would change to,

$$Intensity = \frac{samp - Bound2}{BR} \tag{4.42}$$

Other rules which are followed in these calculations include:

1) The intensity values remain at 0 as long as *samp* is between *Bound1* and *Bound2*.

2) The intensity value is directly assigned as 1 if it is between *Bound3* and *Bound4*.

3) If the samples do enter the fault region but are found to move beyond this region too, then their intensity values are reduced by a factor of $\frac{Boun4-samp}{BR}$.

4) If the samples are found to cross the bounds of normal operation but move in the direction opposite to that of the fault region, then they are assigned a value of 0.1 to indicate that a fault has occurred, but is not related to the fault in the chart.

In industrial settings, it is not advisable to arrive at a conclusion based on a single sample, therefore, one would take the average of *'num'* number of samples before the current one to provide the final intensity value on a bar plot. The value of *'num'* is chosen based on the convenience of the user. A sample plot on how the bar plot presentation would appear is given in Figure 4.8.

Figure 4.8: Control chart like monitoring scheme with fault intensity bar plots

In, this plot, one can clearly see that the monitored samples are found to have a strong affinity to fault 1 shown in red. This can also be noticed in the first control chart at the top of the Figure where the samples have crossed over from the lower zone, which is the normal zone to the upper zone.

Thus, with these intensity calculations, a complete explanation of the CA-WPSLDA methodology has been concluded and a complete summary of the procedure is provided below:

1) A CA model of the normal operating condition is first developed; it is then used on historical datasets to detect faults using $T^2$ and $Q$ statistics.

2) The data related to the faults detected are then combined with the normal operating data used to create the initial CA model at a very high cumulative inertia (say 95%).

3) The combined dataset is then subjected to CA for two main purposes; firstly dimension reduction and secondly preliminary discrimination.

90

4) The row scores of this new combined model are then subjected to Weighted Pairwise Scatter Linear Discriminant Analysis (WPSLDA) to push any clusters that may have been too close or may have overlapped with one another. If the clusters are already further apart, then there will be no need to apply WPSLDA to the combined model.

5) WPSLDA is then applied in a pairwise fashion for the row scores each of the fault related datasets along with the normal operating data. The LDA vectors obtained for each pairwise calculation represent the fault directions for each fault.

6) These pairwise LDA vectors are used to develop a control chart where the boundaries are marked for the operating condition as well as the fault.

7) Intensity calculations are performed based on the position of the monitoring sample in the chart to predict its chances of being part of a certain fault. This intensity value is shown in the form of a bar plot for each sample.



Figure 4.9: CA-WPSLDA methodology

**4.3 Comparison of integrated methodology to LDA**

In order to compare the integrated methodology, we compare initially the results of the combined CA model developed in section 4.2.1 to LDA. The samples selected by CA monitoring will then be subjected to LDA under supervised conditions.
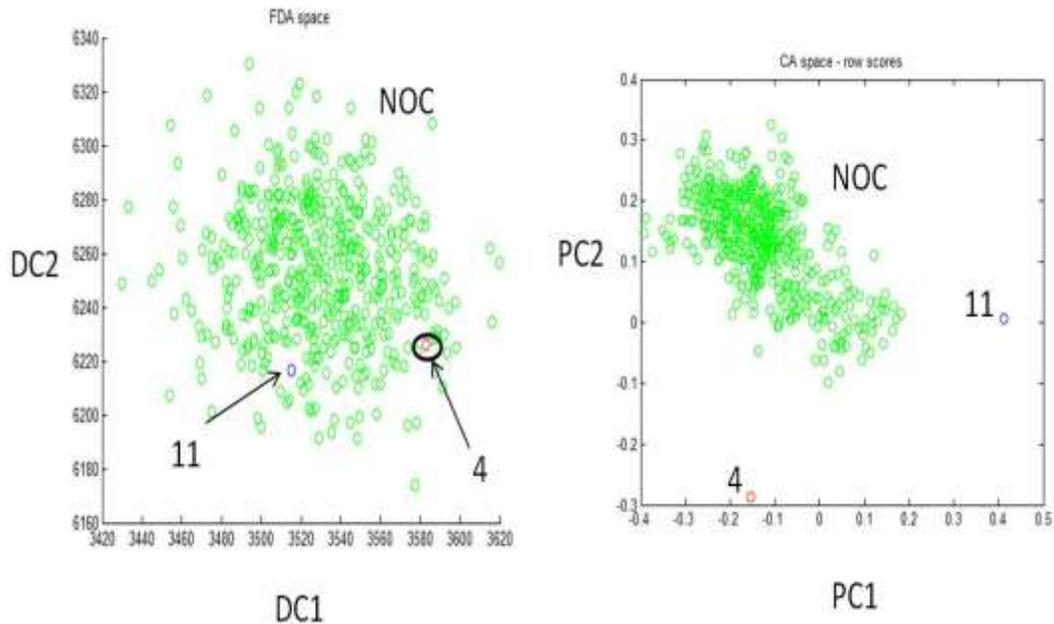


Figure 4.10: Comparison between CA and LDA

It is very clear from Figure 4.10 that in this case, CA is much more superior to LDA. There may be certain cases where the number of CA dimensions would be greater than 2 for the combined model and in this case it is better to apply WPS-LDA to these scores to reduce the dimensions further and improve separation if possible. Thus the integrated CA-WPSLDA methodology is found to be far more efficient as compared to PCA (from section 4.2.2) and LDA in terms of discrimination due to the application of the WPSLDA methodology over the CA space.

**4.4 Application to simulated case studies**

The integrated methodology has been applied to simulated case studies of the Quadruple tank system and the Depropanizer process. The faults are the same as the ones described in Table 3.2 and Table 3.5. The intensity values will be shown in the form of curves for convenience in both the cases.

**4.4.1 Quadruple tank system**

The five classes involved in the development of the model will include the normal operating condition, faults 1, 2, 3, and 4. The faults 5, 6, 7, and 8 are then tested using the algorithm to see if the nature of the faults can be predicted.

Table 4.2: Quadruple tank system – model faults and symbols

| Datasets | Symbol |
|---|---|
| Normal Condition | Green circle |
| Fault 1 | Red Circle |
| Fault 2 | Blue Circle |
| Fault 3 | Black circle |
| Fault 4 | Cyan circle |

**4.4.2 Depropanizer Process**

In this system, the first 9 faults are used to develop the integrated model while faults 10, 11, 12, 13, 14, and 15 are monitored by the CA-WPSLDA methodology and the results are obtained. The description of the faults can be obtained from Table 3.5.

Table 4.3: DPP – model faults and symbols

| Datasets | Symbol |
|---|---|
| Normal Condition | Green circle |
| Fault 1 | Red Circle |
| Fault 2 | Blue Circle |
| Fault 3 | Black circle |
| Fault 4 | Cyan circle |
| Fault 5 | Red Cross |
| Fault 6 | Yellow Circle |
| Fault 7 | Magenta Circle |
| Fault 8 | Black cross |
| Fault 9 | Magenta Cross |

## 4.5 Results and Discussion

### 4.5.1 Quadruple tank system

The final CA and WPSLDA models for the quadruple tank system are developed and the results are as shown in the Figures below.



Figure 4.11: Number of PCs for combined CA model – Quadruple tank system



Figure 4.12: First 2 PCs of final combined CA model – Quadruple tank system

Figure 4.13: Final WPSLDA model – Quadruple tank system

In this case one will find that all the four clusters do separate very well. The number of PCs for the combined CA model is chosen at a cumulative inertia level of 95%. This is because the data contained in these classes could be spaced far apart, and information might be lost by treating some of the samples as noise. We use the WPSLDA model to reduce the number of dimensions and fit all our information into just two dimensions. The four control charts are then developed and then applied for monitoring purposes.



Figure 4.14: CA-WPSLDA methodology – monitoring – fault 5

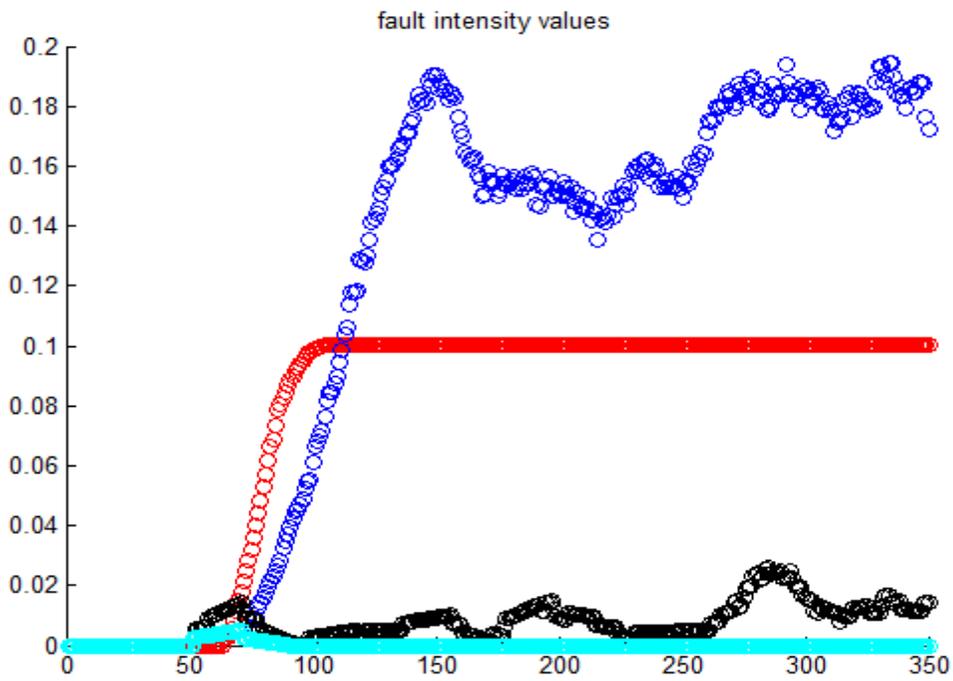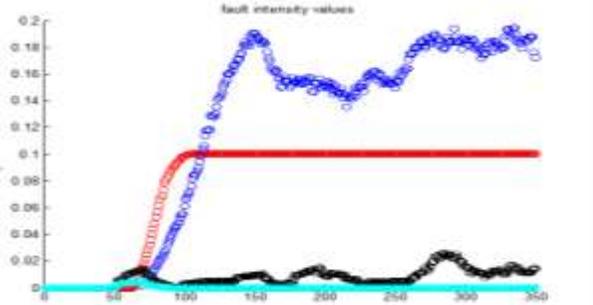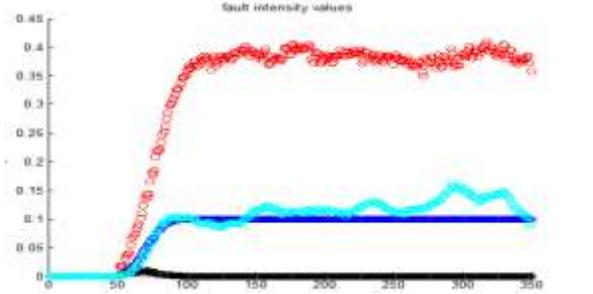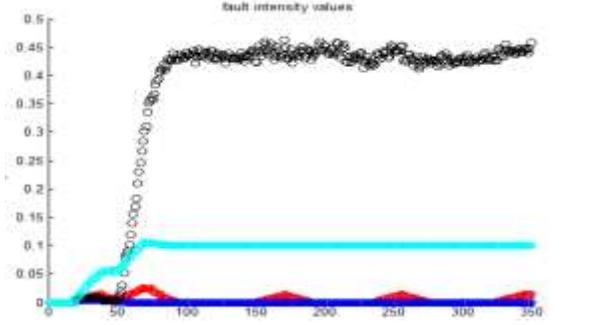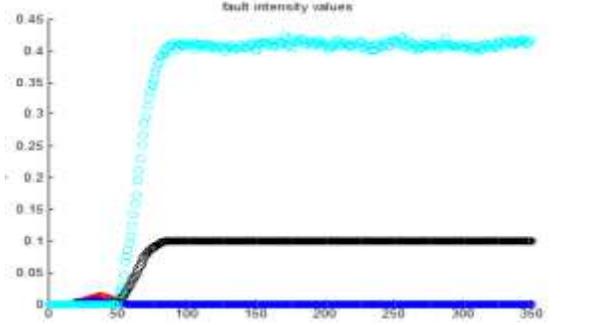Figure 4.15: CA-WPSLDA methodology – control charts – fault 5



Figure 4.16: CA-WPSLDA methodology – intensity values – fault 5 (x-axis: sample number; y-axis: fault intensity)

Table 4.4: Quadruple tank system – CA-WPSLDA methodology results

| Fault | Results – fault intensity values (x-axis: sample number; y-axis: fault intensity) | Description of results |
|---|---|---|
| 5 |  | Clear fault affinity is shown around the 150[th] sample towards fault 2 at a value between 14 – 20 %. Fault 1 could be related or is just displaying the presence of a fault. |
| 6 |  | Highest fault intensity is associated with fault 1 at a value of 40% at the 90[th] sample while other faults show only a maximum intensity of 15 % |
| 7 |  | Highest fault affinity is related to fault 3 at a value of 45 % starting from the 90[th] sample while others are 10% or less. |
| 8 |  | Highest fault affinity is related to fault 3 at a value of 40 % starting from the 90[th] sample while others are 10% or less. |

In Figure 4.14, the fault regions seem to be represented by straight lines as their fault regions are very narrow as compared to the normal region. From the Figure, it is also clear that only fault 2 (represented by blue circles) has some approach towards its region while the other drift away from the region of normal operation, but in the direction opposite to that of the fault region; hence their intensity calculations would be negligible. The intensity values shown in Figure 4.15 support the control charts where only the intensity values of fault 2 show a variation between 15 and 20%, while fault1 shows a variation of 10% which could be an approach towards the fault or just an indicator that the fault has left the normal region and may be proceeding in the opposite direction. Intensity values of faults 3 and 4 also convey the same information but at much lower intensities. Thus the only conclusion for fault 5 is that out of the two contributing faults of 1 and 2, only fault 2 is identified by the CA-WPSLDA method as being associated with fault 5. Fault 6 clearly shows that it is associated with fault 1 which is true, as both fault 1 and 6 are related to a leakage in tank 1 and the leakage co-efficient in fault 6 is 40% of the value in fault 1, which is also the fault intensity value shown in Table 4.2. Fault 7 which is related to a positive bias in height $h_1$ with bias value of 0.4 is shown to be clearly related to fault 3 which has a negative bias in the sensor related to height $h_1$ with the same bias value of 0.4. Fault 8 related to a positive bias in height $h_2$ was clearly found to be related to fault 4, which was also related to a bias in $h_2$ but in the negative direction. The absolute value of bias taken in this case was also 0.4.
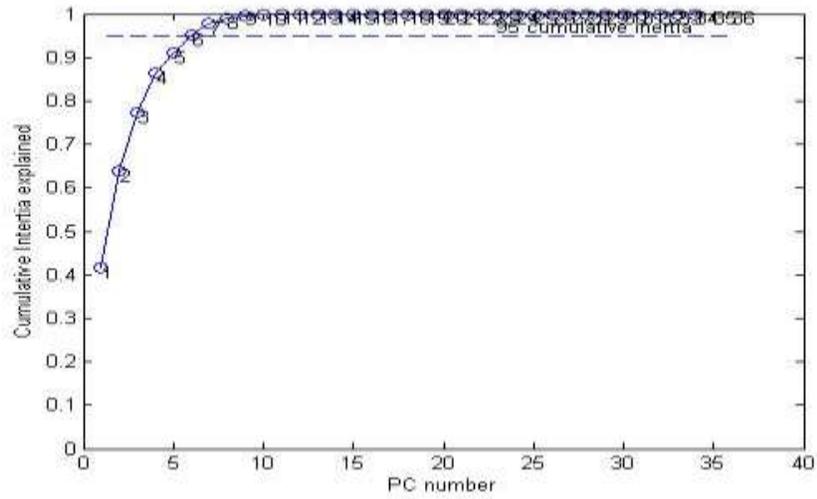
**4.5.2 Depropanizer Process**



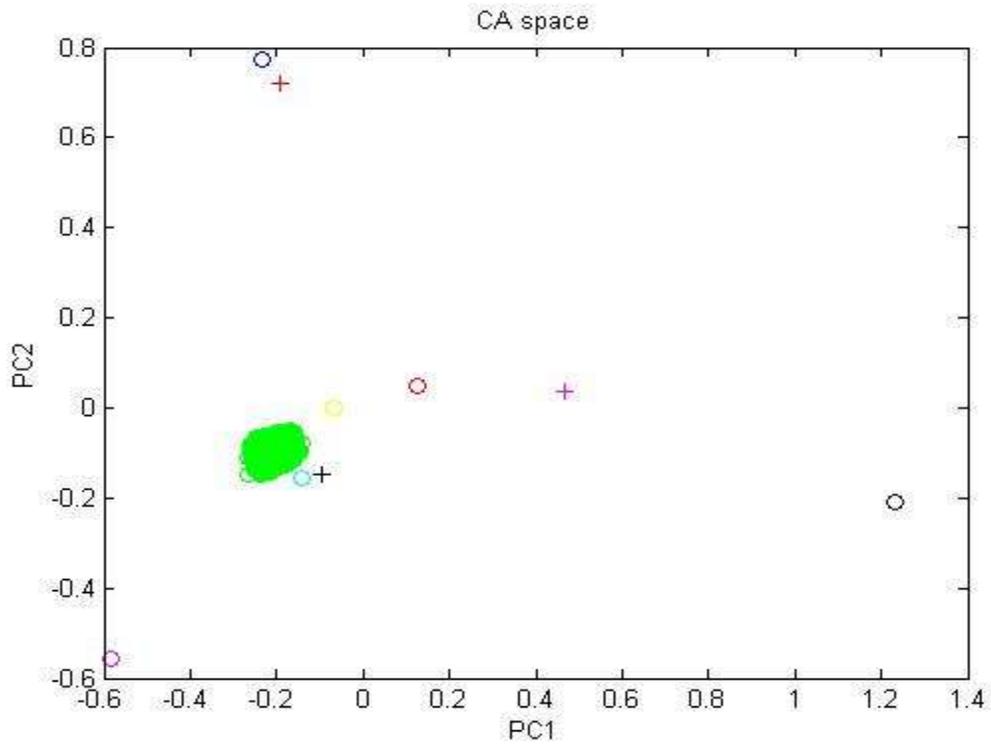Figure 4.17: Number of PCs combined CA model – Depropanizer process



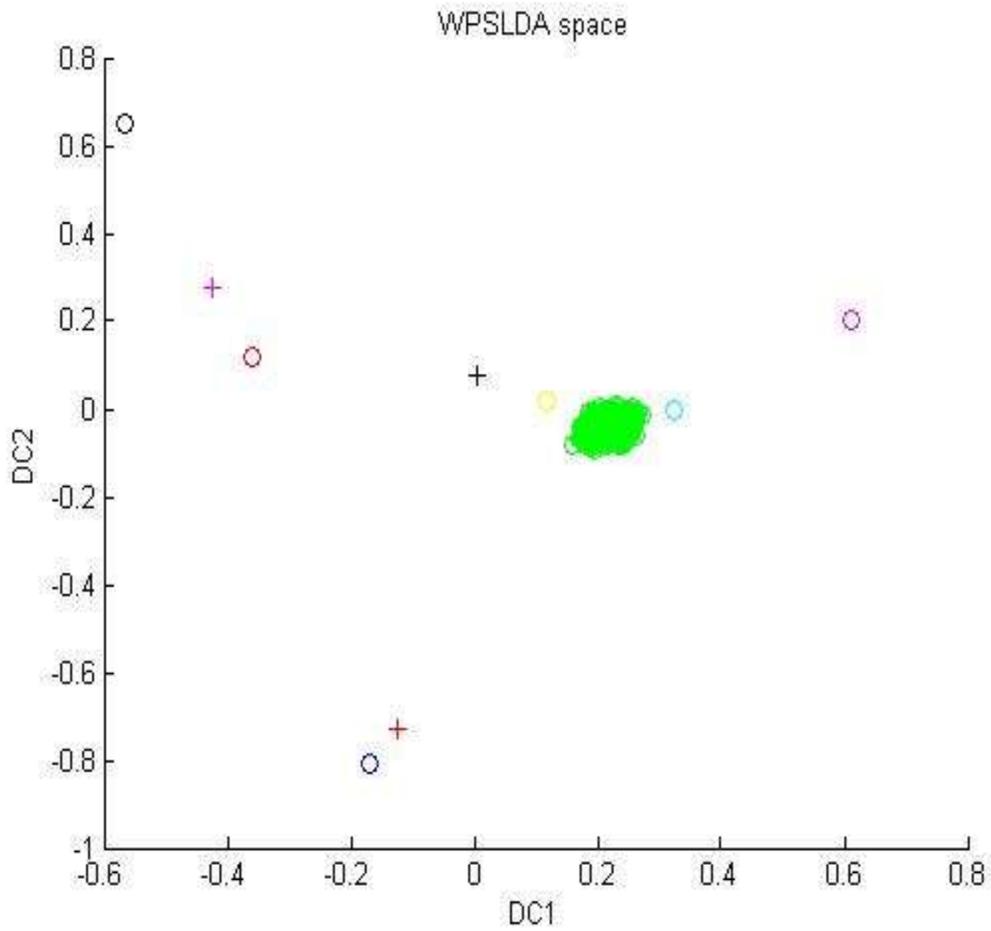Figure 4.18: First 2 PCs of final combined CA model - Depropanizer process

Figure 4.19: Final WPSLDA model – Depropanizer process

The combined CA model was developed with 5 PCs as shown in Figure 4.16 and from Figures 4.17 and 4.18. We can understand that WPSLDA has been effective in moving clusters related to faults 4, 6, and 8 further away from the NOC as compared to the usual CA model.
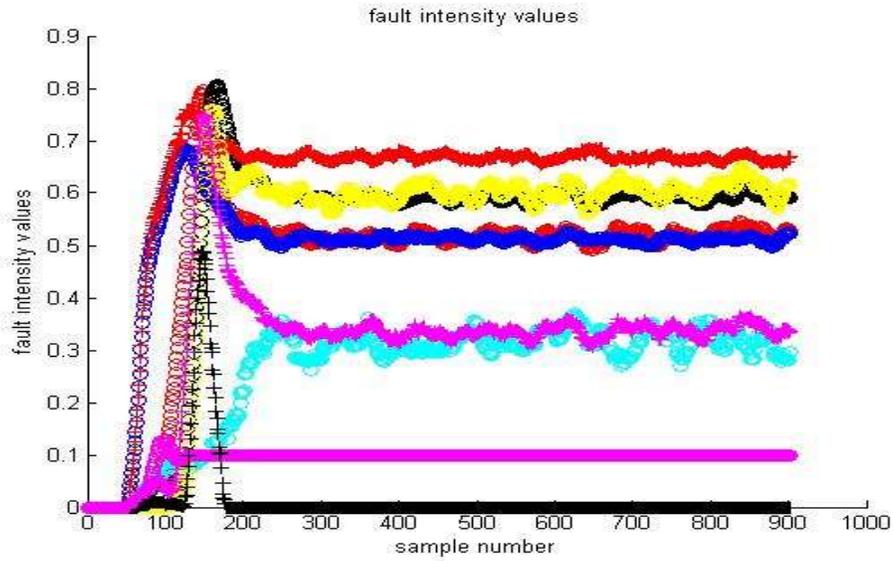
Figure 4.20: Depropanizer process Fault 10 fault intensity



Figure 4.21: Depropanizer process Fault 10 – Individual significant fault intensity values

Figure 4.22: Depropanizer process Fault 11 - Fault intensity values



Figure 4.23: Depropanizer process Fault 11 – Individual significant fault intensity values

Figure 4.24: Depropanizer process Fault 12 – Fault intensity values
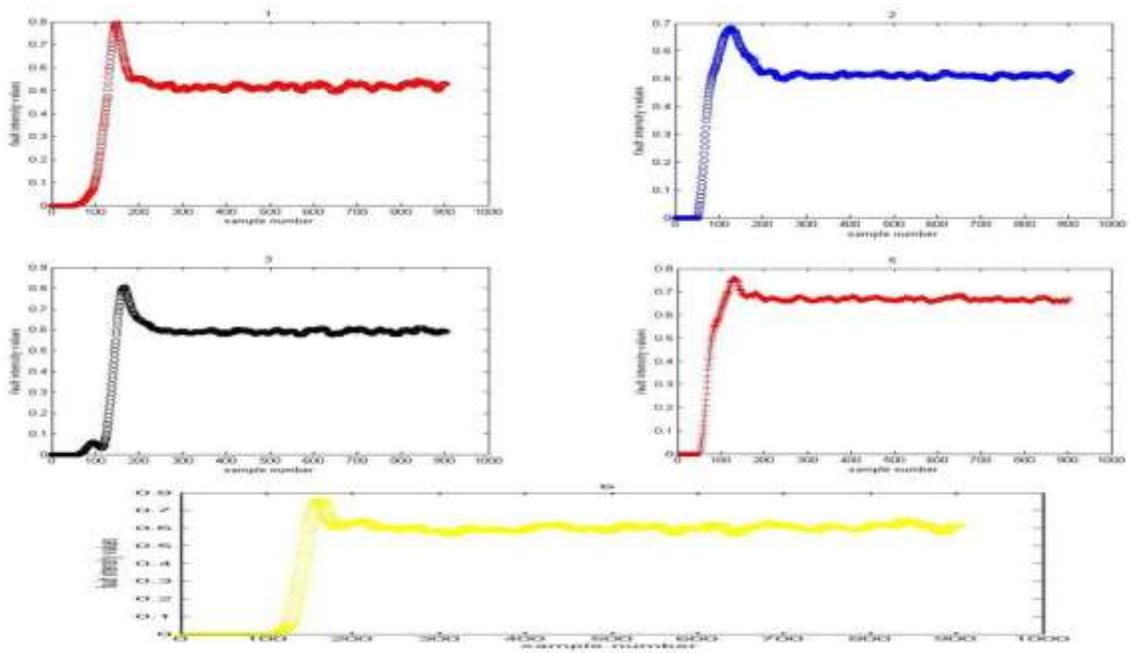


Figure 4.25: Depropanizer process Fault 12 – Individual significant fault intensity values

Figure 4.26: Depropanizer process Fault 13 – Fault intensity values



Figure 4.27: Depropanizer process Fault 13 – Individual significant fault intensity values

Figure 4.28: Depropanizer process Fault 14 – Fault intensity values



Figure 4.29: Depropanizer process Fault 14 – Individual significant fault intensity values
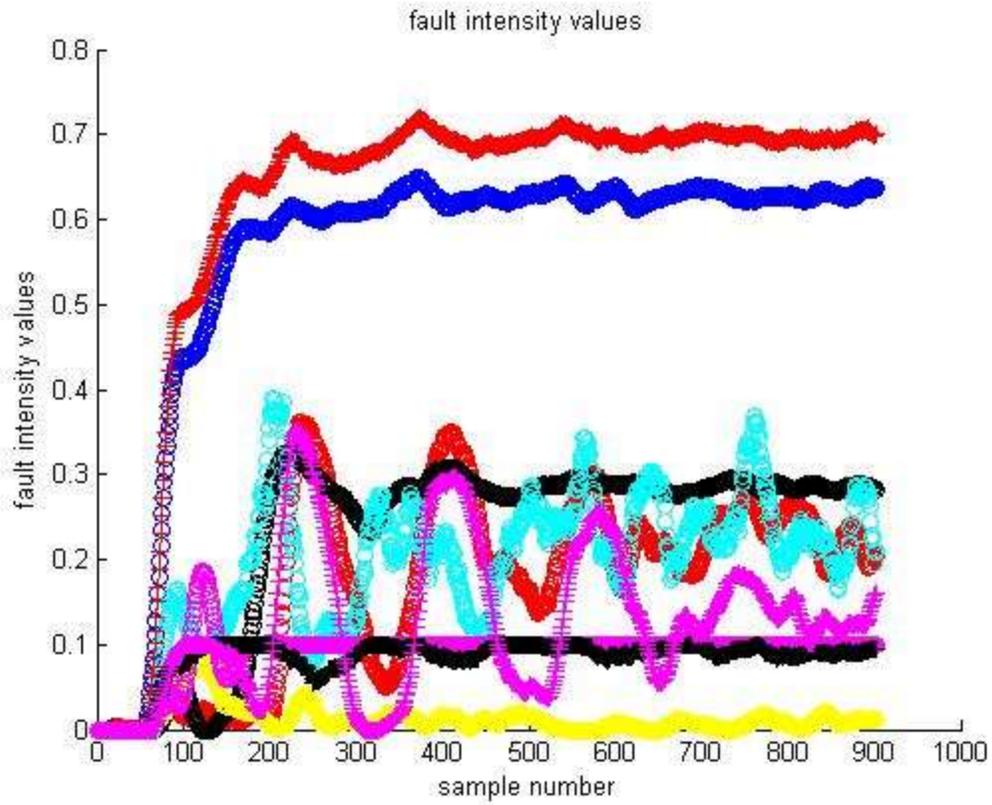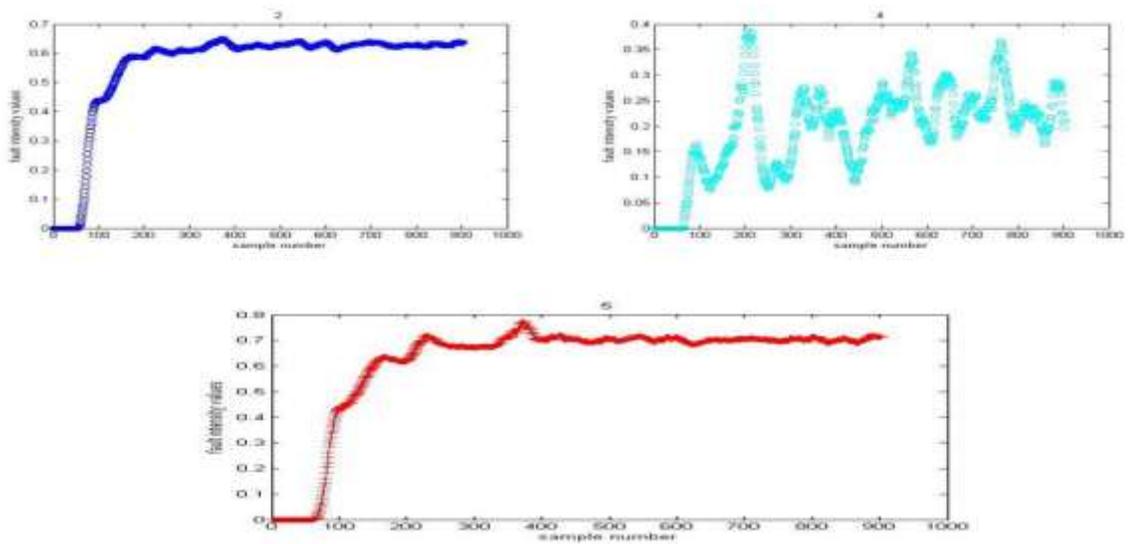
Figure 4.30: Depropanizer process Fault 15 – Fault intensity values



Figure 4.31: Depropanizer process Fault 15 – Individual significant fault intensity values
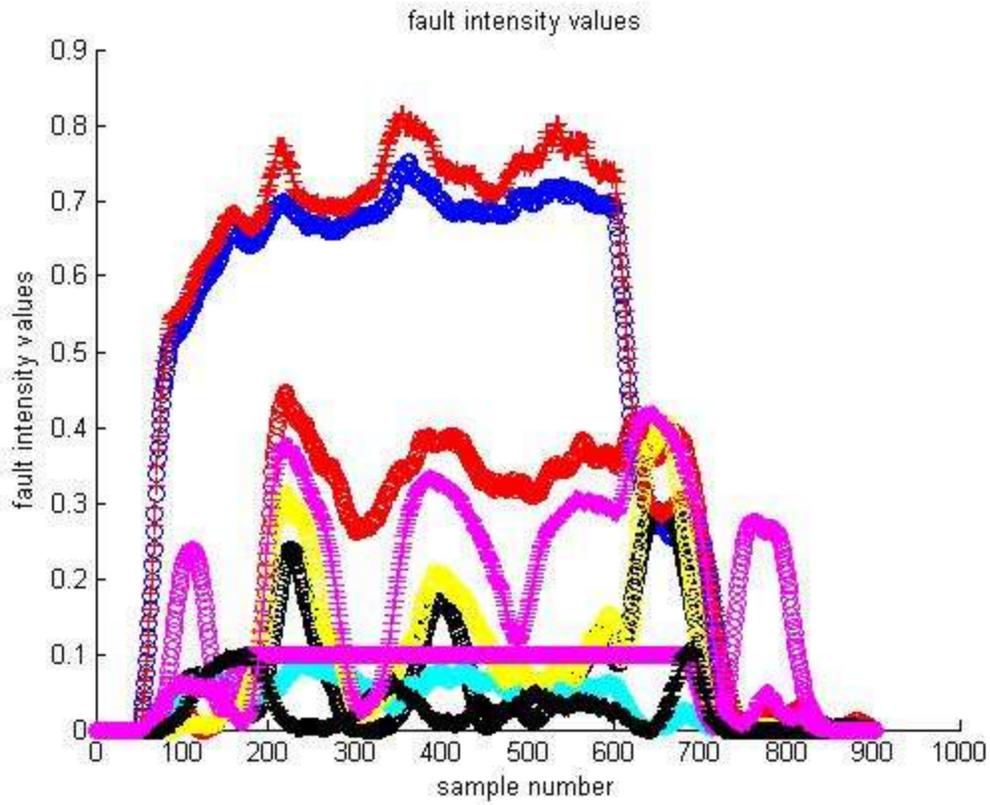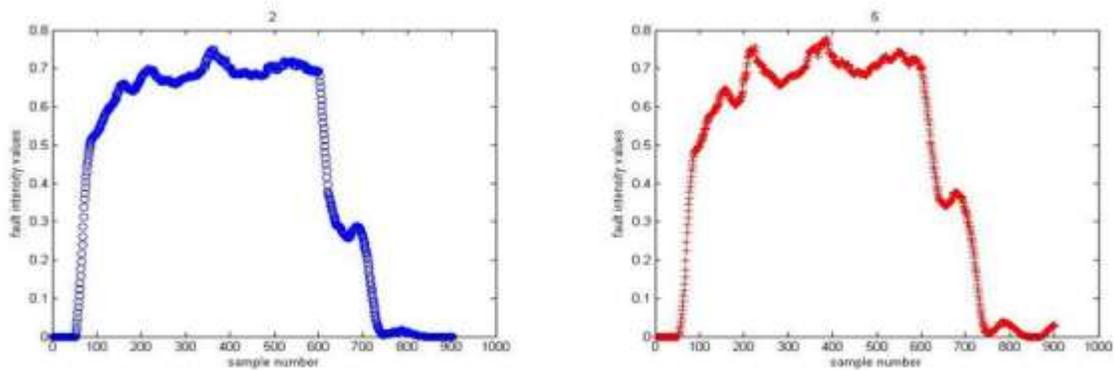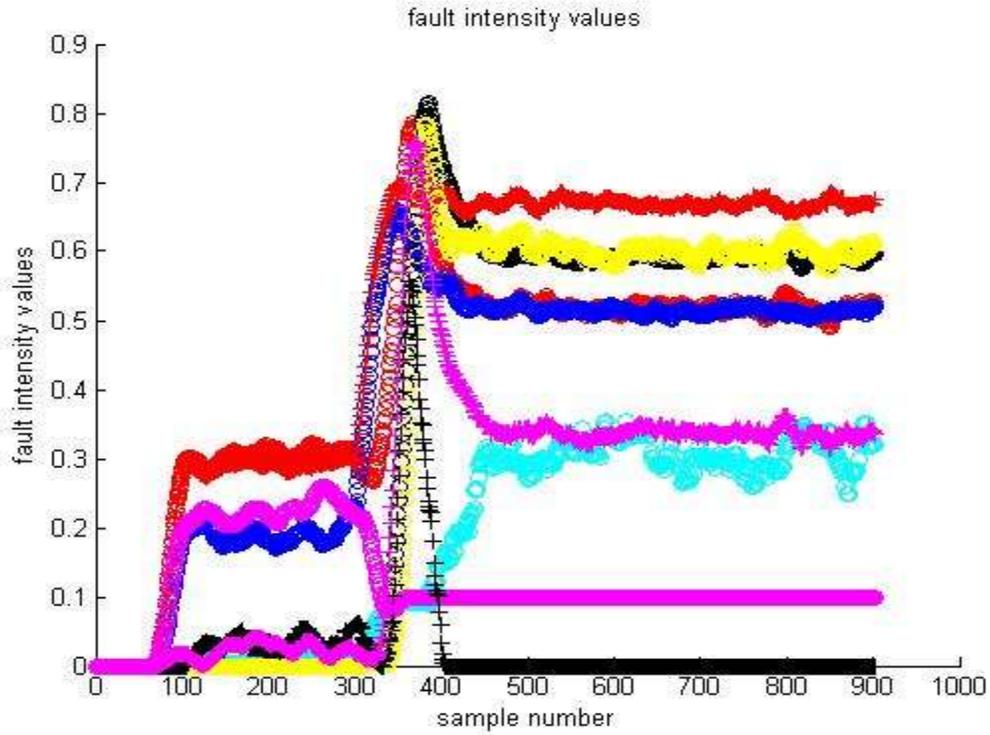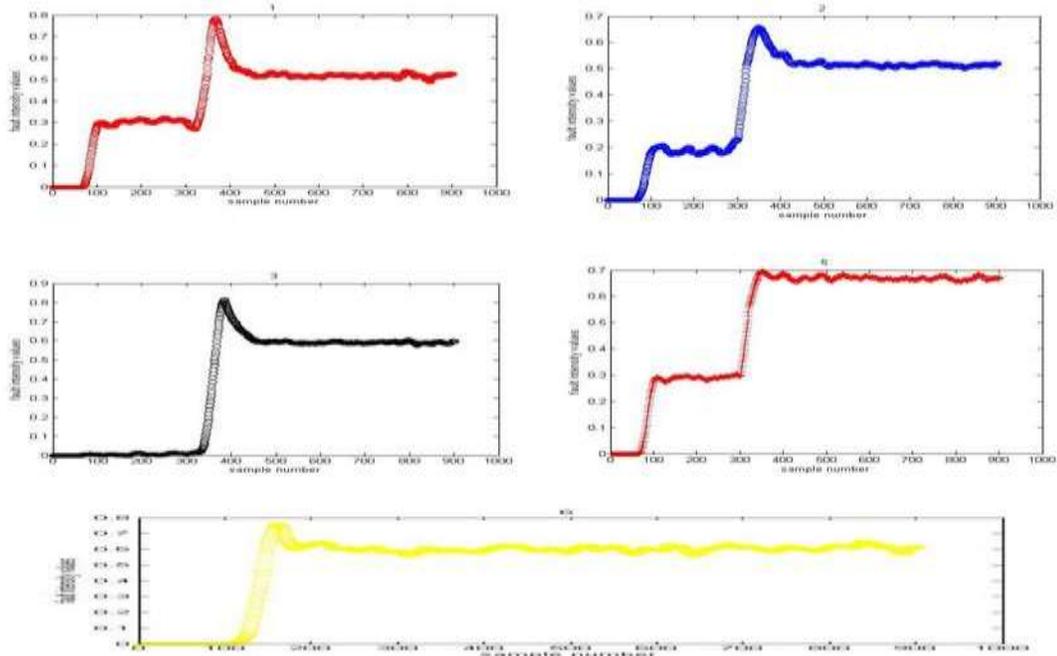
Table 4.5: Depropanizer Process – CA-WPSLDA methodology results

| Fault | Results – fault intensity values | Description of results |
|---|---|---|
| 10 |  | High affinity shown by fault 5 (0.7) followed by secondary contributions from 3 and 6 (0.6). Main affinity is towards fault 5. |
| 11 |  | High affinity shown towards faults 5 and fault 2. Secondary presence is noticed from fault 4 but it has low values (0.2 - 0.4). |
| 12 |  | High affinity is shown towards faults 5 and 2 followed by drops in intensity indicating possible deactivation of fault. |
| 13 |  | High affinity towards fault 5 followed by 6,3,1, and 2. Main variable responsible seems to be fault 5. |
| 14 |  | High affinity towards fault 8 which falls after $600^{th}$ sample indicating deactivation of fault. Fault 7 has a short term contribution after that. |
| 15 |  | Main variable responsible seems to be fault 8, followed closely by fault 9, and then there is drop in intensity values indicating deactivation of fault |

According to the results provided fault 10 seems to have maximum relation to fault 5. Fault 10 is actually a simultaneous occurrence of fault1 and fault2. The results revealed in the CA-WPSLDA methodology are partially correct as fault 2 and fault 5 seem to be very close to each other, this was confirmed by the contribution plot values from Table 3, where both PCA and CA methods showed the same main contribution variables and almost same plots as shown in Figure 4. Fault 2 is the failure of the feed control valve to the tower while fault 5 is related to the loss of feed to the tower.



Figure 4.32: Contribution plots of fault 2 and 5 as calculated in chapter 3

Fault 11 which is actually the simultaneous occurrence of fault 4 and 5 was also better indicated by the methodology as compared to fault 10, where the method shows that there is clearly a strong affinity to fault 5 or 2, and fault 4 shows minor but consistent presence throughout the analysis. Fault 12, a staggered occurrence of fault 2 and 6 only indicates the strong presence of fault 2 or 5 while fault 13 which is a staggered occurrence of fault 1 and 2 only indicates the strong affinity of fault5 but is closely followed by 4 other variables leading to ambiguity in the

109

results. Fault 14 is the occurrence of fault 8 with variable intensity and is rightly indicated as shown in fig 4. Fault 15 which is the occurrence of fault 9 is not shown as the main reason for the occurrence but is only shown as a secondary reason. Therefore, an overall conclusion would be that only one of the faults was most properly indicated while three others which involved two original model faults were partially indicated and one fault (fault 15) was unable to be identified. The main reason attributed to these results could be possible overcrowding of the space and the close relationship between two faults. Another reason could be attributed to the weighted scaling technique employed bringing in the need for a better scaling technique.

# 5. CONCLUSIONS AND RECOMMENDATIONS

## 5.1 Conclusions

From the methods and results described and provided in chapters 2, 3 and 4, it is clear that multivariate statistical techniques such as PCA, PLS and CA are efficient in detection and diagnosis. PCA was found to have both advantages as well as disadvantages in its detection and diagnosis – it offered high detection rates while also resulting in high false alarm rates and more number of contribution variables to consider. Thus, arriving at a correct diagnosis may prove difficult with PCA. CA, on the other hand, was found to be more reliable based on application to several case studies; its only drawback was high detection delays. CA displayed superior discriminative ability which makes it a prime candidate for the development of a comprehensive fault identification methodology that includes multiple fault identifiability. The CA-WPSLDA methodology proposed here showed positive results and promises to work well for novel fault identifiability. Thus it can be said that CA exhibited a strong ability to provide robustness, multiple fault identifiability and novel identifiability in fault monitoring and diagnosis. Therefore, it can be concluded that CA is a powerful potential tool which should be investigated more closely to construct superior process monitoring techniques for the process industry.

## 5.2 Recommendations for Future Work

Based on the results obtained, there are two major areas which could be worthwhile future projects. The first is to develop an improved statistic for CA so as to reduce detection delays associated with it. The currently used Q statistic is found to be a major reason for the high detection delays noticed. The PVR (Principal Component Variable Residual) and CVR (Common Variable Residual) statistics developed by splitting the Q statistic into two parts based

on multiple correlation (Wang *et al.,* 2002) is found to be promising in this regard. This lead could be developed further. The second possible area for future work would be to investigate replacing the WPSLDA technique with a more powerful discriminative tool such as Pareto discriminant analysis (Abou-Moustafa *et al.,* 2010) to separate the fault and normal clusters.

# REFERENCES

1. Abou-Moustafa, K.T., de la Torre, F., and Ferrie, F. P., *Pareto discriminant analysis*. in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2010.

2. Baldassarre, M.T., Caivano, D., Kitchenham,B., Visaggio,G., *Systematic review of statistical process control: An experience report.* in *11th International Conference on Evaluation and Assessment in Software Engineering* 2007. UK.

3. Baxter, M.J., *Exploratory Multivariate Analysis in Archaeology.* 1994, Edinburgh: Edinburgh University Press.

4. Bersimis, S., Psarakis, S., Panaretos, J., *Multivariate statistical process control charts: an overview.* Quality and Reliability Engineering International, 2007. 23(5): p. 517-543.

5. Boik, R., *An efficient algorithm for joint correspondence analysis.* Psychometrika, 1996. 61(2): p. 255-269.

6. Bro, R. and A.K. Smilde, *Centering and scaling in component analysis.* Journal of Chemometrics, 2003. 17(1): p. 16-33.

7. Carroll, J.D., Green, P. E., Schaffer, C. M., *Reply to Greenacre's Commentary on the Carroll-Green-Schaffer Scaling of Two-Way Correspondence Analysis Solutions.* Journal of Marketing Research, 1989. 26(3): p. 366-368.

8. Chang, C.C., & Yu, C. C., *On-line fault diagnosis using the signed directed graph.* Industrial and Engineering Chemistry Research, 1990 29(7): p. 1290-1299.

9. Chester, D., Lamb, D., Dhurjati, P., *Rule-based computer alarm analysis in chemical process plants.* in *In Proceedings of 7th Micro-Delcon.* 1984.

10.     Cheung, J.T.Y., Stephanopoulos, G., *Representation of process trends--Part I. A formal representation framework.* Computers & Chemical Engineering, 1990. 14(4-5): p. 495-510.

11.     Chiang, L.H., E.L. Russell, and R.D. Braatz, *Fault diagnosis in chemical processes using Fisher discriminant analysis, discriminant partial least squares, and principal component analysis.* Chemometrics and Intelligent Laboratory Systems, 2000. 50(2): p. 243-252.

12.     Choi, S.W., Lee, I-B., *Nonlinear dynamic process monitoring based on dynamic kernel PCA.* Chemical Engineering Science, 2004. 59(24): p. 5897-5908.

13.     Chow, E.Y., A.S. Willsky, *Analytical redundancy and the design of robust failure detection systems* IEEE Transactions on Automatic Control., 1984. 29(7): p. 603-614.

14.     Clausen, S.-E., *Applied Correspondence Analysis: an introduction (Quantitative applications in the social sciences).* 1998, Sage Publications: Thousand Oaks, California, USA.

15.     Clouse, R.A., *Interpreting Archaeological Data through Correspondence Analysis.* Historical Archaeology, 1999. 33(2): p. 90-107.

16.     de Jong, S., *SIMPLS: An alternative approach to partial least squares regression.* Chemometrics and Intelligent Laboratory Systems, 1993. 18(3): p. 251-263.

17.     De Kleer, J. and J.S. Brown, *A qualitative physics based on confluences.* Artificial Intelligence, 1984. 24(1-3): p. 7-83.

18.     Denney, D.W., MacKay, J., MacHattie,T., Flora,C., and Mastracci, E., *Application of pattern recognition techniques to process unit data.*, in *Canadian Society of Chemical Engineering Conference.* 1985: Sarnia, Ontario, Canada.

19.     Detroja, K.P., Gudi,R.D., Patwardhan, S.C., Roy, K., *Fault detection and isolation using correspondence analysis.* Industrial and Engineering Chemistry Research, 2006. 45(1): p. 223-235.

20.     Detroja, K.P., Gudi, R. D., Patwardhan, S. C., *Plant-wide detection and diagnosis using correspondence analysis.* Control Engineering Practice, 2007. 15(12): p. 1468-1483.

21.     Ding, C. et al., *Unsupervised Learning: Self-aggregation in Scaled Principal Component Space*, in *Principles of Data Mining and Knowledge Discovery*, T. Elomaa, H. Mannila, and H. Toivonen, Editors. 2002, Springer Berlin / Heidelberg. p. 79-112.

22.     Dong, D., McAvoy, T. J., *Batch tracking via nonlinear principal component analysis.* AIChE Journal, 1996. 42(8): p. 2199-2208.

23.     Downs, J.J., and Vogel, E. F., *A plant-wide industrial process control problem.* Computers & Chemical Engineering, 1993. 17(3): p. 245-255.

24.     Frank, P.M., and Wunnenberg, J., *Robust fault diagnosis using unknown input observer schemes.*, in *Fault diagnosis in dynamic systems: theory and applications*, R. J. Patton, Editor. 1989, Prentice Hall: NY, USA.

25.     Frank, P.M., *Fault diagnosis in dynamic systems using analytical and knowledge-based redundancy-a survey and some new results.* Automatica, 1990. 26(3): p. 459-474.

26.     Fussell, J.B., Powers, G.J., Bennetts, R. G., *Fault Trees Analysis-A State of the Art Discussion.* IEEE Transactions on Reliability, 1974. R-23(1): p. 51-55.

27.     Geladi, P., Kowalski, B.R., *Partial least-squares regression: a tutorial.* Analytica Chimica Acta, 1986. 185: p. 1-17.

28.     Gertler, J., *Intelligent supervisory control*, in *Artificial intelligence handbook.*, A.E.N.J.R. Davis, Editor. 1989, Research Triangle Park: NC, USA.

29.     Gertler, J. *Analytical redundancy methods in fault detection and isolation*. in *Proceedings of IFAC/IAMCS symposium on safe process*. 1991.

30.     Gertler, J., *Analytical redundancy methods in fault detection and isolation-survey and synthesis*, in *IFAC symposium on online fault detection and supervision in the chemical process industries*. 1992.

31.     Gertler, J., *Residual generation in model-based fault diagnosis.* Control-theory and advanced technology, 1993. 9(1): p. 259-285.

32.     Gertler, J., *Fault detection and diagnosis in engineering systems*. 1998: Marcel Dekker.

33.     Greenacre, M.J., *Correspondence analysis in medical research.* Statistical Methods in Medical Research, 1992. 1(1): p. 97-117.

34.     Greenacre, M. and T. Hastie, *The Geometric Interpretation of Correspondence Analysis.* Journal of the American Statistical Association, 1987. 82(398): p. 437-447.

35.     Greenacre, M.J., *Theory and Applications of Correspondence Analysis*. 1984, Academic Press, London, UK.

36.     Greenacre, M.J., *Correspondence analysis of multivariate categorical data by weighted least-squares.* Biometrika, 1988. 75(3): p. 457-467.

37.     Greenacre, M.J., *Correspondence Analysis in Practice,*. 1993, London: Academic Press.

38.     He, Q.P., S.J. Qin, and J. Wang, *A new fault diagnosis method using fault directions in Fisher discriminant analysis.* AIChE Journal, 2005. 51(2): p. 555-571.

39.     He, X.B., et al., *Variable-weighted Fisher discriminant analysis for process fault diagnosis.* Journal of Process Control, 2009. 19(6): p. 923-931.

40.     Hill, M.O., *Correspondence Analysis: A Neglected Multivariate Method.* Journal of the Royal Statistical Society. Series C (Applied Statistics), 1974. 23(3): p. 340-354.

41. Hill, M.O., Gauch, H. G., *Detrended correspondence analysis: An improved ordination technique.* Plant Ecology, 1980. 42(1): p. 47-58.

42. Himmelblau, D.M., *Fault detection and diagnosis in chemical and petrochemical processes / David M. Himmelblau.* Chemical engineering monographs ; v. 8. 1978, Elsevier Scientific Pub. Co. USA.

43. Hotelling, H., *The economics of exhaustible resources.* Journal of Political Economy, 1931. 39: p. 137-175.

44. Huang, H.-P., C.-C. Li, and J.-C. Jeng, *Multiple Multiplicative Fault Diagnosis for Dynamic Processes via Parameter Similarity Measures.* Industrial & Engineering Chemistry Research, 2007. 46(13): p. 4517-4530.

45. Hunter, J.S., *The Exponentially Weighted Moving Average.* Journal of Quality Technology, 1986. 18: p. 203-210.

46. Iri, M.A., K. O'Shima, E. Matsuyama, H., *An algorithm for diagnosis of system failures in the chemical process.* Computers & Chemical Engineering, 1979. 3(1-4): p. 489-493.

47. Isermann, R., *Process fault diagnosis based on dynamic models and parameter estimation methods*, in *Fault diagnosis in dynamic systems: theory and applications*, P.M.F. R. J. Patton and R. N. Clark, Editors. 1989, Prentice Hall: NY.

48. Iwasaki, Y., Simon,H.A., *Causality in device behavior.* Artif. Intell., 1986. 29(1): p. 3-32.

49. Jackson, J.E., *A User's Guide to Principal Components*. 1991, NY: Wiley.

50. Janusz, M.E., Venkatasubramanian, V., *Automatic generation of qualitative descriptions of process trends for fault detection and diagnosis.* Engineering Applications of Artificial Intelligence, 1991. 4(5): p. 329-339.

51.     Jemwa, G.T., Aldrich, C., *Kernel-based fault diagnosis on mineral processing plants.* Minerals Engineering, 2006. 19(11): p. 1149-1162.

52.     Jiang, Z., X. He, and Y. Yang. *Key variable identification using discriminant analysis*. in Proceedings of the 27th Chinese *Control Conference, CCC* 2008.

53.     Johansson, K.H., *The quadruple-tank process: a multivariable laboratory process with an adjustable zero.* IEEE Transactions on Control Systems Technology, 2000. 8(3): p. 456-465.

54.     Jollife, I., *Principal Component Analysis*. 1986, Verlag, New York: Springer.

55.     Kaspar, M.H. and W. H. Ray, *Dynamic PLS modelling for process control.* Chemical Engineering Science, 1993. 48(20): p. 3447-3461.

56.     Khare, S., Bavdekar, V., Kadu S.C, Detroja, K. and Gudi, R.D., *Scaling and Monitoring Issues in Monitoring and fault detection and diagnosis.* In Proceedings of DYCOPS, 2007.

57.     Kramer, M.A., *Nonlinear principal component analysis using autoassociative neural networks.* AIChE Journal, 1991. 37(2): p. 233-243.

58.     Kresta, J.V., Macgregor, J. F., Marlin, T. E., *Multivariate statistical monitoring of process operating performance.* The Canadian Journal of Chemical Engineering, 1991. 69(1): p. 35-47.

59.     Kuipers, B., *Qualitative simulation.* Artificial Intelligence, 1986. 29(3): p. 289-338.

60.     Kumamoto, H., Ikenchi, K., Inoue, K., Henley, E. J., *Application of expert system techniques to fault diagnosis.* The Chemical Engineering Journal, 1984. 29(1): p. 1-9.

61. Lakshminarayanan, S., Shah, S. L., Nandakumar, K., *Modeling and control of multivariable processes: Dynamic PLS approach.* AIChE Journal, 1997. 43(9): p. 2307-2322.

62. Lapp, S.A., Powers, G. J., *Computer-aided Synthesis of Fault-trees.* Reliability, IEEE Transactions on, 1977. R-26(1): p. 2-13.

63. Lee, J.-M., Yoo, C.K., Lee, I-B.,, *Fault detection of batch processes using multiway kernel principal component analysis.* Computers & Chemical Engineering, 2004. 28(9): p. 1837-1847.

64. Li, W., Shah, S.L. *Fault detection and isolation in non-uniformly sampled systems*. in *IFAC DYCOPS 7*. Cambridge, MA, USA.

65. Li, W., Yue, H. H., Valle-Cervantes, S., Qin, S. J.,, *Recursive PCA for adaptive process monitoring.* Journal of Process Control, 2000. 10(5): p. 471-486.

66. Li, W. and S. Shah, *Structured residual vector-based approach to sensor fault detection and isolation.* Journal of Process Control, 2002. 12(3): p. 429-443.

67. Li, Y., Y. Gao, and H. Erdogan, *Weighted pairwise scatter to improve linear discriminant analysis*, in *In ICSLP-2000*. 2000: Beijing, China. p. 608 - 611.

68. Lyman, P.R. and C. Georgakis, *Plant-wide control of the Tennessee Eastman problem.* Computers & Chemical Engineering, 1995. 19(3): p. 321-331.

69. MacGregor, J.F., Kourti, T., *Statistical process control of multivariate processes.* Control Engineering Practice, 1995. 3(3): p. 403-414.

70. Mason, R., Young, J., *Multivariate Statistical Process Control with Industrial Applications*. ASA-SIAM. 2002.

71.     Niida, K., Itoh,J., Umeda,T., Kobayashi,S., Ichikawa, A., *Some Expert System Experiments in Process Engineering.* Chemical Engineering Research and Design 1986. 64a: p. 372-380.

72.     Nomikos, P., MacGregor, J. F., *Monitoring batch processes using multiway principal component analysis.* AIChE Journal, 1994. 40(8): p. 1361-1375.

73.     Oyeleye, O.O., Kramer, M. A., *Qualitative simulation of chemical process systems: Steady-state analysis.* AIChE Journal, 1988. 34(9): p. 1441-1454.

74.     Patel, S.R., Gudi, R. D., *Improved monitoring and discrimination of batch processes using correspondence analysis*, in *Proceedings of the 2009 conference on American Control Conference*. 2009, IEEE Press: St. Louis, Missouri, USA. p. 3434-3439.

75.     Potter, J.E., Suman, M.C.,. *Thresholdness redundancy management with arrays of skewed instruments.* in *Integrity in Electronic Flight Control Systems*. 1977: AGARDOGRAPH-224.

76.     Pusha, S., Gudi, R., Noronha, S.,, *Polar classification with correspondence analysis for fault isolation.* Journal of Process Control, 2009. 19(4): p. 656-663.

77.     Qin, J.S., *Statistical process monitoring: basics and beyond.* Journal of Chemometrics, 2003. 17(8-9): p. 480-502.

78.     Qin, S.J., McAvoy, T. J., *Nonlinear PLS modeling using neural networks.* Computers & Chemical Engineering, 1992. 16(4): p. 379-391.

79.     Qin, S.J., *Recursive PLS algorithms for adaptive data modeling.* Computers & Chemical Engineering, 1998. 22(4-5): p. 503-514.

80.     Raich, A.C., and Çinar, A. *Statistical process monitoring and disturbance isolation in multivariate continuous processes* in Proceedings of *ADCHEM*. 1994, pages 452-457.

81. Rännar, S., Lindgren, F., Geladi, P. and Wold, S., *A PLS Kernel Algorithm for Data Sets with Many Variables and Fewer Objects. Part 1: Theory and Algorithm*. Journal of Chemometrics, 1194. 8: p. 111–125.

82. Rännar, S., MacGregor, J.F., Wold, S, *Adaptive batch monitoring using hierarchical PCA*. Chemometrics and Intelligent Laboratory Systems, 1998. 41(1): p. 73-81.

83. Rich, S.H., Venkatasubramanian,V., Nasrallah,M., Matteo,C., *Development of a diagnostic expert system for a whipped toppings process.* Journal of Loss Prevention in the Process Industries, 1989. 2(3): p. 145-154.

84. Roberts, S.W., *Control chart tests based on geometric moving averages.* Technometrics, 1959. 1: p. 239-250.

85. Romagnoli, J.A., Stephanopoulos, G.,, *Rectification of process measurement data in the presence of gross errors.* Chemical Engineering Science, 1981. 36(11): p. 1849-1863.

86. Russell, E.L., Chiang, L H., Braatz, R. D., *Fault detection in industrial processes using canonical variate analysis and dynamic principal component analysis.* Chemometrics and Intelligent Laboratory Systems, 2000. 51(1): p. 81-93.

87. Sacks, E., *Qualitative analysis by piecewise linear approximation.* Artificial Intelligence in Engineering, 1988. 3(3): p. 151-155.

88. Schölkopf, B., Smola, A., Müller, K-R., *Kernel principal component analysis*, in *Artificial Neural Networks — ICANN'97*, W. Gerstner, et al., Editors. 1997, Springer Berlin / Heidelberg. p. 583-588.

89. Seasholtz, M.B., Pell, R. J., Gates, K. E., *Comments on the power method.* Journal of Chemometrics, 1990. 4(4): p. 331-334.

90. Shewhart, W.A., *Economic control of quality of manufactured product D. Van Nostrand Co. Inc.,* 1931, USA.

91.     Shiozaki, J., Matsuyama, H., O'Shima, E., Iri, M., *An improved algorithm for diagnosis of system failures in the chemical process.* Computers & Chemical Engineering, 1985. 9(3): p. 285-293.

92.     Simon, H.A., *Models of discovery.* 1977, Reidel Publishing Company, Boston, USA.

93.     Simpson, E.H., *The Interpretation of Interaction in Contingency Tables.* Journal of the Royal Statistical Society. Series B (Methodological), 1951. 13(2): p. 238-241.

94.     Sparks, R.S., *Quality Control With Multivariate Data.* Australian & New Zealand Journal of Statistics, 1992. 34(3): p. 375-390.

95.     ter Braak, C.J.F., *Canonical Correspondence Analysis: A New Eigenvector Technique for Multivariate Direct Gradient Analysis.* Ecology, 1986. 67(5): p. 1167-1179.

96.     ter Braak, C.J.F., *Ordination*, in *Data Analysis in Community and Landscape Ecology*, R.H. Jongman, ter Braak,C .J.F., van Tongeren, O.F.R., Editors. 1987, Pudoc: Wageningen. p. 91-173.

97.     Umeda, T., Kuriyama, T., O'Shima, E., Matsuyama, H., *A graphical approach to cause and effect analysis of chemical processing systems.* Chemical Engineering Science, 1980. 35(12): p. 2379-2388.

98.     Venkatasubramanian, V., Rengaswamy, R., Yin, K., Kavuri, S.N., *A review of process fault detection and diagnosis: Part I: Quantitative model-based methods.* Computers & Chemical Engineering, 2003a. 27(3): p. 293-311.

99.     Venkatasubramanian, V., Rengaswamy, R., Kavuri, S. N., Yin, K., *A review of process fault detection and diagnosis: Part III: Process history based methods.* Computers & Chemical Engineering, 2003b. 27(3): p. 327-346.

100. Venkatasubramanian, V., Rengaswamy, R., Kavuri, S. N., Yin, K., *A review of process fault detection and diagnosis: Part II: Qualitative models and search strategies.* Computers & Chemical Engineering, 2003c. 27(3): p. 313-326.

101. Vijaysai, P., R.D. Gudi, and S. Lakshminarayanan, *Identification on Demand Using a Blockwise Recursive Partial Least Squares Technique.* Industrial & Engineering Chemistry Research, 2003. 42(3): p. 540-554.

102. Wang, H., Song, Z., Wang, H., *Statistical process monitoring using improved PCA with optimized sensor locations.* Journal of Process Control, 2002. 12(6): p. 735-744.

103. Willsky, A.S., *A survey of design methods for failure detection in dynamic systems.* Automatica, 1976. 12(6): p. 601-611.

104. Wise, B.M., Veltkamp,D.J., Ricker,N.L., B.R. Kowalski, Barnes,S., and Arakali,V.,. *Application of multivariate statistical process control (MSPC) to the West Valley slurry-red ceramic melter process*. in *Waste Management '91 Proceedings*. 1991. Tucson, AZ University of Arizona Press.

105. Wold, S., Geladi, P., Esbensen, K., Öhman, J., *Multi-way principal components-and PLS-analysis.* Journal of Chemometrics, 1987. 1(1): p. 41-56.

106. Woodall, W.H., *Controversies and Contradictions in Statistical Process Control.* Journal of Quality Technology, 2000. 32(4): p. 341-350.

107. Woodward, R.H., Goldsmith,P.L., *Cumulative sum techniques: mathematical and statistical techniques for industry.* ICI - Monograph, 1964.

108. Yang, J. and J.-y. Yang, *Why can LDA be performed in PCA transformed space?* Pattern Recognition, 2003. 36(2): p. 563-566.

109.    Zhou, D., Li, G., Qin, S. J., *Total projection to latent structures for process monitoring.* AIChE Journal, 2010. 56(1): p. 168-178.

110.    Benzecri, J.P., *L'Analyse des Donnees, Tome 2: L'Analyse des Correspondance*. 1973, Paris: Dunod.