

DEVELOPMENT OF MULTIVARIATE STATISTICAL PROCESS MONITORING  
USING COMBINATION MLR-PCA METHOD

NORLIZA BINTI ISHAK

Thesis submitted in fulfillment of the requirement  
for the award of the degree of Master of Engineering in Chemical

Faculty of Chemical and Natural Resources Engineering  
UNIVERSITI MALAYSIA PAHANG

OCTOBER 2015

## ABSTRACT

The most popular types of process monitoring systems is Multivariate Statistical Process Monitoring (MSPM) which has the most practical method in handling the complicated large scale processes. This is due to the ability of the system in maximizing the usage of abundant historical process data, in such a way that the original data dimensions are compressed and data variations preserved to certain extent in a set of transformed variables by way of linear combinations. Thus, the composite model is generally flexible regardless of the amount of variables that utilized. In this regard, conventional Principal Component Analysis (PCA) has been widely applied to conduct such compression function particularly for MSPM. However, the conventional PCA is a linear technique which results sometimes inappropriately employed especially in modeling processes that exhibit severe non-linear correlations. Therefore, a new solution is demanded, whereby the number of original variables can be reduced to certain extent (in terms of scales), while it still can maintain the variation as maximally as possible corresponding to the original, which are then transferable into monitoring statistics. One of the potential techniques available in addressing the issue is known as Multiple Linear Regression (MLR). The main objective of the technique is to predict a set of output values (criterion) based from a specified of linear function, which consists of a set of predictor. Therefore, the main multivariate data will be divided into two groups, which are the criterion and predictor categories. The study adopts, Tennessee Eastman Process (TEP) and Multiple Output and Multiple Input (MIMO) Pilot Plant System for demonstration. The general finding is that MLR-PCA normally employs less number of PCs compared to PCA, and thus, this will perhaps reduce complication during diagnosis. By adopting such approach, the monitoring task can be made simpler and perhaps more effective, in the sense that only those selected criterion variables (predicted values) will be taken for monitoring, while preserving the rest of the predictor value trends in the form of linear functions in association with the criterion variables. This study also shows that MLR-PCA works relatively better in terms of fault detection and identification against the conventional system.

## ABSTRAK

Sistem pemantauan proses yang paling popular ialah Pembolehubah Pemantauan Proses Statistik (MSPM) yang dianggap sebagai kaedah yang paling praktikal dalam menanggapi sistem kompleks serta berskala besar. Hal ini adalah kerana keupayaan sistem dalam memaksimumkan penggunaan data proses dalam apa-apa cara bahawa dimensi data asal dimampatkan dan variasi data dikekalkan pada tahap tertentu dalam satu set pembolehubah yang diubah melalui kombinasi linear. Oleh itu, model komposit biasanya fleksibel tanpa mengira jumlah pembolehubah yang digunakan. Dalam hal ini, Analisis Komponen Utama konvensional (PCA) telah digunakan secara meluas untuk menjalankan fungsi mampatan itu terutamanya untuk MSPM. Walau bagaimanapun, PCA konvensional adalah teknik linear yang mana kaedah ini kadang-kadang tidak sesuai digunakan terutama dalam proses model kerana mempamerkan korelasi bukan linear yang tidak normal. Oleh itu, satu penyelesaian baru diperlukan, di mana bilangan pembolehubah asal boleh dikurangkan kepada tahap tertentu (dari segi skala), sementara ia masih boleh mengekalkan perubahan maksima sebagaimana ia mungkin sama dengan data yang asal, yang kemudiannya dipindahkan ke dalam statistik pemantauan. Salah satu teknik yang berpotensi untuk menangani isu ini dikenali sebagai Regresi Linear Berganda (MLR). Objektif utama teknik ini untuk meramalkan satu set nilai hasil (kriteria) berdasarkan dari yang dinyatakan fungsi linear, yang terdiri daripada satu set peramal. Oleh itu, data multivariat utama akan dibahagikan kepada dua kumpulan, iaitu kriteria dan kategori peramal. Kajian ini mengadaptasi Proses Tennessee Eastman (TEP) dan Multiple Output and Multiple Input (MIMO) Loji Pandu sebagai demonstrasi. Kajian mendapati MLR-PCA biasanya menggunakan PC nombor yang minimum berbanding PCA dan sekaligus, hal ini berkemungkinan dapat mengurangkan kerumitan semasa diagnosis. Dengan menggunakan pendekatan ini, tugas pemantauan boleh dibuat lebih mudah dan mungkin lebih berkesan, dalam erti kata lain bahawa hanya pembolehubah kriteria tertentu sahaja (nilai meramalkan) akan diambil untuk dianalisis, di samping memelihara seluruh trend nilai peramal dalam bentuk fungsi linear bersama dengan pembolehubah kriteria. Kajian ini juga menunjukkan MLR-PCA berfungsi lebih baik dari segi pengesanan ralat dan mengenal pasti ralat berbanding kaedah konvensional.

## TABLE OF CONTENT

	<b>Page</b>
<b>TITLE PAGE</b>	i
<b>SUPERVISOR’S DECLARATION</b>	ii
<b>STUDENT’S DECLARATION</b>	iii
<b>DEDICATION</b>	vi
<b>ACKNOWLEDGE</b>	v
<b>ABSTRACT</b>	vi
<b>ABSTRAK</b>	vii
<b>TABLE OF CONTENTS</b>	viii
<b>LIST OF TABLES</b>	xi
<b>LIST OF FIGURES</b>	xii
<b>LIST OF ABBREVIATIONS</b>	xiv
<b>LIST OF SYMBOLS</b>	xvi
<b>CHAPTER 1      INTRODUCTION</b>	<b>1</b>
1.1      Introduction	1
1.2      Research Background	2
1.3      Problem Statement	3
1.4      Research Objectives	4
1.5      Research Scope	5
1.6      Rationale and Significant	6
1.7      Report Organization	7

<b>CHAPTER 2</b>	<b>LITERATURE REVIEW</b>	<b>8</b>
2.1	Introduction	8
2.2	Fundamental of MSPM	9
	2.2.1 Principal Component Analysis	11
2.3	Industrial Application of MSPM	13
	2.3.1 Applications of MSPM Based on PCA Extension	15
	2.3.2 Applications of MSPM Based on Multivariate Compression Techniques	18
	2.3.2.1 Partial Least Square	18
	2.3.2.2 Independent Component Analysis	20
	2.3.2.3 Factor Analysis	21
	2.3.2.4 Principal Component Regression	22
	2.3.2.5 Multidimensional Scaling	23
2.4	Multiple Linear Regression	25
<b>CHAPTER 3</b>	<b>METHODOLOGY</b>	<b>31</b>
3.1	Introduction	31
3.2	Framework I : PCA-MSPM System Procedures	32
	3.2.1 Phase I Procedures	33
	3.2.2 Phase II Procedures	38
3.3	Framework II : PCA-MLR MSPM Framework	40
	3.3.1 Phase I Procedures	41
	3.3.2 Phase II Procedures	43
3.4	Case Study 1: Tennessee Eastman Process	44
3.5	Case Study 2: MIMO Training System Pilot Plant	47

<b>CHAPTER 4</b>	<b>RESULTS AND DISCUSSION</b>	<b>52</b>
4.1	Introduction	52
4.2	Case Study 1: Tennessee Eastman Process	53
4.2.1	Phase I Results	53
4.2.1.1	PCA Models	54
4.2.2.2	False Alarm Rate Analysis	56
4.2.2	Phase II Results	57
4.2.2.1	Fault Detection Time	58
4.2.2.2	Fault Identification	61
4.2.2.3	Case Study on Fault 1	63
4.2.2.4	Case Study on Fault 4	71
4.2.2.5	Case Study on Fault 11	77
4.3	Case Study 2: MIMO Training System Plant	83
4.3.1	Phase I Results	83
4.3.1.1	PCA Models	83
4.3.2	Phase II Results	85
4.3.2.1	Fault Detection Time	85
4.3.2.2	Fault Identification	86
4.3.2.3	Case Study on Fault 1(TE (BURNT OUT))	87
4.4	Summary	91
<b>CHAPTER 5</b>	<b>CONCLUSION AND RECOMMENDATION</b>	<b>92</b>
5.1	Conclusion	92
5.2	Recommendation for Future Work	95
<b>CHAPTER 7</b>	<b>REFERENCES</b>	<b>96</b>
<b>APPENDICES</b>		<b>113</b>
APPENDIX A		113
APPENDIX B		115
APPENDIX C		120

**LIST OF TABLES**

<b>Table No.</b>	<b>Title</b>	<b>Page</b>
2.1	Applications of MSPM in industrial processes	13
2.2	Description of Extensions of PCA	16
2.3	Recent Applications of MLR for Data Prediction	29
3.1	Process Faults for Tennessee Eastman	45
3.2	Process List of Variables of TEP for monitoring	46
3.3	PID Trial Values for FIC, FcIC and TIC, FhIC	50
3.4	List of Variables for MIMO Plant	51
4.1	Results of FAR of TEP NOC Data	56
4.2	Fault Detection Times of Monitoring System	59
4.3	Analysis of Fastest Detection Between Conventional PCA and MLR – PCA Systems based on TEP Cases	60
4.4	Fault Identification Results of TEP Cases	62
4.5	Fault Detection for MIMO Plant	86
4.6	Fault Identification for MIMO Process Plant	87
A1	Data of Predictor and Criterion of TE Process	113
A2	Data of Predictor and Criterion of MIMO Pilot Plant	114
B1	Results of Fault Identification for TE Process	115
C1	Normal Operating Condition for MIMO Process Plant	120
C2	Faults Data for MIMO	123

## LIST OF FIGURES

Figure No.	Title	Page
2.1	Procedures of Process Monitoring Systems	10
2.2	The Graphical Interpretation of PCA for Data Compression	12
2.3	The Extension of PCA Model	15
2.4	The Illustrative Partitioning of Data matrix	26
3.1	Generic MSPM Framework based on PCA	33
3.2	Generic MLR-PCA based MSPM Framework	40
3.3	Flowsheet of the TEP System	44
3.4	Diagram for MIMO Plant System	48
4.1	Variance Explained vs Number of PCs Selected	55
4.2	NOC Data of XMEAS 1 (top), XMEAS 23 (middle) and XMEAS 4 (bottom)	64
4.3	Fault 1 Data of XMEAS 1 (top), XMEAS 23 (middle) and XMEAS 4 (bottom)	65
4.4	Progression of Fault Detection Results of $T^2$ for Fault 1	77
4.5	Progression of Fault Detection Results of SPE for Fault 1	78
4.6	Diagram A: Contribution plots of Conventional PCA for F1; Diagram B: Contribution plots of MLR-PCA for F1; Diagram C: Contribution plot of Predictor for F1	70
4.7	NOC data for XMV 10 (top) and XMEAS 9 (bottom)	71
4.8	Fault data of XMV 10 (top) and XMEAS 9 (bottom)	72
4.9	Progression Results of Fault 4 based on $T^2$	73
4.10	Progression Results of Fault 4 based on SPE	74



<b>Figure No.</b>	<b>Title</b>	<b>Page</b>
4.11	Diagram A: Contribution plots of Conventional PCA for F4; Diagram B: Contribution plots of MLR-PCA for F4; Diagram C: Contribution plot of Predictor for F4	76
4.12	NOC data for XMV 10 (top) and XMEAS 9 (bottom)	77
4.13	Fault data of XMV 10 (top) and XMEAS 9 (bottom)	78
4.14	Fault Detection for Fault 11 based on $T^2$	80
4.15	Fault Detection for Fault 11 based on SPE	81
4.16	Diagram A: Contribution plots of Conventional PCA for F11; Diagram B: Contribution plots of MLR-PCA for F11; Diagram C: Contribution plot of Predictor for F11	82
4.17	Accumulated Variances vs PCs of MIMO Plant	84
4.18	Progression of MV1 Th of NOC and Fault 1	88
4.19	Progression of Fault Detection of PCA and MLR-PCA.	89
4.20	Contribution Plot for Fault 1	91

**LIST OF ABBREVIATIONS**

<b>CMDS</b>	Classical Multidimensional Scaling
<b>CSTRwR</b>	Continuous Stirred Reactor with Recycle
<b>CVA</b>	Canonical Variate Analysis
<b>DKPCA</b>	Dynamic Kernel Principal Component Analysis
<b>DPCA</b>	Dynamic Principal Component Analysis
<b>FA</b>	Factor Analysis
<b>FAC</b>	False Alarm Cases
<b>FAR</b>	False Alarm Rate
<b>FDT</b>	Fault Detection Time
<b>FI</b>	Fault Identification
<b>FIC</b>	Flow rate Indicator Controller
<b>ICA</b>	Independent Component Analysis
<b>ICs</b>	Independent Components
<b>IDV</b>	Fault Cases
<b>HPLC</b>	High Performance Liquid Chromatography
<b>KPCA</b>	Kernel Principal Component Analysis
<b>LPLS</b>	Log-contrast Partial Least Square
<b>MBPCA</b>	Multi Block Principal Component Analysis
<b>MDS</b>	Multidimensional Scaling
<b>MEWMA</b>	Multivariate Exponentially - Weighted Moving Average
<b>MIMO</b>	Multiple Input Multiple Output
<b>MLR</b>	Multiple linear Regression
<b>MPLS</b>	Multi-way Partial Least Square
<b>MSE</b>	Mean Squared Error
<b>MSKPCA</b>	Multi Scale kernel Principal Component Analysis
<b>MSPC</b>	Multivariate Statistical Process Control
<b>MSPCA</b>	Multi Scale Principal Component Analysis
<b>MSPM</b>	Multivariate Statistical Process Monitoring

<b>NCD</b>	Number of Fault Cases Detected
<b>ND</b>	No detection
<b>NFD</b>	Number of Fastest Detection
<b>NIPAL</b>	Nonlinear Iterative Partial Least Square
<b>NLPCA</b>	Non-linear Principal Component Analysis
<b>NOC</b>	Normal Operating Condition
<b>PCA</b>	Principal Component Analysis
<b>PCR</b>	Principal Component Regression
<b>PCs</b>	Principal Components
<b>PLS</b>	Partial Least Square
<b>RPLS</b>	Recursive Partial Least Square
<b>SPC</b>	Statistical Process Control
<b>SPE</b>	Squared Prediction Error
<b>SPLS</b>	Sparse Partial Least Square
<b>SVD</b>	Singular Value Decomposition
<b>T<sup>2</sup></b>	Hotteling's
<b>TEP</b>	Tennessee Eastman Process
<b>TIC</b>	Temperature Indicator Controller
<b>XMEAS</b>	Process Measurements
<b>XMV</b>	Process Manipulated Variables

## LIST OF SYMBOLS

<b>A</b>	Dissimilarity matrix for MDS.
<b>A</b>	Unknown mixing matrix in ICA.
<i>a</i>	Number of PCs that retained in the PCA model.
<b>B</b>	Regression coefficients of PCR.
<b>B<sub>MDS</sub></b>	Double-centered matrix.
<b>b</b>	Regression coefficient for MLR.
<b>C</b>	Variance - covariance matrix.
<i>c<sub>m,m</sub></i>	Variance - covariance matrix at row <i>m</i> and column <i>m</i> .
<b>E</b>	Residual matrix in ICA.
$\tilde{e}_i$	The <i>i</i> <sup>th</sup> row vector in residual matrix.
$F_{A,n-A,\alpha}$	F distribution index with <i>A</i> and <i>n-A</i> degrees of freedom at $\alpha$ confident limit.
<b>I<sub>m</sub></b>	Identity matrix.
<i>l</i>	Number of selected principal components.
$lim_\alpha$	Control limits for $T^2$ and SPE statistics.
<i>m</i>	Number of samples for MDS.
<i>m</i>	Total number of variables.
$\bar{m}$	Means of the $T^2$ and SPE statistics.
<i>(MLR Statistics)<sub>j</sub></i>	MLR statistics (SPE for the criterion variables) at a particular sampling time 'j'.
<i>n</i>	Number of samples.
<b>P</b>	Loadings factor.
<b>P</b>	Score matrix for PCA.
<i>p</i>	Scores matrix for PLS.
$P_{i,j}$	The <i>i</i> <sup>th</sup> element for principal component <i>j</i> .
<b>P<sub>a</sub></b>	Score matrix that containing the first <i>a</i> score vectors.
<i>(PredictorStatistics)<sub>j</sub></i>	Predictor statistics at a particular sampling time 'j'.
<b>S</b>	Independent component matrix in ICA.
$SPE_\alpha$	$\alpha$ confidence limits for SPE parameter.

$SPE_i$	The contribution of the $i$ th variable to SPE.
$(SPE_i)_j$	Contribution of the $i$ th variable to MLR (SPE) statistics at particular sampling time ' $j$ '.
<b>T</b>	Orthogonal scores of the predictor variables.
$\mathbf{T}_{n \times l}$	Matrix of the multivariate scores.
<b>t</b>	Latent component matrix in PLS.
$\mathbf{T}_{(a)}$	Decomposing of data matrix into orthogonal scores.
$T_\alpha$	$A$ confidence limits for $T^2$ parameter.
$T^2$	T squared value
$T_i^2$	T squared value at sample $i$ .
<b>V</b>	Eigenvector matrix for PCA.
<b>V</b>	Loading matrix (a set of selected eigenvectors).
$\mathbf{V}_1$	Corresponding sets of eigenvectors in MDS.
$\mathbf{V}_a$	Loading matrix containing the first $a$ loading vectors.
$v$	Variances of the $T^2$ and SPE statistics.
$w$	Weight vectors.
$w_i$	Aggregating weight.
<b>X</b>	Predictor variables matrix.
<b>X</b>	Data matrix for ICA.
$\mathbf{X}_{n \times m}$	Data matrix where $n$ and $m$ are the number of samples and variables.
$\mathbf{X}_{\text{PCR}}$	Data scores obtained by the first principal components (PCs) in PCR.
$\mathbf{X}_{\text{PLS}}$	Predictor variables for PLS.
$\tilde{\mathbf{X}}$	Standardized multivariate data.
$\hat{\mathbf{X}}$	Data matrix produces by the decompression of the data scores.
$\tilde{\mathbf{X}}_{(a)}$	Variance of first principal components.
$x_{i,j}$	The ' $i$ 'th individual origin data of variable ' $j$ '.
$\tilde{x}$	Standardized data.
$\bar{x}_j$	Mean of variable ' $j$ '.
$(X_i)_j$	Contribution of the $i$ th variable to PCA statistics at particular sampling time ' $j$ '.
<b>Y</b>	Criterion variables.

$\mathbf{Y}_{MLR}$	New predicted criterion matrix.
$\mathbf{Z}$	Cartesian coordinate matrix.
$z_\alpha$	Standard normal deviate corresponding to the upper $(1-\alpha)$ percentile.
$\alpha$	Level of confidence limits.
$\sigma_j$	Standard variation of variable ' $j$ '.
$\lambda_j$	Eigenvalue corresponds to principal component $j$ .
$\Lambda$	Eigenvalues matrix for PCA.
$\Lambda$	Loading matrix for FA.
$\Lambda_1^{\frac{1}{2}}$	Diagonal matrix with all positive elements.
$\Psi$	Diagonal matrix for FA.
$\mathbf{1}_m$	Vector with element that equal to 1.

## **CHAPTER 1**

### **INTRODUCTION**

#### **1.1 INTRODUCTION**

Monitoring is a continuous real-time task of determining the possible conditions of a physical system, recognizing and indicating inconsistencies of the behavior (Isermann, 2011). The application of statistical method in monitoring is widely uses in chemical-based industries that involves with series of unit operations in order to convert the input materials into desired products following the qualitative and quantitative specifications of the customers. This can be considered as highly challenging as the process subjects to be affected by various unstable conditions over the time of operation. Simply avoiding or slow response in such situations may result with the decadence of product quality and even leads to catastrophic events as well as risking the profitability of the company. Thus, it has always been imperative to have a systematic mechanism which can routinely manage all of these abnormal situations automatically. These problems can be addressed quite effectively by using the process monitoring system. The method normally functioned to conduct fault detection, fault identification and fault diagnosis tasks.

## 1.2 RESEARCH BACKGROUND

In general, there are two popular types of process monitoring systems available for industrial application, which are univariate and multivariate monitoring systems. Since most modern industrial processes are involving multivariate in nature especially in measurements on a number of characteristics, instead one single characteristic, univariate method provides little information regarding the mutual interactions and also do not function well for multivariable processes with highly correlated variables (Nomikos and MacGregor, 1995 and Qin, 2012). The conventional univariate system, such as Statistical Process Control (SPC), has been mainly criticized for its limitation (particularly in the context of chemical-based operation), whereby it is only being operative under univariate analysis setting as well as a large number of control charts is always needed to be monitored concurrently (Bersismis et al., 2007). Besides, it always ignores the implications of harmonization between the output and input variables. Thus, the multivariate system such as Multivariate Statistical Process Monitoring (MSPM) can be regarded as the most practical method for handling complicated and large scale systems (Chiang et al., 2001).

As the process is in multivariate nature, the system will typically develop a model that correlates all of the variables simultaneously by using a set of normal operating condition (NOC) data that obtained from the historical process archive. In the other words, the system can utilize maximally all the process data stored for better use (Zhao et al., 2004). Besides, the system also has been popularly perceived as an advanced technique from the traditional SPC methodology. Unlike SPC, MSPM can extract useful information in terms of inter-related variable variations and represent it by using simplified parameters (normally in terms of multivariate scores and monitoring statistics). By applying the method, the monitoring operation can be executed much simpler (in the sense that only a small number of control charts are required) as well as critically consider the effect of all changes that contributed from various variables concurrently.



At the fundamental level application, there are two types of monitoring charts typically employed – Hotelling's  $T^2$  and Squared Prediction Errors (SPE). The first represents conceptually the magnitude of deviation of the current sample from the center, whereas, the second analyzes the consistency of the current sample correlation according to the NOC model development. Both have been used complementary, whereby, control limits for both statistics are also computed accordingly. The main task would be to observe the progressions of both statistics on a control chart (usually Shewhart-type control chart) that constructed respectively. When the process is normal, all the statistics will remain below the control limit lines. However, whenever a fault event takes place in the process, the corresponding statistics will move away from the normal region until a point where it goes beyond the control limits specified (sooner or later). If the abnormal trend persists over a period of time, the system will then initiate the alarm which signifies that one or a combination of faulty event(s) has (have) actually occurred in the process (fault detection). Contribution plot is then applied, purposely to identify the main possible variables that either contributing or being affected from that particular abnormal event that detected. Lastly, further investigation is needed to critically sort out and finally diagnose the true source of the problem that related to that particular abnormal existence.

### **1.3 PROBLEM STATEMENT**

MSPM normally utilizes linear-based principal component analysis (PCA) as the main technique of multivariate data compression. However, PCA sometimes is improperly used especially in modeling highly nonlinear processes as a high number of principal components (PCs) are always involved may lead to inefficient and unreliable monitoring performance reflected in false alarms and missed faults (Dong and McAvoy, 1996 and Žvokelj et al., 2011). If large variable are involved, then the PCs may also be selected considerably. As a result, Zhang et al., (1997) introduced non-linear PCA which based on the combination of neural network and principal curve but the computation is very demanding and it always requires a massive amount of data for creating the optimized NOC model (Yunus, 2012).

In other applications, Yunus and Zhang (2010a, b and c) as well as Yunus (2012) have developed three main frameworks of the MSPM by using the classical multidimensional scaling (CMDS) approach. Even though some improvements can be observed in terms of fault detection efficiency, those approaches employed different score projection (by way of variable scores) as opposed to the conventional PCA (by means of sample scores). It is argued that the CMDS approach cannot be effectively applied when it involves with a very large number of the variables as the variable scores cannot be reproduced as precisely as possible according to the pre-specified NOC configurations. This shows that, both techniques (PCA and MDS) suffered from technical difficulties particularly when handling large amount of variables in monitoring. Therefore, a new solution is demanded, whereby the number of original variables can be reduced to a certain extent (in terms of scales), while it can considerably maintain the original variation as largely as possible.

#### **1.4 RESEARCH OBJECTIVES**

One of the potential techniques available to address the issue is known as multiple linear regressions (MLR). The main objective of the technique is to predict a set of output values (criterion) based from a specified set of linear function, which consists of the predictor variables and to reduce the number of the dimensionality. Therefore, the main multivariate data will be divided into two groups, which are the criterion and predictor categories. By adopting such approach, the monitoring task can be made simpler and perhaps more effective, in the sense that only those criterion variables (predicted values) will be taken for monitoring, while preserving the rest of the predictor value trends in the form of linear functions. In light of this, the primary objectives of the study are:

- i. To apply a basic MSPM system using conventional PCA, whereby the monitoring outcome of this system will be used as the benchmark performance in order to assess the credibility of the proposed system.
- ii. To develop and investigate the performance of the MSPM using MLR method against the standard performance of MSPM-PCA. In particular, the monitoring data are divided into two main categories, whereby only the quality variables are used for monitoring.
- iii. To analyze the performance of MLR-PCA by applying real process instrument data.

## **1.5 RESEARCH SCOPES**

In order to accomplish the objectives, the scope of this study is focusing on the several criteria as follows:

- i. This study considers the challenging case studies by using the Tennessee Eastman Process (TEP) and also the data from the Multiple Input and Multiple Output (MIMO) Plant for monitoring.
- ii. The monitoring platform will be developed by using MATLAB 7.8.0:347 (R2009a) software.
- iii. In applying the MLR method, the main multivariate data will be divided into two groups, which are the criterion and predictor categories.
- iv. A number of comparative analyses between the proposed and conventional method are essential is evaluating the credibility of the new system's performance, particularly on the ground of number of cases detected, number of faster detection, fault detection and false alarm rate.
- v. Fault identification is also conducted to complement the fault detection results.

## 1.6 RATIONALE AND SIGNIFICANT

It is expected that the monitoring performance can be performed efficiently as well as effectively using a small number of dimensions in comparison to the conventional monitoring outcomes. This is simply because the modified data contains less number of variables (in term of the criterion value of the MLR model) in contra to the original magnitude. As a result, fewer PCs or dimensions are expected to be used in the PCA models respectively. At the same time, it is also assumed that the original variations of the predictor variables can be significantly preserved and re-produced by the criterion variables during monitoring operation.

As a whole, the main contributions of this study are:

- i. MLR technique can actually represent the behavior of the original data, in such a way it has the potential to be used in process monitoring.
- ii. Newly system of MSPM framework which integrates MLR and PCA techniques under a single application will create the process system less complicated.
- iii. The approach technique perhaps will improve the sensitivity of the fault detection performance comparative to the conventional approach particularly using less number of compressed PCs.
- iv. Process monitoring based MLR framework potentially reduce the complexity during the fault identification as well as diagnosis operations.

## 1.7 REPORT ORGANIZATION

This thesis has been divided into seven main chapters. Chapter 1 represents the introduction of the study, including research background, the objectives, goals and contributions, the remaining chapters in thesis are organized as follows;

Chapter 2 begins with some theory regarding the fundamental of MSPM, which emphasizing on the theoretical aspect of PCA. Then, other different multivariate methods are discussed, which consists of PLS, MDS, MLR, and others. The discussion is then followed by explaining the corresponding issues and extensions, which have been conducted in process monitoring.

While, Chapter 3 explains the methodology of developing the monitoring systems based on conventional and enhanced techniques. All the related procedures will be explained in details. In this chapter also includes the process description and the types of fault cases on the case studies that chosen to be investigated - Tennessee Eastman Process (TEP) and MIMO Plant.

Chapter 4 commences with the analysis of the data collection, analysis and the results of the modeling framework for both methods for TEP and MIMO Plant.

Chapter 5 concludes the findings in this work and suggests possible directions for future research. The main consideration of this chapter is to provide the corresponding to the research objectives in Chapter 1. Lastly, Chapter 6 lists all the references used in this research.

## **CHAPTER 2**

### **LITERATURE REVIEW**

#### **2.1 INTRODUCTION**

Maintaining product quality at the highest level is of wide spread concern in the process industries nowadays. This basically means that highly proportion of specifications product is always desirable, and eventually, it has sparked motivation for the reduction of abnormal variability in the normal operation routine to the lowest level as possible. This situation also has led to an increase in the use of process monitoring technique, which traditionally involving Statistical Process Control (SPC) (Papazoglou, 1998; Montgomery, 2009 and Kumar, 2013). However, the traditional SPC technique cannot be applied effectively in a multivariate nature of complex processes, which typically consist of huge matrix of several characteristics rather than a vector measurement. A number of limitations on the usage of SPC have been discussed comprehensively by several authors including (MacGregor and Kourti, 1995; Montgomery, 1996; Papazoglou, 1998; Jackson, 2005 and Behbahani et al., 2012).

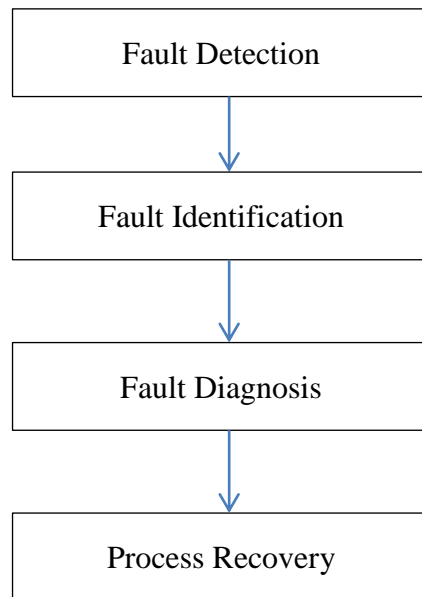
In addressing this particular issues, Multivariate Statistical Process Monitoring (MSPM), has been found the most suitable approach as well as progressively develop and widely applied for process monitoring presently. This could be perhaps due to its ability in providing systematic diagnostic tool for the extensive monitoring of a large operation data by means of off-line as well as on-line operation, which typically includes fault detection, identification, diagnosis and also process recovery. Besides, as the system depends heavily on the statistical approach, therefore, it is also applicable to wide variety of processes (batch and continuous) regardless of size and environment.

Thus, this chapter outlines the generic overview of the MSPM system, which consisting of four main sections. The first is on the introduction and subsequently followed by two main sections on developing the MSPM framework. In particular, the first describes a comprehensive overview on the fundamentals as well as the normal procedures of implementing MSPM system. Meanwhile, the second explains on the extensions and also the relative importance of the system with respect to the other industrial monitoring techniques. At the end of the chapter, a critical review is provided especially to highlight the significance of the study in contributing to the current MSPM scopes.

## 2.2 FUNDAMENTAL OF MSPM

According to Bersismis et al., (2005) Multivariate Statistical Process Monitoring (MSPM), which also known as Multivariate Statistical Process Control (MSPC), refers to a systematic statistical modeling procedure that consistently analyzes the performance of the process (continuous and batch) under investigation, particularly in deciding whether the current operating condition is actually normal or abnormal. This has been performed typically by means of observing the progression of monitoring statistics, namely  $T^2$  (magnitude of deviation) and SPE (consistency) on the control charts (Wise and Gallagher, 1996; Martin et al., 1996; Cinar et al., 2007 and MacGregor and Cinar, 2012). Chiang et al., (2001) stated that there are four basic steps typically involved in any of process monitoring procedures as depicted in **Figure 2.1**.

The first element is fault detection where the aim is to initiate the alarm whenever fault(s) has (have) been detected by the monitoring system (especially in the case of consistent violation on the control limits is evidenced). Next, fault identification is a set of specific procedures that highlights all the main variables which have shown great influence on the fault that detected in the first stage. In the third step, another vast set of procedures will be conducted specifically in diagnosing on those identified variables that captured previously, particularly in defining the true nature of the cause that contribute to the fault (fault diagnosis). The last element in the process monitoring system is process recovery whereby this process takes the remedial actions in eliminating the causes that lead to the occurrence of faults.



**Figure 2.1:** Procedures of Process Monitoring Systems (Chiang et al., 2001)



## **CHAPTER 1**

### **INTRODUCTION**

#### **1.1 INTRODUCTION**

Monitoring is a continuous real-time task of determining the possible conditions of a physical system, recognizing and indicating inconsistencies of the behavior (Isermann, 2011). The application of statistical method in monitoring is widely uses in chemical-based industries that involves with series of unit operations in order to convert the input materials into desired products following the qualitative and quantitative specifications of the customers. This can be considered as highly challenging as the process subjects to be affected by various unstable conditions over the time of operation. Simply avoiding or slow response in such situations may result with the decadence of product quality and even leads to catastrophic events as well as risking the profitability of the company. Thus, it has always been imperative to have a systematic mechanism which can routinely manage all of these abnormal situations automatically. These problems can be addressed quite effectively by using the process monitoring system. The method normally functioned to conduct fault detection, fault identification and fault diagnosis tasks.

## 1.2 RESEARCH BACKGROUND

In general, there are two popular types of process monitoring systems available for industrial application, which are univariate and multivariate monitoring systems. Since most modern industrial processes are involving multivariate in nature especially in measurements on a number of characteristics, instead one single characteristic, univariate method provides little information regarding the mutual interactions and also do not function well for multivariable processes with highly correlated variables (Nomikos and MacGregor, 1995 and Qin, 2012). The conventional univariate system, such as Statistical Process Control (SPC), has been mainly criticized for its limitation (particularly in the context of chemical-based operation), whereby it is only being operative under univariate analysis setting as well as a large number of control charts is always needed to be monitored concurrently (Bersismis et al., 2007). Besides, it always ignores the implications of harmonization between the output and input variables. Thus, the multivariate system such as Multivariate Statistical Process Monitoring (MSPM) can be regarded as the most practical method for handling complicated and large scale systems (Chiang et al., 2001).

As the process is in multivariate nature, the system will typically develop a model that correlates all of the variables simultaneously by using a set of normal operating condition (NOC) data that obtained from the historical process archive. In the other words, the system can utilize maximally all the process data stored for better use (Zhao et al., 2004). Besides, the system also has been popularly perceived as an advanced technique from the traditional SPC methodology. Unlike SPC, MSPM can extract useful information in terms of inter-related variable variations and represent it by using simplified parameters (normally in terms of multivariate scores and monitoring statistics). By applying the method, the monitoring operation can be executed much simpler (in the sense that only a small number of control charts are required) as well as critically consider the effect of all changes that contributed from various variables concurrently.

At the fundamental level application, there are two types of monitoring charts typically employed – Hotelling's  $T^2$  and Squared Prediction Errors (SPE). The first represents conceptually the magnitude of deviation of the current sample from the center, whereas, the second analyzes the consistency of the current sample correlation according to the NOC model development. Both have been used complementary, whereby, control limits for both statistics are also computed accordingly. The main task would be to observe the progressions of both statistics on a control chart (usually Shewhart-type control chart) that constructed respectively. When the process is normal, all the statistics will remain below the control limit lines. However, whenever a fault event takes place in the process, the corresponding statistics will move away from the normal region until a point where it goes beyond the control limits specified (sooner or later). If the abnormal trend persists over a period of time, the system will then initiate the alarm which signifies that one or a combination of faulty event(s) has (have) actually occurred in the process (fault detection). Contribution plot is then applied, purposely to identify the main possible variables that either contributing or being affected from that particular abnormal event that detected. Lastly, further investigation is needed to critically sort out and finally diagnose the true source of the problem that related to that particular abnormal existence.

### **1.3 PROBLEM STATEMENT**

MSPM normally utilizes linear-based principal component analysis (PCA) as the main technique of multivariate data compression. However, PCA sometimes is improperly used especially in modeling highly nonlinear processes as a high number of principal components (PCs) are always involved may lead to inefficient and unreliable monitoring performance reflected in false alarms and missed faults (Dong and McAvoy, 1996 and Žvokelj et al., 2011). If large variable are involved, then the PCs may also be selected considerably. As a result, Zhang et al., (1997) introduced non-linear PCA which based on the combination of neural network and principal curve but the computation is very demanding and it always requires a massive amount of data for creating the optimized NOC model (Yunus, 2012).

## **CHAPTER 3**

### **METHODOLOGY**

#### **3.1 INTRODUCTION**

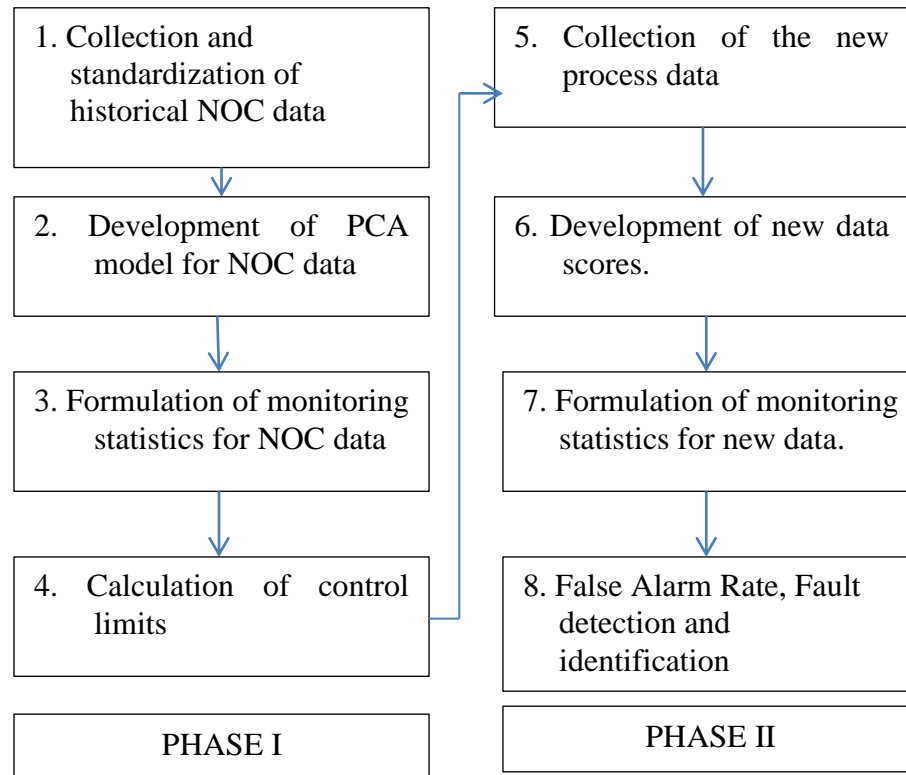
This chapter explains the methodology of the conventional as well as the proposed method. The former utilizes Multivariate Statistical Process Monitoring (MSPM) System, which basically developed based on the procedures of Principal Component Analysis (PCA) technique. Meanwhile, the structure of the proposed system applies the integration of Multiple Linear Regression (MLR) into the original MSPM framework and that to be used for data prediction. This newly introduced monitoring system is still depending on the PCA to mainly compress the multivariate data.

In general, this chapter is divided into four main sections starting briefly with the introduction. Next, the PCA-MSPM procedures (Framework I) is explained and subsequently followed by the new procedures of MLR-PCA-MSPM (Framework II). All the model development for PCA and MLR-PCA are from the MATLAB toolbox. Then, the description of the Tennessee Eastman Process (TEP) case study which has been utilized to demonstrate the capability of the monitoring system that implemented in this study which the data was generated from a simulation work due to Chiang et al., (2001).

This chapter also presents description of the real case study, Multiple Input Multiple Output (MIMO) training system pilot plant which has been setup in Universiti Malaysia Pahang, Malaysia. As to ensure the true strength of the proposed monitoring system, MIMO plant was chosen, which conceptually represents the real plant operation data that originally produced through this project. The process contains several of fault cases, which are found suitable to be applied in this research.

### **3.2 FRAMEWORK I: PCA-MSPM SYSTEM PROCEDURES**

The complete original procedures of the conventional MSPM framework can be obtained from Macgregor and Kourti (1995) as well as Raich and Cinar (1996), whereby it can be separated into two main phases as shown in **Figure 3.1** (Yunus, 2012). From **Figure 3.1**, the first phase is related to the model development of NOC data whereas the second facilities for monitoring of the new process data.



**Figure 3.1:** Generic MSPM Framework based on PCA.

### 3.2.1 PHASE I PROCEDURES

The first step of the first phase basically involves with collection of NOC data  $\mathbf{X}_{m \times n}$  ( $m$ : variables,  $n$ : samples), which normally conducted off-line based on the historical process data archive. The data are then standardized to zero mean and unit variance by using equation (3.1) until equation (3.3). All the procedures are due to MacGregor and Kourti, (1995); Nomikos and MacGregor, (1995) and Jackson, (2005). PCA typically utilizes variance-covariance or a matrix correlation measure of the normal operating condition (NOC) data matrix as the basis in developing the compressed multivariate configuration. According to Chiang et al., (2001), data standardization relates to capturing the data variation that extracted from the mean of data variables and scales it to unit variance. The data variables can be standardized by applying equation (3.1):