

**COMPARING THREE METHODS OF HANDLING MULTICOLLINEARITY  
USING SIMULATION APPROACH**

**NORLIZA BINTI ADNAN**

**UNIVERSITI TEKNOLOGI MALAYSIA**

COMPARING THREE METHODS OF HANDLING MULTICOLLINEARITY  
USING SIMULATION APPROACH

NORLIZA BINTI ADNAN

A thesis submitted in fulfilment of the  
requirements for the award of the degree of  
Master of Science (Mathematics)

Faculty of Science  
Universiti Teknologi Malaysia

MEI 2006

*To*

*My loving and supportive parents,*

*Hj Adnan and Pn Ramnah*

*My siblings,*

*Hairil, Aizam, Haffiz and Sarah*

*and*

*Pn Rusnah and family, and all my supportive friends*

## ACKNOWLEDGEMENT

*“In the name of ALLAH the All-Merciful, The All-Compassionate. All praise be to ALLAH for giving me the strength and courage to complete this thesis”*

A very great gratitude and appreciation expressed to all those who make a part to the successful cease of this thesis either directly or indirectly. Particularly, I wish to express my sincere appreciation to my main thesis supervisor, Dr. Maizah Hura binti Ahmad, for encouragement, guidance, critics and motivation. I am also very thankful to my co-supervisor Dr. Robiah binti Adnan for her guidance, advices and motivation. Without their continued support and interest, this thesis would not have been the same as presented here. An appreciation also goes to the Universiti Teknologi Malaysia for financial support.

Lots of gratitude and special thanks also to my family for their support, love and encouragement throughout my study. My sincere appreciation also extends to all my lovely friends.

## ABSTRACT

In regression, the objective is to explain the variation in one or more response variables, by associating this variation with proportional variation in one or more explanatory variables. A frequent obstacle is that several of the explanatory variables will vary in rather similar ways. This phenomenon called multicollinearity, is a common problem in regression analysis. Handling multicollinearity problem in regression analysis is important because least squares estimations assume that predictor variables are not correlated with each other. The performances of ridge regression (RR), principal component regression (PCR) and partial least squares regression (PLSR) in handling multicollinearity problem in simulated data sets are compared to help and give future researchers a comprehensive view about the best procedure to handle multicollinearity problems. PCR is a combination of principal component analysis (PCA) and ordinary least squares regression (OLS) while PLSR is an approach similar to PCR because a component that can be used to reduce the number of variables need to be constructed. RR on the other hand is the modified least square method that allows a biased but more precise estimator. The algorithm is described and for the purpose of comparing the three methods, simulated data sets where the number of cases was less than the number of observations were used. The goal was to develop a linear equation that relates all the predictor variables to a response variable. For comparison purposes, mean square errors (MSE) were calculated. A Monte Carlo simulation study was used to evaluate the effectiveness of these three procedure. The analysis including all simulations and calculations were done using statistical package S-Plus 2000 software.

## ABSTRAK

Objektif bagi regresi ialah untuk menerangkan variasi bagi satu atau lebih pembolehubah bersandar dengan cara menghubungkan variasi ini berkadar dengan satu atau lebih pembolehubah tak bersandar. Halangan yang sering berlaku ialah apabila wujudnya kebersandaran antara pembolehubah-pembolehubah tak bersandar. Fenomena ini dipanggil multikolinearan. Mengawal dan mengatasi masalah multikolinearan di dalam analisis regresi adalah penting kerana kaedah penganggaran kuasa dua terkecil menganggap bahawa pembolehubah tak bersandar tidak berkorelasi antara satu sama lain. Perbandingan antara penggunaan kaedah regresi permatang (RR), regresi komponen berkepentingan (PCR) dan regresi sebahagian kuasa dua terkecil (PLSR) di dalam mengawal masalah multikolinearan dilakukan menggunakan set data yang disimulasi bagi membantu dan memberi satu pendekatan kepada para pengkaji yang akan datang tentang pemilihan kaedah terbaik bagi mengawal masalah multikolinearan. Kaedah regresi permatang ialah pengubahsuaian kaedah kuasa dua terkecil (LS) yang memasukkan pemalar kepincangan,  $\delta$  di dalam penganggar kuasa dua terkecil. Regresi komponen berkepentingan pula merupakan gabungan analisis komponen berkepentingan (PCA) dengan kaedah kuasa dua terkecil biasa (OLS) sementara kaedah regresi sebahagian kuasa dua terkecil adalah hampir sama dengan kaedah regresi komponen berkepentingan di mana komponen baru perlu dibina untuk mengurangkan bilangan pembolehubah. Algoritma bagi setiap kaedah turut diterangkan dan untuk tujuan perbandingan bagi setiap kaedah, set data bagi kes bilangan pembolehubah tak bersandar lebih kecil dari bilangan pemerhatian. Perbandingan keberkesanan bagi ketiga-tiga kaedah tersebut menggunakan ralat min kuasa dua (MSE). Kaedah simulasi Monte Carlo digunakan untuk menilai keberkesanan ketiga-tiga kaedah yang dibincangkan. Semua simulasi dan pengiraan dilakukan dengan menggunakan pakej statistik S-PLUS 2000.

## TABLE OF CONTENTS

CHAPTER	SUBJECT	PAGE
	<b>COVER</b>	i
	<b>DECLARATION</b>	ii
	<b>DEDICATION</b>	iii
	<b>ACKNOWLEDGEMENT</b>	iv
	<b>ABSTRACT</b>	v
	<b>ABSTRAK</b>	vi
	<b>TABLE OF CONTENTS</b>	vii
	<b>LIST OF TABLES</b>	x
	<b>LIST OF FIGURES</b>	xiv
	<b>LIST OF SYMBOLS</b>	xvii
	<b>LIST OF APPENDICES</b>	xix
<b>1</b>	<b>INTRODUCTION</b>	
	1.1 Background	1
	1.2 The Problem of Multicollinearity	2
	1.3 Statement of Problem	5
	1.4 Research Objectives	5
	1.5 Scope of Research	5
	1.6 Summary and Outline of Research	6

## 2 LITERATURE REVIEW

2.1	Introduction	7
2.2	Ordinary Least Squares Regression	7
2.3	Multicollinearity Problem in Regression Analysis	10
2.3.1	Explanation of Multicollinearity	10
2.3.2	Effects of Multicollinearity in Least Squares Regression	13
2.3.3	Multicollinearity Diagnostics	23
	2.3.3.1 Informal Diagnostics	24
	2.3.3.2 Formal Methods	25
2.4	Concluding Remarks	32

## 3 METHODS FOR HANDLING MULTICOLLINEARITY

3.1	Introduction	34
3.2	Partial Least Squares Regression	35
3.2.1	The construction of $k$ components	37
3.2.2	Regress the response into $k$ components	40
3.3	Principal Components Regression	49
3.3.1	The construction of $k$ components	50
3.3.2	Regress the response into $k$ components	53
3.3.3	Bias in Principal Components Coefficient	54
3.4	Ridge Regression	61



<b>4</b>	<b>SIMULATION AND ANALYSIS</b>	
4.1	Introduction	71
4.2	Generating Simulated Data Sets	72
4.3	Performance Measures	76
4.4	Simulation Results	78
4.4.1	Partial Least Squares Regression	82
4.4.2	Principal Component Regression	104
4.4.3	Ridge Regression	116
<b>5</b>	<b>COMPARISONS ANALYSIS AND DISCUSSIONS</b>	
5.1	Introduction	125
5.2	Performance on Classical Data	125
5.3	Performance on Simulated Data Sets	127
5.4	Comparison Analysis	134
5.5	Discussions	142
<b>6</b>	<b>SUMMARY, CONCLUSIONS AND FUTURE RESEARCH</b>	
6.1	Introduction	145
6.2	Summary	145
6.3	Significant Findings and Conclusions	147
6.4	Future Research	149

**REFERENCES**

**APPENDICES**

## LIST OF TABLES

TABLE NO.	TITLE	PAGE
3.1	Tobacco Data	45
3.2	Variance Inflation Factors for Tobacco Data	46
3.3	PLS weights vectors, $r_k$ for the PLS Components	46
3.4	Loadings for the PLS Components	47
3.5	PLS Components	47
3.6	RMSE values for $k$ components	48
3.7	MSE values for $k$ components	49
3.8	The correlation matrix	58
3.9	The eigenvalues of the correlation matrix	58
3.10	Matrix of eigenvectors	58
3.11	The principal components	59
3.12	Values of $C_\delta$ and the regression coefficients for various values of $\delta$	69
4.1	Factors and levels for the simulated data sets	72
4.2	The specific values for $x_{ip}$ for each sets of $p$ regressors	74
4.3	The response variable, $y$ for each sets of $p$ regressors	74
4.4	The VIF's values for the choice of the variance, $\Sigma$ for noise matrix ( $\Delta$ )	75
4.5	The VIF's values for each sets of $p$ regressors	76
4.6	The values of correlation for $p = 2$ regressors	79
4.7	The values of correlation for $p = 4$ regressors	79
4.8	The values of correlation for $p = 6$ regressors	79

4.9	Regression model for $p = 2$ regressors for $n = 100$	80
4.10	Regression model for $p = 4$ regressors for $n = 100$	80
4.11	Regression model for $p = 6$ regressors for $n = 100$	81
4.12	PLS weights, $\mathbf{r}_i$ for the PLS Components for $p = 2$ regressors	83
4.13	PLS weights, $\mathbf{r}_i$ for the PLS Components for $p = 4$ regressors	83
4.14	PLS weights, $\mathbf{r}_i$ for the PLS Components for $p = 6$ regressors	83
4.15	PLS weights, $\mathbf{r}_i$ for the PLS Components for $p = 50$ regressors	84
4.16	PLS loadings, $\mathbf{p}_i$ for the PLS Components for $p = 2$ regressors	85
4.17	PLS loadings, $\mathbf{p}_i$ for the PLS Components for $p = 4$ regressors	85
4.18	PLS loadings, $\mathbf{p}_i$ for the PLS Components for $p = 6$ regressors	85
4.19	PLS loadings, $\mathbf{p}_i$ for the PLS Components for $p = 50$ regressors	86
4.20	Correlations between each PLS components and $y$	87
4.21	Correlation between each components	89
4.22	RMSE values for $p = 2$ data sets	97
4.23	RMSE values for $p = 4$ data sets	97
4.24	RMSE values for $p = 6$ data sets for first five PLS components	97
4.25	RMSE values for $p = 50$ data sets for first five PLS components	97
4.26	Regression model using all the PLS Scores for $p = 2$ and $n = 100$	99
4.27	Regression model using all the PLS Scores for $p = 4$ and $n = 100$	100
4.28	Regression model using selected PLS Scores for $p = 4$ ( $k_{opt} = 1$ )	101
4.29	Regression model using selected PLS Scores ( $k_{opt} = 5$ ) for $p = 6$ and $n = 100$	102
4.30	The eigenvalues of the correlation matrix	105
4.31	Matrix of eigenvectors	106
4.32	Percentage of variance explained and the eigen values of the $p = 2$	107

4.33	Percentage of variance explained and the eigen values of the $p = 4$	108
4.34	Percentage of variance explained and the eigen values of the $p = 6$	108
4.35	Percentage of variance explained and the eigen values of the $p = 50$	109
4.36	Regression model using all the PC scores for $p = 2$ and $n = 100$	111
4.37	Regression model using selected PC scores ( $PC_1$ ) for $p = 2$ and $n = 100$	111
4.38	Regression model using all the PC scores for $p = 4$ and $n = 100$	112
4.39	Regression model using selected PC scores ( $PC_1$ ) for $p = 4$ and $n = 100$	113
4.40	Regression model using all the PC scores for $p = 6$ and $n = 100$	114
4.41	Regression model using selected PC scores ( $PC_1$ ) for $p = 6$ and $n = 100$	114
4.42	Values of $\delta$ , $C_\delta$ and coefficient vectors employed for $p = 2$ and $n = 100$	119
4.43	Values of $\delta$ , $C_\delta$ and coefficient vectors employed for $p = 4$ and $n = 100$	120
4.44	Values of $\delta$ , $C_\delta$ and coefficient vectors employed for $p = 6$ and $n = 100$	121
5.1	Performance of PLS, PC and RR methods on classical data sets	126
5.2	Cross-validation of PLS, PC and RR methods, $R^2$ for $p = 2$	127
5.3	Cross-validation of PLS, PC and RR methods, $R^2$ for $p = 4$	130
5.4	Cross-validation of PLS, PC and RR methods, $R^2$ for $p = 6$	130
5.5	Cross-validation of PLS, PC and RR methods, $R^2$ for $p = 50$	130

5.6	MSE values for $p = 2$ and specified $n = 20, 30, 40, 60, 80$ and 100	131
5.7	MSE values for $p = 4$ and specified $n = 20, 30, 40, 60, 80$ and 100	131
5.8	MSE values for $p = 6$ and specified $n = 20, 30, 40, 60, 80$ and 100	135
5.9	MSE values for $p = 50$ and specified $n = 60, 80$ and 100	135
5.10	Summary of the performances of the three methods with $p = 2$	142
5.11	Summary of the performances of the three methods with $p = 4$	142
5.12	Summary of the performances of the three methods with $p = 6$	143
5.13	Summary of the performances of the three methods with $p = 50$	143

## LIST OF FIGURES

FIGURE NO.	TITLE	PAGE
1.1	Multicollinearity in simple linear regression	3
1.2	Multicollinearity in multiple linear regression	4
2.1	Picket Fences illustrations	15
2.2	The choice of VIF value against the R-square value.	29
3.1	Steps in SIMPLS algorithm	42
3.2	Steps in Principal Component Regression algorithm	56
3.3	The sampling distribution of biased and unbiased estimator	61
3.4	Steps in Ridge Regression algorithm	67
3.5	Plot of $C_\delta$ against $\delta$	70
4.1	Flowchart summarizing performance assessment of methodology	77
4.2	Correlation between $x_1$ and $x_2$ for $p = 2$ regressors	88
4.3	Correlation between first and second components for $p = 2$ regressors	89
4.4	Correlation between each components for $p = 4$ regressors	90
4.5	Correlation between each components for $p = 6$ regressors	92
4.6	X- and Y-scores for $p = 2$ regressors (First and second components)	94
4.7	X- and Y-scores for $p = 4$ regressors (All components)	94
4.8	Correlation between first five components of $p = 6$ regressors data set and response variable, $y$	95
4.9	Plot of $\delta$ against $C_\delta$ for $p = 2$ and $n = 100$	117

4.10	Plot of $\delta$ against $C_\delta$ for $p = 4$ and $n = 100$	117
4.11	Plot of $\delta$ against $C_\delta$ for $p = 6$ and $n = 100$	118
4.12	Plot of $\delta$ against $C_\delta$ for $p = 50$ and $n = 100$	118
5.1	Plot of regression coefficients against number of regressors	126
5.2	Plot of $R^2$ against $m = 100$ replications for $p = 2$ for specified $n = 20$	128
5.3	Plot of $R^2$ against $m = 100$ replications for $p = 2$ for specified $n = 100$	128
5.4	Plot of $R^2$ against $m = 100$ replications for $p = 4$ for specified $n = 20$	131
5.5	Plot of $R^2$ against $m = 100$ replications for $p = 4$ for specified $n = 100$	131
5.6	Plot of $R^2$ against $m = 100$ replications for $p = 6$ for specified $n = 20$	132
5.7	Plot of $R^2$ against $m = 100$ replications for $p = 6$ for specified $n = 100$	132
5.8	Plot of $R^2$ against $m = 100$ replications for $p = 50$ for specified $n = 60$	133
5.9	Plot of $R^2$ against $m = 100$ replications for $p = 50$ for specified $n = 100$	133
5.10	Plot of MSE values against $m = 100$ replications for $p = 2$ for $n = 20$	138
5.11	Plot of MSE values against $m = 100$ replications for $p = 2$ for $n = 100$	138
5.12	Plot of MSE values against $m = 100$ replications for $p = 4$ for $n = 20$	139
5.13	Plot of MSE values against $m = 100$ replications for $p = 4$ for $n = 100$	139
5.14	Plot of MSE values against $m = 100$ replications for $p = 6$ for $n = 20$	140

5.15	Plot of MSE values against $m = 100$ replications for $p = 6$ for $n = 100$	140
5.16	Plot of MSE values against $m = 100$ replications for $p = 50$ for $n = 60$	141
5.17	Plot of MSE values against $m = 100$ replications for $p = 50$ for $n = 100$	141



## LIST OF SYMBOLS

$y$	Response (dependent) variable
$x$	Predictor (independent) variable
$\beta$	Parameter (regression coefficient), known constant
$\varepsilon$	Error
$\sigma^2$	Variance
$Y$	Matrix of observations
$X$	Matrix of predictors
$\beta$	Vector of parameters
$\varepsilon$	Vector of error matrix term
$I$	Identity matrix
$\hat{\beta}$	Estimate regression coefficient
$\hat{\beta}$	Vector of estimate regression coefficient
$\hat{y}$	Fitted value
$H$	Hat matrix
$h_{ii}$	$i$ th diagonal element of hat matrix
$e$	Residual term
$\mathbf{e}$	Vector of residual term
$n$	Number of observations
$p$	Number of regressors
$\delta$	Shrinkage parameter
$V$	Matrix of normalized eigenvectors $X'X$
$A$	The diagonal matrix of eigenvalues of $X'X$
$T$	Components for Partial Least Squares Regression

$P$	Matrix of $x$ -loadings
$\mathbf{r}$	PLS weight vectors
$k$	Number of components for PLS and PCR
$R^2$	Coefficient of Determination
$Z$	Components for Principal Component Regression
$\lambda$	Eigenvalue

**ABBREVIATIONS****MEANING**

<i>CI</i>	Condition Indices
cov	Covariance
GCV	Generalized Cross Validation
LS	Least Squares
max	Maximum
MSE	Mean Square Error
OLS	Ordinary Least Squares
PCR	Principal Component Regression
PLSR	Partial Least Squares Regression
RMSE	Root Mean Square Error
RR	Ridge Regression
SMC	Squared Multiple Correlation
VIF	Variance Inflation Factors

**LIST OF APPENDICES**

<b>APPENDIX</b>	<b>TITLE</b>	<b>PAGE</b>
A	S-PLUS codes for data generating function	156
B	S-PLUS codes and simulation results for Partial Least Squares Regression method in Chapter IV	162
C	S-PLUS codes and simulation results for Principal Component Regression method in Chapter IV	182
D	S-PLUS codes and simulation results for Ridge Regression method in Chapter IV	200
E	S-PLUS code for classical data set	215
F	Simulation Results For Chapter V	225

## CHAPTER 1

### INTRODUCTION

#### 1.1 Background

Regression analysis is one of the most widely use of all statistical tools that utilizes the relation between two or more quantitative variables so that one variable can be predicted from the others. The relationship of each predictor to the criterion is measured by the slope of the regression line of the criterion  $Y$  on the predictor. The regression coefficients are the values of these slopes.

The multiple linear regression model relates  $Y$  to  $X_1, X_2, \dots, X_p$  and can be expressed in terms of matrices as  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  where  $\mathbf{y}$  is the  $n \times 1$  vector of observed response values,  $\mathbf{X}$  is the  $n \times p$  matrix of  $p$  regressors,  $\boldsymbol{\beta}$  is the  $p \times 1$  regression coefficients and  $\boldsymbol{\varepsilon}$  is the  $n \times 1$  vector of error terms. The objectives of regression analysis are, (1) to find the estimates of unknown parameters  $\boldsymbol{\beta}$ 's and test of  $\beta_1, \beta_2, \dots, \beta_k$  for the significance of the associated predictors, (2) to use the regression equation to estimate  $Y$  from  $X_1, X_2, \dots, X_p$  and, (3) to measure the error involved in the estimation. The multiple

linear regression model may be used to identify important regressor variables that can be used to predict future values of the response variable.

The method of least squares is used to find the best line that on the average, is closest to all points. In other words, to find the the best estimates of  $\beta$ 's with the least squares criterion which minimizes the sum of squared distances of all points from the actual observation to the regression surface. The name least squares comes from minimizing the squared residuals. From the Gauss-Markov theorem, least squares is always the best linear unbiased estimator (BLUE) and if  $\varepsilon$  is assumed to be normally distributed with mean 0 and variance  $\sigma^2$ , then least squares is the uniformly minimum variance unbiased estimator.

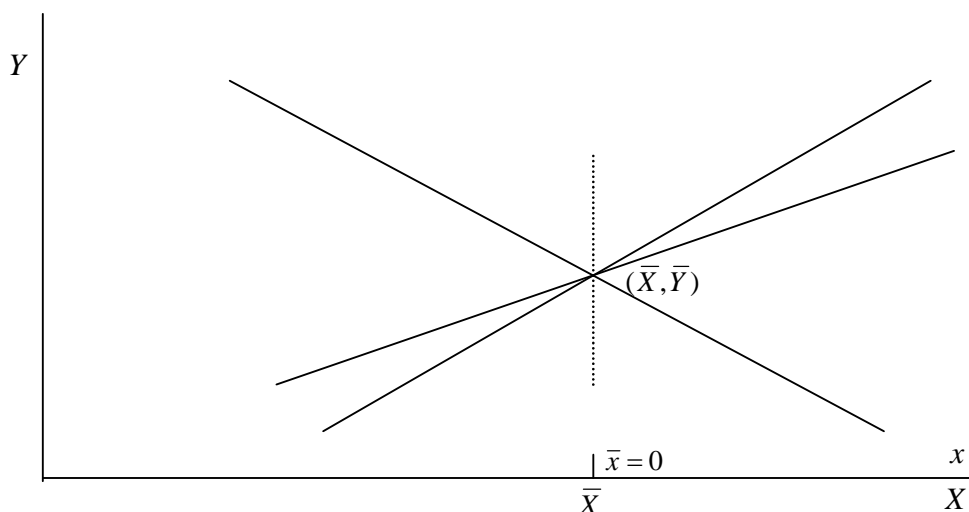
## **1.2 The Problem of Multicollinearity**

In the applications of regression analysis, multicollinearity is a problem that always occurs when two or more predictor variables are correlated with each other. This problem can cause the value of the least squares estimated regression coefficients to be conditional upon the correlated predictor variables in the model. As defined by Bowerman and O'Connell (1990), multicollinearity is a problem in regression analysis when the independent variables in a regression model are intercorrelated or are dependent on each other.

There are a variety of informal and formal methods that have been developed for detecting the presence of serious multicollinearity. One of the most commonly used is the variance inflation factor (VIF) that measures how much the variances of the estimated regression coefficients are inflated compared to when the independent variables are not linearly related (Neter, et. al., 1990). The problem of multicollinearity can be remedied

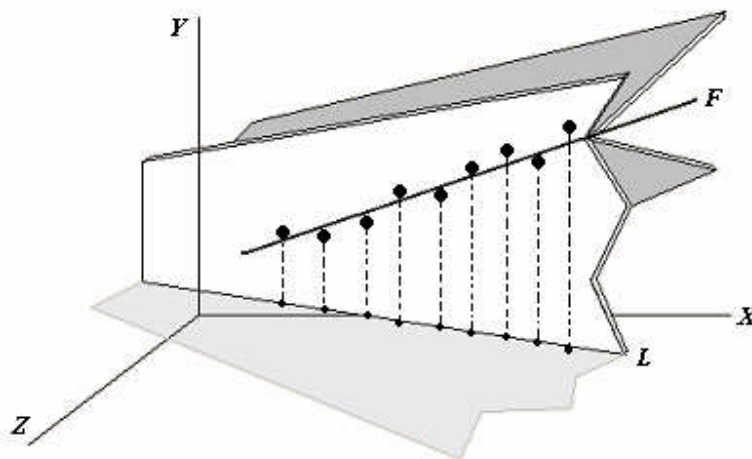
using some method of estimation or some modifications of the method of least squares for estimating the regression coefficients.

The problem of multicollinearity can occur in both simple linear regression and multiple linear regression. Figure 1.1 illustrates the problem of multicollinearity that occur in simple regression (Wannacott and Wannacott, 1981). The figure shows how the estimate  $\hat{\beta}$  becomes unreliable if the  $X_i$ 's were closely bunched, that is, if the regressor  $X$  had little variation. When the  $X_i$ 's are concentrated on one single value  $\bar{X}$ , then  $\hat{\beta}$  is not determined at all. For each line, the sum of squared deviations is the same, since the deviations are measured vertically from  $(\bar{X}, \bar{Y})$ . If  $X_i = \bar{X}$ , then all  $x_i = 0$ , and the term involving  $\hat{\beta}$  is zero. Hence, the sum of squares does not depend on  $\hat{\beta}$  at all. Therefore, when the values of  $X$  show little or no variation, then the effect of  $X$  on  $Y$  can no longer be sensibly investigated. The best fit for  $Y$  for data with multicollinearity was not a line, but rather a point  $(\bar{X}, \bar{Y})$ . In explaining  $Y$ , multicollinearity makes the  $X_i$ 's lose one dimension.



**Figure 1.1** : Multicollinearity in simple regression

Figure 1.2 illustrates the problem of multicollinearity in multiple regression (Wannacott and Wannacott, 1981). All the observed points in the scatter plot lie in the vertical plane running up through  $L$ . In explaining  $Y$ , multicollinearity makes the  $X_i$ 's lose one dimension and the best fit for  $Y$  is not a plane but instead the line  $F$ .



**Figure 1.2** : Multicollinearity in multiple regression

Several approaches for handling multicollinearity problem have been developed such as Principal Component Regression, Partial Least Squares Regression and Ridge Regression. Principal Components Regression (PCR) is a combination of principal component regression analysis (PCA) and ordinary least squares regression (OLS). Partial Least Squares Regression (PLSR) is an approach similar to PCR because one needs to construct a component that can be used to reduce the number of variables. Ridge Regression is the modified least squares method that allow biased estimators of the regression coefficient.



### **1.3 Statement of Problem**

This study will explore the following question :

Which method among Principal Component Regression, Partial Least Squares Regression and Ridge Regression performs best as a method for handling multicollinearity problem in regression analysis ?

### **1.4 Research Objectives**

1. To compare the use of Partial Least Squares Regression, Principal Component Regression and Ridge Regression for handling multicollinearity problem.
2. To study the degree of efficiency among the three methods and hence rank them in terms of their capabilities in overcoming the multicollinearity problem using simulated data sets.

### **1.5 Scope of Research**

For this research, the problem is focused on the analysis of multicollinearity problems in regression analysis using simulated data sets.

In Principal Component Regression, features from Principal Component Analysis (PCA) and multiple regression are combined to handle multicollinearity. The principal components computed a linear combinations of the explanatory variables (regressors). The Partial Least Squares Regression method is similar to Principal Component Regression because it is also a two-step procedure. However, PLSR searches for a set of

components called latent vectors that performs a simultaneous decomposition of  $X$  and  $Y$  with the constraint that those component explain as much as possible of the covariance between  $X$  and  $Y$ . This step generalizes PCA and it is followed by a regression step where the decomposition of  $X$  is used to predict  $Y$ .

The other method considered is the Ridge Regression. It is performed by adding a small biased estimator to the elements on the diagonal of the matrix  $(X'X)$  to be inverted, that is modifications of the least squares estimator.

## **1.6 Summary and Outline of Research**

The goal of this research is to find the best procedure and method to handle multicollinearity problems by comparing the performances of the three methods to determine which method is superior than the others in terms of its practicality and efficiency. Practicality means how effective or convenient a method is in actual use while efficiency means how well the methods work in producing the best regression model and measured by a specified test discussed in Chapter 4. Basically three different methods are put forward, two of which are the methods with two-step procedures where it computes a component(s) and then regressed it on the response variable. The third method is a modification of the least squares method that allows biased estimators of the regression coefficients. The algorithms for each method used in this study are shown in Chapter 3.

Chapter 2 reviews the relevent literature on published work done recently concerning the problems of multicollinearity. Discussion on methods for handling multicollinearity problems in regression analysis is presented in Chapter 3. Chapter 4 describes the simulation work and the analysis of the three methods. Chapter 5 discusses the performances of the three methods and makes comparisons among them. Lastly, Chapter 6 concludes the study and makes recommendations for further research.

variable. The advantages of PCR are that hypothesis testing can be performed, and that complete optimisation is used in determining the PCs.

#### **6.4 Future Research**

The performance of the three methods can also be done by comparing the use of all methods for high-dimensional regressors where  $p > n$ . It is known that the problem of multicollinearity is present in the data set where the number of variables is high compared to the number of observations.

These three methods can also be considered in the handling of multiple outliers problem since according to Engelen et al. (2003), PLS and PCR are very sensitive to the presence of outliers in data set. The robust version of PLS and PCR are known to resist several types of contamination.

## REFERENCES

- Abdi, H., (2003a).** Least Squares. In M. Lewis-Beck, A. Bryman, T. Futing (Eds): *Encyclopedia for research methods for the social sciences*. Thousand Oaks (CA): Sage.
- Abdi, H., (2003b).** Partial Least Squares Regression (PLS regression). In M. Lewis-Beck, A. Bryman, T. Futing (Eds): *Encyclopedia for research methods for the social sciences*. Thousand Oaks (CA): Sage.
- Affi, A.A. and Clark, Y., (1984).** *Computer-aided multivariate analysis*. Lifetime Learning Publications, Belmont, California.
- Barker, M., (1997).** *A Comparisons of Principal Component Regression and Partial Least Squares Regression*. Multivariate Project.
- Bjökström, A., (2001).** *Ridge Regression and Inverse Problem*. Research Report. Stockholm University, Sweden.
- Bjökström, A. and Sundberg, R., (1999).** A Generalized View on Continuum Regression. *Scandinavian Journal Statistics*. **25** : 17 – 30.
- Bowerman, B.L. and O’Connell, R.T. (1990).** *Linear Statistical Models an Applied Approach*. 2nd Ed., Boston : PWS-KENT Publishing Co.

- Boeris, M.S., Luco, J.M. and Olsina, R.A., (2000).** Simultaneous Spectrophotometric Determination of Phenobarbital, Phenytoin and Methylphenobarbital in Pharmaceutical Preparations by Using Partial Least Squares and Principal Component Regression Multivariate Calibration. *Journal of Pharmaceutical and Biomedical Analysis*. **24** : 259 – 271.
- De Jong, S. (1993),** SIMPLS: An Alternative Approach to Partial Least Squares Regression. *Chemometrics and Intelligent Laboratory Systems*. **18** : 251 – 263.
- De Jong, S., Farebrother and R.W., (1994).** Extending the Relationship Between Ridge Regression and Continuum Regression. *Chemometrics and Intelligent Laboratory Systems*. **25** : 179 – 181.
- Dempster, A.P., Schatzoff, M. and Wermuth, N., (1977).** A Simulation Study of Alternatives to Ordinary least Squares. *Journal of the American Statistical Association*. **72** : 77 – 91.
- Dijkstra, T. (1983),** Some Comments on Maximum Likelihood and Partial Least Squares Methods. *Journal of Econometrics*. **22** : 67 – 90.
- Dijkstra, T. (1985),** *Latent Variables in Linear Stochastic Models: Reflections on Maximum Likelihood and Partial Least Squares Methods*. Second Edition, Amsterdam, The Netherlands: Sociometric Research Foundation.
- Engelen, S., Hubert, M., Vanden B.K. and Verboven, S. (2003).** Robust PCR and Robust PLSR: a comparative study. *Theory and Applications of Recent Robust Methods*. edited by M. Hubert, G. Pison, A. Struyf and S. Van Aelst, Series : Statistics for Industry and Technology, Birkhauser, Basel.
- Farebrother, R. W., (1999).** A Class of Statistical Estimators Related to Principal Components. *Linear Algebra and its Applications*. **289** : 121 – 126.

- Filzmoser, P., Croux, C., (2002).** A Projection Algorithm for Regression with Collinearity. *Classification, Clustering and Data Analysis*, 227 – 234.
- Filzmoser, P., Croux, C., (2003).** Dimension Reduction of the Explanatory Variables in Multiple Linear Regression. *Pliska Studia Mathematica Bulgaria*. **14** : 59 – 70.
- Garthwaite, P.H., (1994).** An Interpretation of Partial Least Squares. *Journal of the American Statistical Association*. **89** : 122 – 127.
- Geladi, P., Kowalski, B.R., (1986a).** Partial Least-Squares Regression: A Tutorial. *Analytica Chimica Acta*. **185** : 1 – 17.
- Geladi, P., Kowalski, B.R., (1986b).** An Example of 2-Block Predictive Partial Least Squares Regression with Simulated Data. *Analytica Chimica Acta*. **185** : 19 – 32.
- Gibbons, D.G., (1981).** A Simulation Study of Some Ridge Estimators. *Journal of the American Statistical Association*. **76** : 131 – 139.
- Gruber, M.H.J., (1990).** *Regression Estimators : A Comparative Study*. California : Academic Press
- Hair, J. F., Anderson, R.E., Tatham, R.L. and Black, W.C., (1998).** *Multivariate Data Analysis*. 5th Ed., New Jersey : Prentice Hall.
- Hansen, P.M. and Schjoerring, J.K., (2003).** Reflectance Measurement of Canopy Biomass and Nitrogen Status in Wheat Crops Using Normalized Difference Vegetation Indices and Partial Least Squares. *Remote Sensing of Environment*. **86** : 542 – 553.
- Hawkins, D.M. and Yin, X., (2002).** A Faster Algorithm for Ridge Regression of Reduced Rank Data. *Journal of Computational Statistics & Data Analysis*. **40** : 253 – 262.

- Hocking, R.R., (1996).** *Methods and Applications of Linear Models : Regression and the Analysis of Variance.* USA : John Wiley & Sons.
- Hoerl, R.W., Schuenemeyer, J.H. and Hoerl, A.E., (1986).** A Simulation of Biased Estimation and Subset Selection Regression Techniques. *Technometrics.* **28(4)** : 369 – 380.
- Hoerl, A.E., Kennard, R.W., (1970).** Ridge Regression : Biased Estimation to Nonorthogonal Problems. *Technometrics.* **12** : 56 – 67.
- Hubert, M., Branden .K.V., (2003a).** Robust methods for Partial Least Squares Regression, *Journal of Chemometrics.* **17** : 537-549.
- Hubert, M., Branden .K.V., (2003b).** Robustness properties of a robust Partial Least Squares Regression, *Journal of Chemometrics.* **17** : 537-549.
- Hwang, J.T., Nettleton, D., (2000).** Principal Components Regression With Data-Chosen Components and Related Methods. *Technometrics.* **45** : 70 – 79.
- Jennrich, R.I., (1995).** *An Introduction to Computational Statistics : Regression Analysis.* New Jersey : Prentice Hall.
- Kvalheim, O. M. (1987).** Latent-structure decompositions (projections) of multivariate data. *Chemometrics and Intelligent Laboratory Systems.* **2** : 283 – 290
- Kleinbaum, D.G., Kupper, L.L., Muller, K.E. and Nizam, A., (1998).** *Applied Regression Analysis and Other Multivariable Methods.* USA : DUXBURY Press.
- Li, B., Martin, E.B. and Moris, A.J., (2001).** Box-Tidwell Transformation Based Partial Least Squares Regression. *Computers and Chemical Engineering.* **25** : 1219 – 1233.
- Marquardt, (1970).** Generalized Inverses, Ridge Regression, Biased Linear Estimation and Nonlinear Estimation. *Technometrics.* **12** : 591 – 612.

- Miller, A.J., (1990).** *Subset Selection in Regression*. Chapman & Hall, New York.
- Myers, R.H., (1986).** *Classical and Modern Regression With Applications*. 2nd Ed.  
USA : PWS-KENT Publishing Company.
- Neter, J., Wasserman, W. and Kutner, M.H., (1985).** *Applied Linear Statistical Models*. 2nd Ed. USA : IRWIN.
- Neter, J., Wasserman, W. and Kutner, M.H., (1990).** *Applied Linear Regression Models*. 3rd Ed. USA : IRWIN Book Team.
- Neter, J., Kutner, M.H., Nachtseim, C.J. and Wasserman, W., (1996).** *Applied Linear Statistical Models*. 4th Ed. USA : Irwin Bokk Team.
- Ozcelik, Y., Kulaksiz S. and Cetin, M.C., (2002).** Assessment of the Wear of Diamond Beads in the Cutting of Different Rock Types by the Ridge Regression. *Journal of Materials Processing Technology*. **127** : 392 – 400.
- Rougoor, C.W., Sundaram, R. and Van Arendonk, J.A.M., (2000).** The Relation Between Breeding Management and 305-day Milk Production, Determined via Principal Components Regression and Partial Least Squares. *Livestock Product Science*.  
**66** : 71 – 83.
- Serneels, S. and Van Espen, P.J., (2003).** Sample specific prediction intervals in SIMPLS. In *PLS and related methods*, Ed. Vilares, M., Tenenhaus, M., Coelho, P., Esposito, V., Vinzi and Morineau, A., DECISIA, Levallois Perret (France), 219 – 233.
- Stone, M. and Brooks, R.J., (1990).** Continuum Regression : Cross-Validated Sequentially Constructed Prediction Embracing OLS, PLS and PCR. *Journal of Royal Statistic Society*. **52** : 237 – 269.



- Sundberg, R., (1993).** Continuum Regression and Ridge Regression. *Journal of Royal Statistics Society.* **55** : 653 – 659.
- Sundberg, R., (2002).** Continuum Regression. Article for 2nd Ed. of *Encyclopedia of Statistical Sciences.*
- Tobias, R.D., (1997).** *An Introduction to Partial Least Squares Regression.* Cary, NC : SAS Institute.
- Wannacott, T.H. and Wannacott, R.J., (1981).** *Regression :A Second Course in Statistics.* USA : John Wiley & Sons.
- Wesolowsky, G.O., (1976).** *Multiple Regression and Analysis of Variance.* USA : John Wiley & Sons.
- Wold, H. (1966),** Estimation of Principal Components and Related Models by Iterative Least Squares. In *Multivariate Analysis*, Ed. Krishnaiah, P.R., New York : Academic Press, 391 – 420.
- Wold, S. (1994),** PLS for Multivariate Linear Modeling. *QSAR: Chemometric Methods in Molecular Design. Methods and Principles in Medicinal Chemistry*, ed. H. van de Waterbeemd, Weinheim, Germany: Verlag-Chemie.
- Xie, Y.L. and Kalivas, J.H., (1997a).** Evaluation of Principal Component Selection Methods to form a Global Prediction Model by Principal Component Regression. *Analytica Chimica Acta.* **348** : 19 – 27.
- Xie, Y.L. and Kalivas, J.H., (1997b).** Local Prediction Models by Principal Component Regression. *Analytica Chimica Acta.* **348** : 29 – 38.
- Younger, M.S., (1979).** *A Handbook for Linear Regression.* USA : DUXBURY Press.