

DOCTORAL THESIS

Multivariate statistical diagnostics with application to the growth curve model

Pan, Jianxin

Date of Award:
1996

[Link to publication](#)

General rights

Copyright and intellectual property rights for the publications made accessible in HKBU Scholars are retained by the authors and/or other copyright owners. In addition to the restrictions prescribed by the Copyright Ordinance of Hong Kong, all users and readers must also observe the following terms of use:

- Users may download and print one copy of any publication from HKBU Scholars for the purpose of private study or research
- Users cannot further distribute the material or use it for any profit-making activity or commercial gain
- To share publications in HKBU Scholars with others, users are welcome to freely distribute the permanent URL assigned to the publication

**Multivariate Statistical Diagnostics with Application
to the Growth Curve Model**

PAN Jian-Xin

A thesis submitted in partial fulfillment of the requirements

for the degree of

Doctor of Philosophy

April 1996

Hong Kong Baptist University



gift

8-9-9)

14246053

TH

1996 PM

PREFACE

The growth curve model (GCM) is a generalized multivariate analysis-of-variance model (GMANOVA), which was first proposed by Potthoff and Roy (1964) and has been studied subsequently by many authors such as Rao (1965), Khatri (1966) and von Rosen (1989). This model is useful especially for investigating growth problems on short times in economics, biology, medical research and epidemiology. Also, it is one of the fundamental tools for dealing with longitudinal data especially with serial correlation (Jones, 1993) as well as repeated measures (Crowder and Hand, 1990). However, it is not uncommon to find outliers and influential observations in growth data that significantly affect the estimates of the GCM. Under a specific covariance structure, namely, spherical covariance structure, Liski (1991) presented several measures for detecting outliers and influential observations in the GCM. For this model with other commonly used covariance structures, such as Rao's simple covariance structure (SCS) as well as unstructured covariance scheme (UCS), how multiple outliers and/or influential observations should be detected effectively is still open.

It is the purpose of this thesis to systematically investigate multivariate statistical diagnostics for the GCM. The research will be mainly concentrated on the GCM with the two covariance structures mentioned above. The multivariate diagnostic techniques used in this thesis can be classified into two categories: global influence (also known as case-deletion approach) and local influence, and each of them is employed to diagnose the adequacy of the GCM based on likelihood framework as well as Bayesian framework.

Chapter 1 of this thesis gives a background of statistical diagnostics, a brief introduction to multiple outlier identification in high-dimensional data sets, and a brief review of the GCM as well as its model selection criteria with respect to covariance structure, which serves as a basis of diagnostics development for the GCM. Three real-life biological data sets are analyzed in order to emphasize the use of the selection criteria. Also, the main contributions of this thesis on statistical diagnos-

tics to the GCM are presented in a summarized form in this chapter. In addition, some preparations related to matrix derivatives and matrix-variate distributions are provided for later use.

In Chapter 2, based on likelihood framework, we use case-deletion approach to explore the relation between the multiple individual deletion model and the mean shift regression model, to build up multiple outlier detection criteria, and to construct influential measures based on the generalized Cook's distance and the confidence ellipsoid volume. Also, the influential measures are used to assess a linear combination of the regression coefficients. These diagnostic techniques are applied to the GCM with SCS and UCS, respectively. For illustration, the three biological data sets used in Chapter 1 are analyzed by the use of such techniques for outlier detection and influential observation identification.

Chapter 3 is devoted to discussion of Cook's (1986) likelihood local influence applied to diagnose the adequacy of the GCM with SCS and UCS, respectively. Under these two covariance structures, the observed information matrix as well as the Hessian matrix for the GCM are investigated in detail, which serves as a basis of the local influence assessment for the model. As an ancillary result, the Hessian matrix is shown to be invariant under a one-to-one measurable transformation of parameters. Also, the real-life data sets are analyzed by using the local influence approach.

Based on Bayesian framework, in Chapter 4 we discuss the influence of a subset of observations on the growth fittings by the use of case-deletion approach. Under a non-informative prior distribution, the posterior distributions of the parameters for the GCM with SCS and UCS are derived, respectively. The Kullback-Leibler divergence or Bayesian entropy in information theory is then used for measuring the change of the posterior distributions when the subset of observations is deleted from the data. Also, the numerical examples presented in the pervious chapters are analyzed again by using the methods of this chapter.

From the Bayesian point of view, Chapter 5 is devoted to discussion of local influence approach for the GCM. The fundamental idea is to replace the likelihood displacement with the Kullback-Leibler divergence. For the two commonly used covariance structures, SCS and UCS, the Bayesian Hessian matrices of the GCM are investigated respectively under an abstract perturbation scheme, which play pivotal roles in the Bayesian local influence for the model. Also, some new properties of the Bayesian Hessian matrix are obtained as ancillary results. For illustration, the covariance-weighted perturbation is considered especially and employed to analyze the real-life biological data sets addressed previously.

In addition, more comments and recommendations for further research on this subject are proposed, which could be viewed as a guideline of further investigation for the GCM. Finally, the subroutines for calculating the model selection criteria with respect to covariance structure and all the diagnostic measures obtained in this thesis, written in MATLAB, are presented in Appendix of this dissertation.

Jian-Xin Pan
Department of Mathematics
Hong Kong Baptist University
224 Waterloo Road, Kowloon Tong
Hong Kong

April 15, 1996

CONTENTS

| | |
|---|----|
| Preface | i |
| Declaration | iv |
| Acknowledgments | v |
| Notations and Abbreviations | ix |
| 1. Introduction and Summary | 1 |
| 1.1 General Remarks | 1 |
| 1.1.a Statistical diagnostics | 1 |
| 1.1.b Outlier and influential observation | 3 |
| 1.2 Diagnostics in Multivariate Analysis | 8 |
| 1.2.a Multivariate outliers in high-dimensional data sets | 8 |
| 1.2.b Diagnostics in multivariate models | 11 |
| 1.3 Growth Curve Model | 14 |
| 1.3.a Growth curve model (GCM) | 14 |
| 1.3.b Covariance structure selection for the GCM | 17 |
| 1.4 Summary of This Thesis | 23 |
| 1.4.a Diagnostics for the GCM based on likelihood framework | 25 |
| 1.4.b Diagnostics for the GCM based on Bayesian framework | 26 |
| 1.5 Some Preparations | 28 |
| 1.5.a Matrix operation and matrix derivative | 28 |
| 1.5.b Matrix-variate normal and t distributions | 33 |
| 2. Outlier and Influential Observation for the Growth Curve Model | 38 |
| 2.1 General Remark | 38 |
| 2.1.a Outlier-generating model | 38 |
| 2.1.b Influential observation identification | 39 |
| 2.2 Outlier Detection for the GCM with SCS | 41 |
| 2.2.a Multiple individual deletion model | 42 |
| 2.2.b Mean shift regression model | 44 |
| 2.2.c Multiple outlier detection criteria | 46 |
| 2.2.d Illustrative examples | 49 |

| | | |
|-------|--|-----|
| 2.3 | Influential Observation for the GCM with SCS | 51 |
| 2.3.a | Generalized Cook-type distance | 52 |
| 2.3.b | Measurements based on confidence ellipsoid volume | 54 |
| 2.3.c | Influence assessment on linear combination of \mathbf{B} | 57 |
| 2.3.d | Illustrative example | 60 |
| 2.4 | Outlier Detection for the GCM with UCS | 64 |
| 2.4.a | Multiple individual deletion model | 64 |
| 2.4.b | Mean shift regression model | 67 |
| 2.4.c | Multiple outlier detection criteria | 70 |
| 2.4.d | Illustrative example | 76 |
| 2.5 | Influential Observation for the GCM with UCS | 77 |
| 2.5.a | Generalized Cook-type distance | 78 |
| 2.5.b | Measurements based on confidence ellipsoid volume | 79 |
| 2.5.c | Influence assessment on linear combination of \mathbf{B} | 82 |
| 2.5.d | Illustrative example | 86 |
| 3. | Local Influence Assessment for the Growth Curve Model | 88 |
| 3.1 | General Remark | 88 |
| 3.1.a | Background | 88 |
| 3.1.b | Local influence analysis | 89 |
| 3.2 | Local Influence Assessment for the GCM with SCS | 92 |
| 3.2.a | Observed information matrix | 92 |
| 3.2.b | Hessian matrices | 94 |
| 3.2.c | Covariance-weighted perturbation | 98 |
| 3.2.d | Illustrative examples | 101 |
| 3.3 | Local Influence Assessment for the GCM with UCS | 105 |
| 3.3.a | Observed information matrix | 105 |
| 3.3.b | Hessian matrices | 107 |
| 3.3.c | Covariance-weighted perturbation | 113 |
| 3.3.d | Illustrative example | 114 |
| 4. | Bayesian Influence Analysis for the Growth Curve Model | 117 |
| 4.1 | General Remark | 117 |
| 4.1.a | Bayesian influence analysis | 117 |
| 4.1.b | Kullback-Leibler divergence of matrix-variate distribution | 120 |
| 4.2 | Bayesian Influence Analysis for the GCM with SCS | 122 |
| 4.2.a | Posterior distributions | 122 |

| | | |
|-------|--|-----|
| 4.2.b | Bayesian influence measure | 124 |
| 4.2.c | Illustrative example | 128 |
| 4.3 | Bayesian Influence Analysis for the GCM with UCS | 131 |
| 4.3.a | Posterior distributions | 131 |
| 4.3.b | Bayesian influence measure | 136 |
| 4.3.c | Illustrative example | 144 |
| 5. | Bayesian Local Influence Assessment for the Growth Curve Model | 147 |
| 5.1 | General Remark | 147 |
| 5.1.a | Bayesian local influence | 147 |
| 5.1.b | Bayesian Hessian matrix of matrix-variate distribution | 151 |
| 5.2 | Bayesian Local Influence for the GCM with SCS | 157 |
| 5.2.a | Bayesian Hessian matrix under SCS | 157 |
| 5.2.b | Covariance-weighted perturbation | 160 |
| 5.2.c | Illustrative examples | 162 |
| 5.3 | Bayesian Local Influence for the GCM with UCS | 166 |
| 5.3.a | Bayesian Hessian matrix under UCS | 167 |
| 5.3.b | Covariance-weighted perturbation | 171 |
| 5.3.c | Illustrative example | 175 |
| | Concluding Remarks and Recommendations | 178 |
| | Appendix | 181 |
| | MATLAB subroutines | 181 |
| | A.1 Model selection with respect to covariance structure | 181 |
| | A.2 Global influence measures under SCS | 183 |
| | A.3 Global influence measures under UCS | 185 |
| | A.4 Local influence measures under SCS | 188 |
| | A.5 Local influence measures under UCS | 191 |
| | A.6 Bayesian global influence measures under SCS | 194 |
| | A.7 Bayesian global influence measures under UCS | 197 |
| | A.8 Bayesian Local influence measures under SCS | 200 |
| | A.9 Bayesian Local influence measures under UCS | 203 |
| | References | 205 |
| | Vita | 216 |