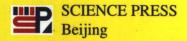
Mathematics Monograph Series 8

Growth Curve Models and Statistical Diagnostics

Jianxin Pan and Kaitai Fang

(生长曲线模型及其统计诊断)



Mathematics Bloograph limits 🗗

Growth Curve Models and Statistical Diagnostics

Jacobs Dec and Kalcal Resp.

1月 中国国际企工及设计设计。



Mathematics Monograph Series 8

Jianxin Pan and Kaitai Fang

1.12

Growth Curve Models and Statistical Diagnostics

(生长曲线模型及其统计诊断)

Responsible Editor: Jiashan Liu, Deping Yan, Yang Zhang

Copyright© 2007 by Science Press Published by Science Press 16 Donghuangchenggen North Street Beijing 100717, China

Printed in Beijing

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior written permission of the copy right owner.

ISBN 978-7-03-019532-6 (Beijing)

To the memory of my mother, to my father in his 70th year, and to my wife Haiyan and my son Kainan for their patience during the writing of this book.

J.X. Pan

Jianxin Pan Centre for Medical Statistics Department of Mathematics Keele University Staffordshire, ST5 5BG United Kingdom

To my wife Tingmei and my two daughters, Ying and Yan, for their constant support.

K.T. Fang

Kaitai Fang Department of Mathematics Hong Kong Baptist University 224 Waterloo Road, Kowloon Hong Kong

Preface

A growth curve model (GCM) is a generalized multivariate analysis-of-variance model (GMANOVA), first summarized by Potthoff and Roy (1964) and studied subsequently by many authors including Rao (1965), Khatri (1966) and von Rosen (1989). The GCM is especially useful for investigating grwoth problems on short time series in areas such as economics, biology, medical research, and epidemiology. It is also a fundamental tool for analyzing longitudinal data especially with serial correlation (Jones, 1993) and repeated measures (Crowder and Hand, 1990). It is not uncommon, however, to find outliers and influential observations in growth data that significantly affect statistical inference in the GCM.

The purpose of this book to introduce the theory of the GCM with particular emphasis on statistical diagnostics, which is mainly based on recent work on diagnostics made by the authors and their collaborators. This book is intended for researchers who are working in the area of theoretical studies related to the GCM as well as multivariate statistical diagnostics, and for applied statisticians working in application of the GCM to practical areas. Hence, on the one hand, we provide theoretical proofs for the most theorems in this book. On the other hand, applications of these techniques to practical data analysis are emphasized; for example, almost every approach discussed in this book is illustrated with practical examples. In addition, the computer programmes for calculating various measurements involved in this book have been written in S-PLUS and GENSTAT. We will put the computer programmes on our web site in due course. A link to the web site can be found in the list of author web pages at the Springer web page, www.springer-ny.com.

The statistical diagnostics considered in this book focuses mainly on GCMs with two specific covariance structures, namely, Rao's simple covariance structure (SCS) and unstructured covariance (UC), since these two covariance structures are very common in practice and some other covariance structures are their special cases. For example, the uniform covariance

ii Preface

structure and random-effects covariance structure are two special cases of SCS. GCMs with other covariance structures can also be analyzed in a similar manner. The multivariate diagnostic techniques addressed in this book are classified into two categories: global influence (also known as case-deletion approach) and local influence; and each of these is used to diagnose the adequacy of GCMs within the likelihood and Bayesian frameworks as well.

Chapter 1 of this book gives a background of statistical diagnostics, a brief introduction to multiple outlier identification in multivariate data sets, a brief review of the GCM, and the related model selection criteria with respect to covariance structure. Also, the main approaches and results on statistical inferences and diagnostics in GCMs are presented in a summarized form in this chapter. In addition, preparatory materials related to matrix derivatives and matrix-variate distributions are provided for later use.

In Chapter 2, the fundamental concepts of GCMs are introduced and several most commonly encountered forms of the models are explained in terms of practical examples in biology, agriculture and medical sciences. The generalized least square estimate (GLSE) and admissibility of estimates on linear combinations of regression coefficients are discussed. We show that the GLSE of the regression coefficient is also the best linear unbiased estimate (BLUE) in the sense of matrix loss functions. We also study here the necessary and sufficient conditions of admissible estimates of linear combinations of the regression coefficients.

Maximum likelihood estimate (MLEs) of the regression coefficient and dispersion component in growth curve models are discussed in Chapter 3. We also study the expectation and variance-covariance matrix of the estimates. In general, the MLE of the regression coefficient is different from the GLSE given in Chapter 2. In fact, the latter is a linear function of the response variable while the former is not. There is indeed a special case, however, in which the MLE is completely identical to the GLSE. In this case, statistical inferences based on the MLE in growth curve models becomes simpler. This special case is nothing but SCS, in which the dispersion component matrix Σ is decomposed as two orthogonal components. This point will be shown with illustrative examples. As an alternative to the MLE, restricted maximum likelihood (REML) estimates are studied in the context of growth curve models with SCS and random effects covariance structure in this chapter. Estimates of the dispersion components are unbiased in this case. Numerical studies are conducted to compare the GLSE, MLE and REML in growth curve models.

Within the likelihood framework in Chapter 4 we use the case deletion technique to explore the relationship between the multiple individual deletion model (MIDM) and the mean shift regression model (MSRM), to build up multiple outlier detection criteria, and to construct influence measurements based on the generalized Cook's distance and the confidence ellipsoid

Preface

volume. Also, influence measurements are used to assess a linear combination of regression coefficients. These diagnostic techniques are applied to GCMs with SCS and UC, respectively. For illustration, some biological, medical, and agricultural data sets are analyzed using these diagnostic techniques for outlier detection and influential observation identification.

Chapter 5 is devoted to discussing how Cook's (1986) likelihood-based local influence technique could be used to diagnose the applied of GCMs with SCS and UC, respectively. With these two specific covariance structures, the observed information matrix and the Hessian matrix are studied; the Hessian matrix serves as a basis of the local influence assessment in these models. As an ancillary result, the Hessian matrix is shown to be invariant under a one-to-one measurable transformation of parameters. Also, the practical data sets analyzed in the previous chapters are reanalyzed using local influence approach discussed in this chapter.

In Chapter 6, within the Bayesian framework, we discuss the influence of a subset of observations on growth fittings in terms of case deletion technique. Under a noninformative prior distribution, the posterior distributions of the parameters in GCMs with SCS and UC are considered, respectively. The Kullback-Leibler divergence is used to measure the change of posterior distributions when the subset of observations is removed from the model. The numerical examples addressed in the previous chapters are analyzed once again using the methods developed in this chapter.

Chapter 7 is devoted to discussion of the local influence approach in the GCM from a Bayesian point of view. The fundamental idea of Bayesian local influence is to replace the likelihood displacement of Cook's local influence with the Kullback-Leibler divergence. For the two commonly used covariance structures, SCS and UC, Bayesian Hessian matrices in the GCM are studied under an abstract perturbation. Those matrices play a pivotal role in the Bayesian local influence. Also, some properties of Bayesian Hessian matrix are considered as ancillary results, and the relationships between likelihood-based local influence and Bayesian local influence are studied. For illustration, the covariance-weighted perturbation is considered especially and employed to analyze several practical data sets.

We would like to thank the following people from Hong Kong Baptist University: Prof. C. F. Ng, Dr. C. W. Tse, and Prof. F. J. Hickernell. We also thank Prof. G. MacKenzie and Prof. P. W. Jones of Keele University, Prof. R. Thompson of IACR-Rothamsted, and Prof. X. R. Wang of Yunnan University. All gave valuable inspiration on our research and constant encouragement. We are grateful to Prof. P. M. Bentler and Prof. D. von Rosen for their encouragement in the writing of this book. We also acknowledge with gratitude P.M. Bentler, D. von Rosen, R. Thompson, G. MacKenzie, and anonymous referees for their reading of the manuscript and their invaluable comments and suggestions. We thank our various coauthors in a series of papers, including D. von Rosen, E.P. Liski, W.K. Fung, and P. Bai for the nice and very stimulating collaborations, from which some

iv Preface

results are reflected in this book. Also, we are indebted to J. S. Liu and X. Y. Ge of Science Press and J. Kimmel and J. Wolkowicki of Springer-Verlag for their hard work on the publication of this book. Mrs J.Drewery of Keele University kindly helped us to prepare the Index pages and her great help is very appreciated. Last but not least, we gratefully acknowledge support from our families.

Jian-Xin Pan's research was partially supported by a grant from Yunnan Science Foundation, a fellowship from Hong Kong Baptist University during his time there as a Ph. D. student, a grant from Agriculture, Environment and Fisheries Department of the Scottish Office during his time as a post-doctoral research fellow in IACR-Rothamsted, and a grant from the Acute (NHC) Trust of North Staffordshire. Kai-Tai Fang's research was partially supported by the Hong Kong University Grant Council. We would like to acknowledge gratefully the generous support of all these institutions.

Keele University Hong Kong Baptist University Jian-Xin Pan Kai-Tai Fang

21st February, 2007

Acronyms

AIC Akaike Information Criterion

ANOVA Analysis of Variance

BIC Bayesian Information Criterion
BLUE Best Linear Unbiased Estimate
BLUP Best Linear Unbiased Prediction

CLS Cyclic Lattice Squares

DFFITS Difference between fitted values EM Expectation–Maximization algorithm

GCM Growth Curve Model GLM Generalized Linear Model

GLMM Generalized Linear Mixed Model

GLP Good Lattice Points

GLSE Generalized Least Square Estimate

GMANOVA Generalized Multivariate Analysis-of-Variance

i.i.d. Independent Identically Distributed random variable

KLD Kullback-Leibler Divergence

LISREL Linear structure regression equation analysis

LSE Least Square Estimate

MAD Median of Absolute Deviation from the median

MCMC Markov Chain Monte Carlo

MD Mahalanobis Distance

MIDM Multiple Individual Deletion Model MLE Maximum Likelihood Estimate MSRM Mean Shift Regression Model

MVE Minimum Volume Ellipsoid estimate OMRM Ordinary Multivariate Regression Model

p.d. Positive Definite matrixp.d.f. Probability Density Function

PP Projection Pursuit

PQL Penalized Quasi-Likelihood estimator REML Restricted Maximum Likelihood

RSS Residual Sum of Squares

SCS Rao's Simple Covariance Structure SVD Singular Value Decomposition

UC Unstructured Covariance

Notations

Space

 I_p

Matrices are denoted by upper case bold letters, column vectors by lower case bold letters, and scalars by lower case letters. For example, \boldsymbol{X} , \boldsymbol{x} and \boldsymbol{x} represent a matrix, a column vector, and a scalar, respectively.

Space	
R^p	p-dimensional Euclidian space
S^q	unit sphere in the q -dimensional Euclidian space \mathbb{R}^q
Scalar	
$J(oldsymbol{X} ightarrow oldsymbol{Y})$	Jacobian of the transformation from X to Y
$\operatorname{tr}(\boldsymbol{A})$	trace of the square matrix \boldsymbol{A}
$\det(\hat{\boldsymbol{A}})$	determinant of the square matrix \boldsymbol{A}
$\Gamma_m(a)$	multivariate gamma function, i.e.,
	$\Gamma_m(a) = \pi^{m(m-1)/4} \prod_{i=1}^m \Gamma(a - (i-1)/2)$
Vector	
1_{p}	$1_{p}=(1,1,\cdots,1)\in R^{p}$, i.e., the <i>p</i> -variate vector of ones
$oldsymbol{d_{ ext{max}}}$	unit eigenvector associated with the largest
	absolute eigenvalue of the Hessian matrix
$ ext{vec}(oldsymbol{A})$	direct operator of the matrix \boldsymbol{A}
$ ext{vec}'(oldsymbol{A})$	transpose of the vector $\text{vec}(\boldsymbol{A})$
$\operatorname{svec}(\boldsymbol{A})$	symmetric direct operator of the symmetric matrix \boldsymbol{A}
$\operatorname{svec}'(\boldsymbol{A})$	transpose of the vector $svec(\mathbf{A})$
$\mathrm{diag}(\!\boldsymbol{A})$	vector formed by the diagonal elements of \boldsymbol{A}
Matrix	
$oldsymbol{A}'$	transpose of the matrix \boldsymbol{A}
A^{-1}	inverse of the nonsingular matrix \boldsymbol{A}
A^+	Moore–Penrose generalized inverse of the matrix \boldsymbol{A}
	5

identity matrix with order $p \times p$

xii Notations

$oldsymbol{E_{ij}}(p,q)$	$p \times q$ matrix with (i, j) th element being one
•,	and others being zero. in short, denoted E_{ij}
$oldsymbol{K}_{pq}$	permutation matrix with order $pq \times pq$;
• •	for $p = q$, denoted \mathbf{K}_{p^2}
$oldsymbol{S_p}$	duplication matrix with order $p^2 \times p(p+1)/2$
P_{A}	$P_{\mathbf{A}} = \mathbf{A}(\mathbf{A}^{\tau}\mathbf{A})^{-1}\mathbf{A}^{\tau}$, the projection matrix of \mathbf{A}
Q_S^-	$\mathbf{Q}_{\mathbf{S}} = \mathbf{S}\mathbf{Q}(\mathbf{Q}^{\tau}\mathbf{S}\mathbf{Q})^{-1}\mathbf{Q}^{\tau}$, a semiprojection matrix.
$ ilde{oldsymbol{A}\otimes oldsymbol{B}}$	Kronecker product of the matrices \boldsymbol{A} and \boldsymbol{B}
A*B	Hadamard product of the matrices \boldsymbol{A} and \boldsymbol{B}
$oldsymbol{\Sigma} > oldsymbol{0}$	Σ is a positive definite matrix
$oldsymbol{\Sigma}^{1/2}$	square root of the matrix Σ
$Cov(\boldsymbol{X})$	variance-covariance matrix of the random matrix X ,
` ,	i.e., $Cov(\boldsymbol{X}) = Cov(vec(\boldsymbol{X}))$

Univariate distribution

$N(\mu,\sigma^2)$	univariate normal distribution with
	expectation μ and variance σ^2
t_{p}	Student's t distribution with p degrees of freedom
$egin{array}{c} t_p \ \chi_p^2 \ F_{n,m} \end{array}$	chi-square distribution with p degrees of freedom
$\dot{F_{n,m}}$	F distribution with n and m as first and
	second degrees of freedom, respectively
$\operatorname{Gamma}(a,b)$	gamma distribution with parameters a and b
$\mathrm{Beta}(a,b)$	beta distribution with parameters a and b
$\Lambda(p,m,n)$	Wilk's distribution with three parameters p, m , and n
$GT^2(m,r,n)$	Hotelling's generalized T^2 distribution
	with three parameters m, r , and n

Vector-variate distribution

$N_{m p}(\!m \mu, \!m \Sigma)$	p-dimensional normal distribution
	with expectation vector $\boldsymbol{\mu}$ and dispersion matrix $\boldsymbol{\Sigma} > 0$
$t_{m p}(m \mu, m \Sigma, u)$	p -dimensional t distribution with location μ
	and dispersion $\Sigma > 0$, and ν degrees of freedom

Matrix-variate distribution

$N_{p,n}(oldsymbol{M}, oldsymbol{\Sigma}, oldsymbol{\Omega})$	matrix-variate normal distribution with location
	matrix M and dispersion matrices $\Sigma>0$ and $\Omega>0$
$W_m(n, \mathbf{\Sigma})$	Wishart distribution with parameters n and $\Sigma > 0$
$t_{oldsymbol{p},oldsymbol{n}}(oldsymbol{M},oldsymbol{\Sigma},oldsymbol{\Omega},oldsymbol{ u})$	matrix-variate t distribution with location
	matrix M and dispersion matrices $\Sigma>0$ and $\Omega>0$
	and ν degrees of freedom

Contents

	nyms ix	
Not	tion x	i
	oter 1	
Int	duction 1	L
1.1	General Remarks 1	Ĺ
	1.1.1 Statistical Diagnostics	
	1.1.2 Outliers and Influential Observation 3	3
1.2	tatistical Diagnostics in Multivariate Analysis 10)
	1.2.1 Multiple Outliers in Multivariate Data)
	1.2.2 Statistical diagnostics in multivariate models	1
1.3	Growth Curve Model (GCM)	3
	1.3.1 A Brief Review	3
	1.3.2 Covariance Structure Selection	9
1.4	Summary 23	3
	1.4.1 Statistical Inference 2 ²	4
	1.4.2 Diagnostics Within a likelihood Framework	5
	1.4.3 Diagnostics Within a Bayesian Framework	
1.5	Preliminary Results	3
	1.5.1 Matrix Operation and Matrix Derivative	3
	1.5.2 Matrix-variate Normal and t Distributions	2
1.6	Further Readings	7
	oter 2	
Ge	eralized Least Square Estimation 38	8
2.1	General Remarks 38	
	2.1.1 Model Definition 38	8
	2.1.2 Practical Examples 4	5
2.2	Generalized Least Square Estimation	2

vi Contents

	2.2.1 Generalized Least Square Estimate (GLSE)2.2.2 Best Linear Unbiased Estimate (BLUE)	58
	2.2.3 Illustrative Examples	
2.3	Admissible Estimate of Regression Coefficient	
	2.3.1 Admissibility	
	2.3.2 Necessary and Sufficient Condition	
2.4	Bibliographical Notes	74
	apter 3	
	ximum Likelihood Estimation	
3.1	Maximum Likelihood Estimation	
	3.1.1 Maximum Likelihood Estimate (MLE)	
	3.1.2 Expectation and Variance-covariance	
	3.1.3 Illustrative Examples	
3.2	(- +)	
	3.2.1 Condition That the MLE Is Identical to the GLSE	
	3.2.2 Estimates of Dispersion Components	
	3.2.3 Illustrative Examples	
3.3		
	3.3.1 Restricted Maximum Likelihood (REMLs) estimate	
	3.3.2 REMLs Estimates in the GCM	
	3.3.3 Illustrative Examples	
3.4	Bibliographical Notes	156
	apter 4	
	scordant Outlier and Influential Observation	
4.1	General Remarks	
	4.1.1 Discordant Outlier-Generating Model	
	4.1.2 Influential Observation	
4.2	Discordant Outlier Detection in the GCM with SCS	
	4.2.1 Multiple Individual Deletion Model (MIDM)	
	4.2.2 Mean Shift Regression Model (MSRM)	
	4.2.3 Multiple Discordant Outlier Detection	
	4.2.4 Illustrative Examples	
4.3	Influential Observation in the GCM with SCS	
	4.3.1 Generalized Cook-type Distance	
	4.3.2 Confidence Ellipsoid's Volume	
	4.3.3 Influence Assessment on Linear Combination	
	4.3.4 Illustrative Examples	185
4.4	Discordant Outlier Detection in the GCM with UC	
	4.4.1 Multiple Individual Deletion Model (MIDM)	
	4.4.2 Mean Shift Regression Model (MSRM)	
	4.4.3 Multiple Discordant Outlier Detection	
	4.4.4 Illustrative Examples	204

Con	tents	vii
4.5	Influential Observation in the GCM with UC 4.5.1 Generalized Cook-type Distance 4.5.2 Confidence Ellipsoid's Volume 4.5.3 Influence Assessment on Linear Combination	207 208 212
4.6	4.5.4 Illustrative Examples	
Ch	apter 5	
	elihood-Based Local Influence	224
	General Remarks	
	5.1.1 Background	224
	5.1.2 Local Influence Analysis	226
5.2	Local Influence Assessment in the GCM with SCS	
	5.2.1 Observed Information Matrix	229
	5.2.2 Hessian Matrix	231
	5.2.3 Covariance-Weighted Perturbation	236
	5.2.4 Illustrative Examples	
5.3	Local Influence Assessment in the GCM with UC	247
	5.3.1 Observed Information Matrix	247
	5.3.2 Hessian Matrix	
	5.3.3 Covariance-Weighted Perturbation	256
	5.3.4 Illustrative Examples	
5.4	Bibliographical Notes	262
Ch	apter 6	
Ba	yesian Influence Assessment	264
6.1	General Remarks	264
	6.1.1 Bayesian Influence Analysis	264
	6.1.2 Kullback–Leibler Divergence	267
6.2	Bayesian Influence Analysis in the GCM with SCS	269
	6.2.1 Posterior Distribution	
	6.2.2 Bayesian Influence Measurement	
	6.2.3 Illustrative Examples	. 277
6.3	Bayesian Influence Analysis in the GCM with UC	. 286
	6.3.1 Posterior Distribution	. 286
	6.3.2 Bayesian Influence Measurement	293
	6.3.3 Illustrative Examples	. 301
6.4	Bibliographical Notes	. 305

viii	Contents
Y ====	Contente

	resian Local Influence	
7.1	General Remarks	
	7.1.1 Bayesian Local Influence	
	7.1.2 Bayesian Hessian Matrix	
7.2	Bayesian Local Influence in the GCM with SCS	
	7.2.1 Bayesian Hessian Matrix	320
	7.2.2 Covariance-Weighted Perturbation	323
	7.2.3 Illustrative Examples	326
7.3	Bayesian Local Influence in the GCM with UC	336
	7.3.1 Bayesian Hessian Matrix	337
	7.3.2 Covariance-Weighted Perturbation	342
	7.3.3 Illustrative Examples	347
7.4	Bibliographical Notes	351
۸n	pendix	
Дþ	to and around in Albin bounds	353
-	ta sets used in this book	000
Da	ferences	
Da Re		361

Chapter 1

Introduction

Statistical diagnostics is one of the most useful techniques in statistical science. The aim of diagnostics is to detect outliers that deviate from the postulated model, to identify influential observations that have large effects on the statistical inference drawn from the postulated model, and to validate the chosen statistical model. The main theme of this book is to comprehensively explore multivariate diagnostic techniques, which are specifically suitable for diagnosing the adequacy of multivariate models, with particular emphasis on the application to growth curve models. The approaches employed are case-deletion and local influence within the likelihood and Bayesian frameworks. We give a brief introduction to statistical diagnostics in Section 1.1 and the associated multivariate techniques in Section 1.2. Section 1.3 is devoted to a brief review of growth curve models as well as model selection criteria with respect to covariance structures. In Section 1.4, the main approaches and results in this book on statistical diagnostics for growth curve models are outlined in a summarized form. Some preparatory materials related to matrix derivatives and matrix-variate distributions are given in Section 1.5 for later use.

1.1 General Remarks

1.1.1 Statistical diagnostics

In statistical science, statistical models play an important role in analyzing data, making statistical inferences, and making future predictions. As long as a random sample or data set is drawn from a practical phenomenon, statisticians often need to build up a "good" statistical model for fitting the